

The Internet Protocol Journal

June 1998

Volume 1, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
What Is a VPN?—Part I	2
SSL: Foundation for Web Security	20
Call for Papers	30
Book Reviews	31
Fragments	35

FROM THE EDITOR

Welcome to the first edition of *The Internet Protocol Journal* (IPJ). This publication is designed to bring you in-depth technical articles on current and emerging Internet and intranet technologies. We will publish technology tutorials, as well as case studies on all aspects of internetworking.

Our first article is a detailed look at *Virtual Private Networks* (VPNs). Many organizations are turning to VPNs as a cost-effective way to implement enterprise networking, but the industry has not yet settled for a single approach, nor even a single definition of the VPN concept. The article by Paul Ferguson and Geoff Huston is in two parts. Part II will follow in our second issue, due out in September.

When the Internet Protocol suite (TCP/IP) was first designed, security was not a major consideration. Indeed, the primary goal in the early days of networking was sharing of information among academics and researchers. Today, TCP/IP is being used for mission-critical applications and for the emerging area of electronic commerce. As a result, security mechanisms are being added at all levels of the protocol stack. In this issue, we take a closer look at the *Secure Sockets Layer* (SSL), which is used for Web transactions. William Stallings explains how SSL works and how it is becoming the standard for Web security.

If you want to learn about computer networks, many options are available, including conferences, journals, standards documents, Web sites, glossaries and, of course, books. Our *Fragments* page gives you some pointers for further reading, and every issue will include at least one book review.

A detailed description of the scope of this journal can be found on page 30 in our *Call for Papers*. We want your input in this new publication. Please send comments, suggestions or questions to ipj@cisco.com. You may also use this address to request a complimentary copy of the next issue of IPJ. If you would like to write an article, send me e-mail and I will send you author guidelines.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

To reserve your complimentary
copy of the next issue of
The Internet Protocol Journal,
please complete and return the
attached postage-paid card.

What Is a VPN? — Part I

by Paul Ferguson, Cisco Systems
and Geoff Huston, Telstra

The term “VPN,” or *Virtual Private Network*, has become almost as recklessly used in the networking industry as has “QoS” (Quality of Service) to describe a broad set of problems and “solutions,” when the objectives themselves have not been properly articulated. This confusion has resulted in a situation where the popular trade press, industry pundits, and vendors and consumers of networking technologies alike generally use the term VPN as an offhand reference for a set of different technologies. This article provides a common-sense definition of a VPN, and an overview of different approaches to building one.

“The wonderful thing about virtual private networks is that its myriad definitions give every company a fair chance to claim that its existing product is actually a VPN. But no matter what definition you choose, the networking buzz-phrase doesn’t make sense. The idea is to create a private network via tunneling and/or encryption over the public Internet. Sure, it’s a lot cheaper than using your own frame relay connections, but it works about as well as sticking cotton in your ears in Times Square and pretending nobody else is around.”^[1]

A Common-Sense Definition

As *Wired Magazine* notes in the quotation, the myriad definitions of a VPN are less than helpful in this environment. Accordingly, it makes sense to begin this examination of VPNs to see if it is possible to provide a common-sense definition of a VPN. Perhaps the simplest method of attempting to arrive at a definition for VPNs is to look at each word in the acronym individually, and then tie each of them together in a simple, common-sense, and meaningful fashion.

Let’s start by examining the word “network.” This term is perhaps the least difficult one for us to define and understand, because the commonly accepted definition is fairly uncontroversial and generally accepted throughout the industry. A network consists of any number of devices that can communicate through some arbitrary method. Devices of this nature include computers, printers, routers, and so forth, and they may reside in geographically diverse locations. They may communicate in numerous ways because the electronic signaling specifications, and data-link, transport, and application-layer protocols are countless. For the purposes of simplicity, let’s say that a “network” is a collection of devices that can communicate in some fashion, and can successfully transmit and receive data among themselves.

The term “private” is fairly straightforward, and is intricately related to the concept of “virtualization” insofar as VPNs are concerned, as we’ll discuss in a moment. In the simplest of definitions, “private” means communications between two (or more) devices is, in some

fashion, secret—that the devices that are not participating in the “private” nature of communications are not privy to the communicated content, and that they are indeed completely unaware of the private relationship altogether. Accordingly, data privacy and security (data integrity) are also important aspects of a VPN that need to be considered when implementing any particular VPN.

Another means of expressing this definition of “private” is through its antonym, “public.” A “public” facility is one that is openly accessible, and is managed within the terms and constraints of a common public resource, often via a public administrative entity. By contrast, a private facility is one where access is restricted to a defined set of entities, and third parties cannot gain access. Typically, the private resource is managed by the entities who have exclusive right of access. Examples of this type of private network can be found in any organizational network that is not connected to the Internet, or to any other external organizational network, for that matter. These networks are private because there is no external connectivity, and thus no external network communications.

Another important aspect of privacy in a VPN is through its technical definition. For example, privacy in an addressing and routing system means that the addressing used within a VPN community of interest is separate and discrete from that of the underlying shared network, and from that of other VPN communities. The same holds true for the routing system used within the VPN and that of the underlying shared network. The routing and addressing scheme within a VPN should, in general, be self-contained, but this scenario degenerates into a philosophical discussion of the context of the term “VPN.” Also, it is worthwhile to examine the differences between the “peer” and “overlay” models of constructing VPNs—both of which are discussed in more detail later under the heading “Network-Layer VPNs.”

“Virtual” is a concept that is slightly more complicated. *The New Hacker’s Dictionary* (formerly known as the Jargon File)^[2] defines virtual as:

virtual /adj./ [via the technical term “virtual memory,” prob. from the term “virtual image” in optics] 1. Common alternative to {logical}; often used to refer to the artificial objects (like addressable virtual memory larger than physical memory) simulated by a computer system as a convenient way to manage access to shared resources. 2. Simulated; performing the functions of something that isn’t really there. An imaginative child’s doll may be a virtual playmate. Oppose {real}.

Insofar as VPNs are concerned, the second definition is perhaps the most appropriate comparison for virtual networks. The “virtualization” aspect is one that is similar to what we briefly described previously as private, but the scenario is slightly modified—the private communication is now conducted across a network infrastructure that

is shared by more than a single organization. Thus, the private resource is actually constructed by using the foundation of a logical partitioning of some underlying common, shared resource rather than by using a foundation of discrete and dedicated physical circuits and communications services. Accordingly, the private network has no corresponding private physical communications system. Instead, the private network is a virtual creation that has no physical counterpart.

The virtual communications between two (or more) devices is because the devices that are not participating in the virtual communications are not privy to the content of the data, and they are also altogether unaware of the private relationships between the virtual peers. The shared network infrastructure could, for example, be the global Internet and the number of organizations or other users not participating in the virtual network may literally number into the thousands or even millions.

A VPN can also said to be a discrete network^[3]:

(discrete \dis*crete", a. [L. discretus, p.p. of discernere. See Discreet.]
1. Separate; distinct; disjunct).

The discrete nature of VPNs allows both privacy and virtualization. Although VPNs are not completely separate, intrinsically, the distinction is that they operate in a discrete fashion across a shared infrastructure, providing exclusive communications environments that do not share any points of interconnection.

The combination of these terms produces VPN—a private network, where the privacy is introduced by some method of virtualization. A VPN could be built between two end systems or between two organizations, between several end systems within a single organization or between multiple organizations across the global Internet, between individual applications, or any combination.

It should be noted that there is really no such thing as a nonvirtual network, if the underlying common public transmission systems and other similar public infrastructure components are considered to be the base level of carriage of the network. What separates a VPN from a truly private network is whether the data transits a shared versus a nonshared infrastructure. For instance, an organization could lease private line circuits from various telecommunications providers and build a private network on the base of these private circuit leases, but the circuit-switched network owned and operated by the telecommunications companies are actually circuits connected to their *Digital Access and Crossconnect Systems* (DACSS) network and subsequently their fiber-optics infrastructure. This infrastructure is shared by any number of organizations through the use of multiplexing technologies. Unless an organization is actually deploying private fiber and layered transmission systems, any network is layered with “virtualized” connectivity services in this fashion.

A VPN doesn't necessarily mean communications isolation, but rather the controlled segmentation of communications for communities of interest across a shared infrastructure.

The common and somewhat formal characterization of the VPN, and perhaps the most straightforward and strict definition, follows:

A VPN is a communications environment in which access is controlled to permit peer connections only within a defined community of interest, and is constructed through some form of partitioning of a common underlying communications medium, where this underlying communications medium provides services to the network on a nonexclusive basis.

A simpler, more approximate, and much less formal description follows:

A VPN is private network constructed within a public network infrastructure, such as the global Internet.

It should also be noted that although VPNs may be constructed to address any number of specific business needs or technical requirements, a comprehensive VPN solution provides support for dial-in access, support for multiple remote sites connected by leased lines (or other dedicated means), the ability of the VPN service provider (SP) to "host" various services for the VPN customers (for example, Web hosting), and the ability to support not just intra-, but also inter-VPN connectivity, including connectivity to the global Internet.

VPN Motivations

There are several motivations for building VPNs, but a common thread is that they all share the requirement to "virtualize" some portion of an organization's communications—in other words, make some portion (or perhaps all) the communications essentially "invisible" to external observers, while taking advantage of the efficiencies of a common communications infrastructure.

The base motivation for VPNs lies in the economics of communications. Communications systems today typically exhibit the characteristic of a high fixed-cost component, and smaller variable-cost components that vary with the transport capacity, or bandwidth, of the system. Within this economic environment, it is generally financially attractive to bundle numerous discrete communications services onto a common, high-capacity communications platform, allowing the high fixed-cost components associated with the platform to be amortized over a larger number of clients. Accordingly, a collection of virtual networks implemented on a single common physical communications plant is cheaper to operate than the equivalent collection of smaller, physically discrete communications plants, each servicing a single network client.

Therefore, if aggregation of communications requirements leads to a more cost-effective communications infrastructure, why not pool all these services into a single public communications system? Why is there still the requirement to undertake some form of partitioning within this common system that results in these “virtual private” networks?

In response to this question, the second motivation for VPNs is that of communications privacy, where the characteristics and integrity of communications services within one closed environment is isolated from all other environments that share the common underlying plant. The level of privacy depends greatly on the risk assessment performed by the subscriber organization—if the requirement for privacy is low, then the simple abstraction of discretion and network obscurity may serve the purpose. However, if the requirement for privacy is high, then there is a corresponding requirement for strong security of access and potentially strong security applied to data passed over the common network.

History

This article cannot do justice to the concept of VPNs without some historical perspective, so we need to look at why VPNs are an evolving paradigm, and why they will continue to be an issue of confusion, contention, and disagreement. This examination is important because opinions on VPN solutions are quite varied, as well as how they should be approached.

Historically, one of the precursors to the VPN was the *Public Data Network* (PDN), and the current familiar instance of the PDN is the global Internet. The Internet creates a ubiquitous connectivity paradigm, where the network permits any connected network entity to exchange data with any other connected entity. The parallels with the global *Public Switched Telephone Network* (PSTN) are, of course, all too obvious—where a similar paradigm of ubiquitous public access is the predominate characteristic of the network.

The Public Data Network has no inherent policy of traffic segregation, and any modification to this network policy of permitting ubiquitous connectivity is the responsibility of the connecting entity to define and enforce. The network environment is constructed using a single addressing scheme and a common routing hierarchy, which allows the switching elements of the network to determine the location of all connected entities. All these connected entities also share access to a common infrastructure of circuits and switching.

However, the model of ubiquity in the “Internet PDN” does not match all potential requirements, especially the need for data privacy. For organizations that wish to use this public network for private purposes within a closed set of participants (for example, connecting a set of geographically separated offices), the Internet is not always a palatable possibility. Numerous factors are behind this mismatch, including issues of Quality of Service (QoS), availability and reliability, use of

public addressing schemes, use of public protocols, site security, and data privacy and integrity (the possibility of traffic interception). Additionally, a corporate network application may desire more stringent levels of performance management than are available within the public Internet, or indeed may wish to define a management regime that differs from that of the underlying Internet PDN.

Service-Level Agreements

It is worthwhile at this point to briefly examine the importance of *Service-Level Agreements* (SLAs) in regards to the deployment of VPNs. SLAs are negotiated contracts between VPN providers and their subscribers; they contain the service criteria to which the subscriber expects specific services to be delivered. The SLA is arguably the only binding tool at the subscriber's disposal with which to ensure that the VPN provider delivers the service(s) to the level and quality as agreed, and it is in the best interest of the subscribers to monitor the criteria outlined in the SLA for compliance. However, SLAs present some challenging technical issues for both the provider and the subscriber.

For the subscriber, the challenge is to devise and operate service measurement tools that can provide a reasonable indication as to what extent the SLA is being honored by the provider. Also, it should be noted that a subscriber may use an SLA to bind one or more providers to a contractual service level, but if the subscriber's VPN spans multiple providers' domains, the SLA must also encompass the issue of provider interconnection and the end-to-end service performance.

For the provider, the challenge lies in honoring multiple SLAs from a number of service providers. In the case of an Internet PDN provider, the common mode of best-effort service levels is not conducive to meeting SLAs, given the unpredictable nature of the host's resource allocation mechanisms. In such environments, the provider either has to ensure that the network is generously engineered in terms of the ratio of subscriber access capacity to internal switching capacity, or the provider can deploy service differentiation structures to ensure that minimum resource levels are allocated to each SLA subscriber. It must be noted that the former course of action does tend to reduce the benefit of aggregation of traffic, which in turn has an ultimate cost implication, while the latter course of action has implications in terms of operational management complexity and scalability of the network.

Alternatives to the VPN

The alternative to using the Internet as a VPN today is to lease circuits, or similar dedicated communications services, from the public network operators (the local telephone company in most cases), and create a completely private network. It is a layering convention that allows us to label this as "completely private," because these dedicated communications services are (at the lower layers of the protocol

stack) again instances of virtual private communications systems constructed atop a common transmission bearer system. Of course, this scenario is not without precedent, and it must be noted that most of the early efforts in data networking, and many of the current data networking architectures, do not assume a deployment model of ubiquitous public access.

It is interesting to note that this situation is odd, when you consider that the inherent value of an architecture where ubiquitous public access over a chaotic collection of closed private networks had been conclusively demonstrated in the telephony marketplace since the start of the 20th century. Although the data communications industry appears to be moving at a considerable technological pace, the level of experiential learning, and consequent level of true progress as distinct from simple motion, still leaves much to be desired!

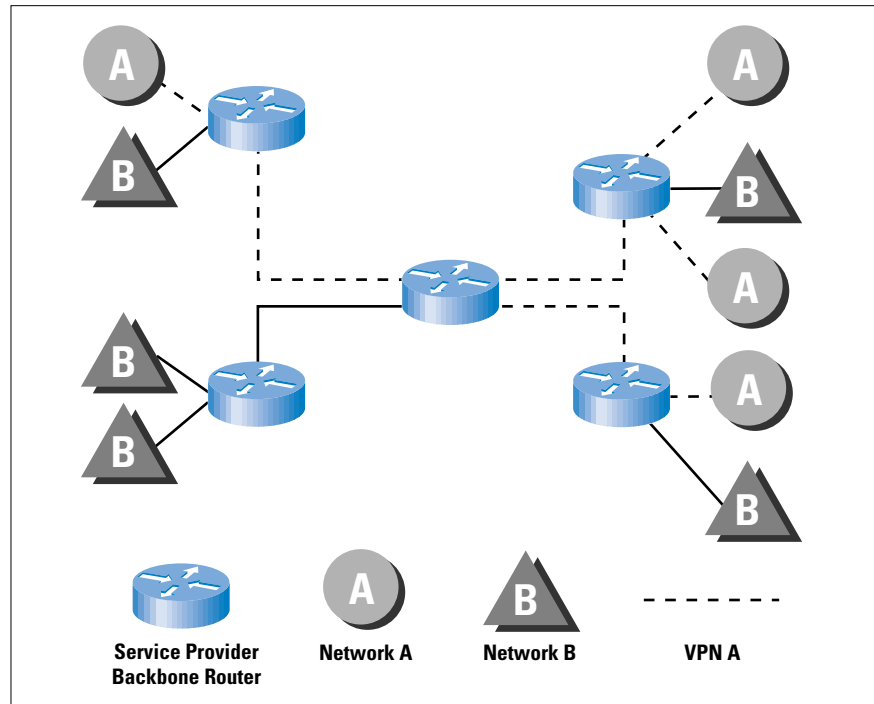
Instead of a public infrastructure deployment, the deployment model used has been that of a closed (or private) network environment where the infrastructure, addressing scheme, management, and services were dedicated to a closed set of subscribers. This model matched that of a closed corporate environment, where the network was dedicated to serve a single corporate entity as the sole client. This precursor to the VPN, which could be called the private data network, was physically constructed using dedicated local office wiring and dedicated leased circuits (or private virtual circuits from an underlying switching fabric such as X.25) to connect geographically diverse sites.

However, this alternative does have an associated cost, in that the client now has to manage the network and all its associated elements, invest capital in network switching infrastructure, hire trained staff, and assume complete responsibility for the provisioning and ongoing maintenance of the network service. Such a dedicated use of transport services, equipment, and staff is often difficult to justify for many small-to-medium sized organizations, and whereas the functionality of a private network system is required, the expressed desire is to reduce the cost of the service through the use of shared transport services, equipment, and management. Numerous scenarios can address this need, ranging from outsourcing the management of the switching elements of the network (managed network services), to outsourcing the capital equipment components (leased network services), to outsourcing the management, equipment, and transport elements to a service provider altogether.

An Example VPN

In the simple example illustrated in Figure 1, Network “A” sites have established a VPN (depicted by the dashed lines) across the service provider’s backbone network, where Network “B” is completely unaware of its existence. Both Networks “A” and “B” can harmoniously coexist on the same backbone infrastructure.

Figure 1:
A Virtual Private
Network of
"A" Sites



This type of VPN is, in fact, the most common type of VPN—one that has geographically diverse subnetworks that belong to a common administrative domain, interconnected by a shared infrastructure outside their administrative control (such as the global Internet or a single service provider backbone). The principal motivation in establishing a VPN of this type is that perhaps most of the communications between devices within the VPN community may be sensitive (again, a decision on the level of privacy required rests solely on a risk analysis performed by the administrators of the VPN), yet the total value of the communications system does not justify the investment in a fully private communications system that uses discrete transmission elements.

On a related note, the level of privacy that a VPN may enjoy depends greatly on the technology used to construct the VPN. For example, if the communications between each VPN subnetwork (or between each VPN host) is securely encrypted as it transits the common communications infrastructure, then it can be said that the privacy aspect of the VPN is relatively high.

In fact, the granularity of a VPN implementation can be broken down further to a single end-to-end, one-to-one connectivity scenario. Examples of these types of one-to-one VPNs are single dialup users who establish a VPN connection to a secure application, such as an online banking service, or a single user establishing a secure, encrypted session between a desktop and server application, such as a purchasing transaction conducted on the World Wide Web. This type of one-to-one VPN is becoming more and more prevalent as secure electronic commerce applications become more mature and are further deployed in the Internet. (See article starting on page 20.)

It is interesting to note that the concept of virtualization in networking has also been considered in regard to deploying both research and production services on a common infrastructure. The challenge in the research and education community is one in which there is a need to satisfy both network research and production requirements. VPNs have also been considered as a method to segregate traffic in a network such that research and production traffic behave as “ships in the night,” oblivious to one another’s existence, to the point that major events (for example, major failures, instability) within one community of interest are completely transparent to the other. This concept is further documented in MORPHnet^[4].

It should also be noted that VPNs may be constructed to span more than one host communications network, so that the “state” of the VPN may be supported on one or more VPN provider networks. This scenario is perhaps at its most robust when all the providers explicitly support the resultant distributed VPN environment, but other solutions that do not necessarily involve knowledge of the overlay VPN are occasionally deployed with mixed results.

Types of VPNs

The confusion factor comes into play in the most basic discussions regarding VPNs, principally because there are actually several different types of VPNs, and depending on the functional requirements, several different methods of constructing each type of VPN are available. The process of selection should include consideration of what problem is being solved, risk analysis of the security provided by a particular implementation, issues of scale in growing the size of the VPN, and the complexity involved in implementation of the VPN, as well as ongoing maintenance and troubleshooting.

To simplify the description of the different types of VPNs, they are broken down in this article into categories that reside in the different layers of the TCP/IP protocol suite; Link Layer, Network Layer, Transport Layer, and Application Layer.

Network-Layer VPNs

The network layer in the TCP/IP protocol suite consists of the IP routing system—how reachability information is conveyed from one point in the network to another. There are a few methods to construct VPNs within the network layer—each is examined in the following paragraphs. A brief overview of non-IP VPNs is provided in Part II of this article.

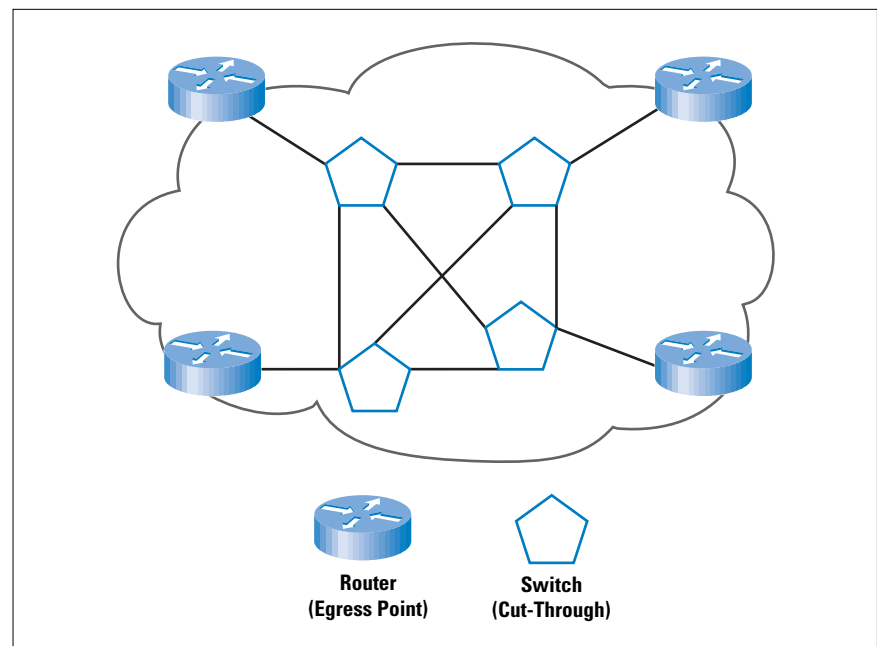
A brief overview of the differences in the “peer” and “overlay” VPN models is appropriate at this point. Simply put, the “peer” VPN model is one in which the network-layer forwarding path computation is done on a hop-by-hop basis, where each node in the intermediate data transit path is a peer with a next-hop node. Traditional routed networks are examples of peer models, where each router in the network

path is a peer with its next-hop adjacencies. Alternatively, the “overlay” VPN model is one in which the network-layer forwarding path is not done on a hop-by-hop basis, but rather, the intermediate link-layer network is used as a “cut-through” to another edge node on the other side of a large cloud. Examples of “overlay” VPN models include ATM, Frame Relay, and tunneling implementations.

Having drawn these simple distinctions between the peer and overlay models, it should be noted that the overlay model introduces some serious scaling concerns in cases where large numbers of egress peers are required because the number of adjacencies increases in direct proportion to the number of peers—the amount of computational and performance overhead required to maintain routing state, adjacency information, and other detailed packet forwarding and routing information for each peer becomes a liability in very large networks. If all the egress nodes in a cut-through network become peers in an effort to make all egress nodes one “Layer 3” hop away from one another, the scalability of the VPN overlay model is limited quite remarkably.

For example, as the simple diagram (Figure 2) illustrates, the routers that surround the interior switched infrastructure represent egress peers, because the switches in the core interior could be configured such that all egress nodes are one Layer 3 hop away from one another, creating what is commonly known as a “cut-through.” This scenario forms the foundation of an overlay VPN model.

Figure 2:
A Cut-Through VPN

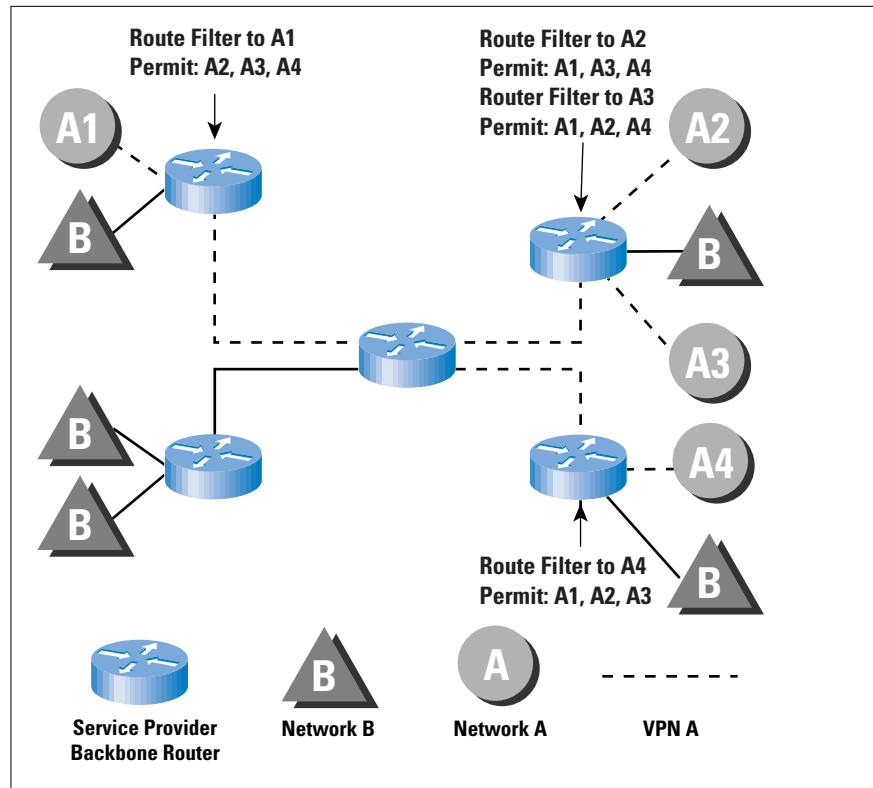


Alternatively, if the switches in the interior are replaced with routers, then the routers positioned at the edge of the cloud become peers with their next-hop router nodes, not other egress nodes. This scenario forms the foundation of the peer VPN model.

Controlled Route Leaking

“Controlled route leaking” (or *route filtering*) is a method that could also be called “privacy through obscurity” because it consists of nothing more than controlling route propagation to the point that only certain networks receive routes for other networks that are within their own community of interest. This model can be considered a “peer” model, because a router within a VPN site establishes a routing relationship with a router within the VPN provider’s network, instead of an edge-to-edge routing peering relationship with routers in other sites of that VPN. Although the common underlying Internet generally carries the routes for all networks connected to it, this architecture assumes that only a subset of such networks form a VPN. The routes associated with this set of networks are filtered such that they are not announced to any other set of connected networks, and all other non-VPN routes are not announced to the networks of the VPN. For example, in Figure 1, if the SP routers “leaked” routing information received from one site in Network “A” to only other sites in Network “A,” then sites not in Network “A” (for instance, sites in Network “B”) would have no explicit knowledge of any other networks which were attached to the service provider’s infrastructure (as shown in Figure 3). Given this lack of explicit knowledge of reachability to any location other than other members of the same VPN, privacy of services is implemented by the inability of any of the VPN hosts to respond to packets which contain source addresses from outside the VPN community of interest.

Figure 3:
Controlled Route
Leaking



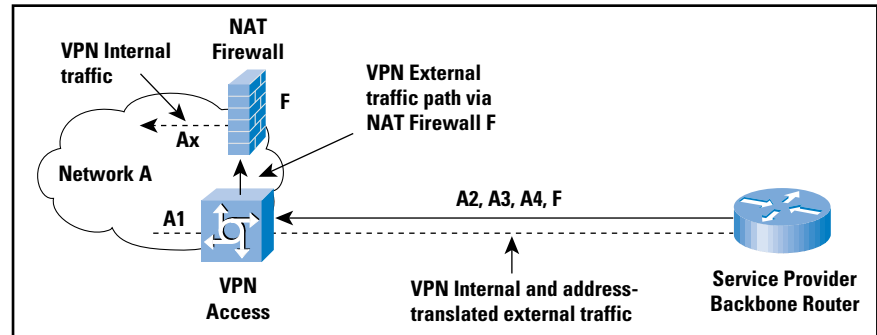
This use of partial routing information is prone to many forms of misconfiguration. One potential problem with route leaking is that it is extremely difficult, if not impossible, to prohibit the subscriber networks from pointing default to the upstream next-hop router for traffic destined for networks outside their community of interest. From within the VPN subscriber's context, this action may be reasonable, in that "default" for the VPN is reachability to all other members of the same VPN, and pointing a default route to the local egress path is, within a local context, a reasonable move. Thus, it is no surprise that this is a common occurrence in VPNs in which the customer configures and manages the customer premise equipment (CPE) routers. If the SP manages the configuration of the CPE routers, then this is rarely a problem. Otherwise, the SP might be wise to place traffic filters on first-hop routers to prohibit all traffic destined for networks outside the VPN community of interest.

It should also be noted that this environment implicitly assumes a common routing core. A common routing core, in turn, implies that each VPN must use addresses that do not clash with those of any other VPN on the same common infrastructure, and cannot announce arbitrary private addresses into the VPN. Another, perhaps less obvious, side effect of this form of VPN structure is that it is not possible for two VPNs to have a single point of interconnection, nor is it possible for a VPN to operate a single point of interconnection to the public Internet in such an environment. (This single point would be a so-called "gateway," where all external traffic is passed through a control point that can enforce some form of access policy and record a log of external transactions.) The common routing core uses a single routing paradigm, based solely on destination address.

It should also be noted that this requirement highlights one of the dichotomies of VPN architectures. VPNs must assume that they operate in a mutually hostile environment, where any vulnerability that exposes the private environment to access by external third parties may be exploited in a hostile fashion. However, VPNs rarely are truly isolated communications environments, and typically all VPNs do have some form of external interface that allows controlled reachability to other VPNs and to the broader public data network. The trade-off between secure privacy and the need for external access is a constant feature of VPNs.

Implementation of inter-VPN connectivity requires the network to route externally originated packets to the VPN interconnection point, and if they are admitted into the VPN at the interconnection point, the same packet may be passed back across the network to the ultimate VPN destination address. Without the use of *Network Address Translation* (NAT) technologies at the interconnection point of ingress into the VPN, this kind of communications structure is insupportable within this architecture (Figure 4).

Figure 4:
Segregating VPN
traffic via address
translation



In general, the technique of supporting private communities of interest simply by route filtering can at best be described as a primitive method of VPN construction, which is prone to administrative errors, and admits an undue level of insecurity and network inflexibility. Even with comprehensive traffic and route filtering, the resulting environment is not totally robust. The operational overhead required to support complementary sets of traditional routing and traffic filters is a relevant consideration, and this approach does not appear to possess the scaling properties desirable to allow the number of VPNs to grow beyond the bounds of a few hundred, using today's routing technologies.

Having said that, however, a much more scalable approach is to use *Border Gateway Protocol* (BGP) *communities*^[5] as a method to control route propagation. The use of BGP communities scales much better than alternative methods with respect to controlling route propagation and is less prone to human misconfiguration. Briefly, the use of the BGP communities attribute allows a VPN provider to “mark” BGP *Network-Layer Reachability Information* (NLRI) with a community attribute, such that configuration control allows route information to propagate in accordance with a community profile.

Because traffic from different communities of interest must traverse a common shared infrastructure, there is no significant data privacy in the portion of the network where traffic from multiple communities of interest share the infrastructure. Therefore, it can be said that although connected subnetworks—or rather, subscribers to the VPN service—may not be able to detect the fact that there are other subscribers to the service, multiple interwoven streams of subscriber data traffic pass unprotected in the core of the service provider's network.

Tunneling

Sending specific portions of network traffic across a tunnel is another method of constructing VPNs. Some tunneling methods are more effective than others. The most common tunneling mechanisms are *Generic Routing Encapsulation* (GRE)^[6] tunneling between a source and destination router, router-to-router or host-to-host tunneling protocols such as *Layer 2 Tunneling Protocol* (L2TP)^[7] and *Point-to-Point Tunneling Protocol* (PPTP)^[8], and *Distance Vector Multicast Routing Protocol* (DVMRP)^[9] tunnels.

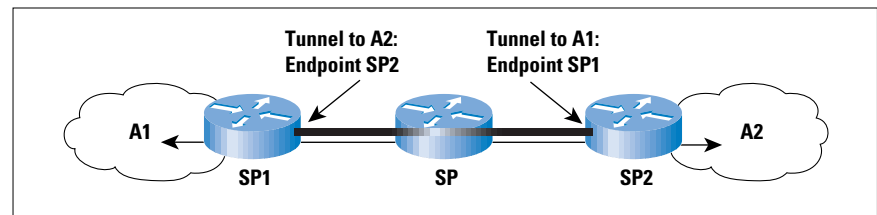
Tunneling can be considered an overlay model, but the seriousness of the scaling impact depends on whether the tunnels are point-to-point or point-to-multipoint. Point-to-point tunnels have fewer scaling problems than do point-to-multipoint tunnels, except in situations where a single node begins to build multiple point-to-point tunnels with multiple endpoints. Although a linear scaling problem is introduced at this point, the manageability of point-to-point tunnels lies solely in the administrative overhead and the number of the tunnels themselves. On the other hand, point-to-multipoint tunnels use “cut-through” mechanisms to make greater numbers of endpoints one hop away from one another and subsequently introduce a much more serious scaling problem.

Although the *Multicast Backbone* (Mbone) itself could literally be considered a global VPN, and although DVMRP tunnels are still widely used by organizations to connect to the Mbone, it really is not germane to the central topic of VPNs, because the focus of this article is on unicast traffic.

Traditional Modes of Tunneling

GRE tunnels, as mentioned previously, are generally configured between a source (*ingress*) router and a destination (*egress*) router, such that packets designated to be forwarded across the tunnel (already formatted with an encapsulation of the data with the “normal” protocol-defined packet header) are further encapsulated with a new header (the GRE header), and placed into the tunnel with a destination address of the tunnel endpoint (the new next-hop). When the packet reaches the tunnel endpoint, the GRE header is stripped away, and the packet continues to be forwarded to the destination, as designated in the original IP packet header (Figure 5).

Figure 5:
Tunneling across a
Service Provider



GRE tunnels are generally point-to-point—that is, there is a single source address for the tunnel and usually only a single destination tunnel endpoint. However, some vendor implementations allow the configuration of point-to-multipoint tunnels—that is, a single source address and multiple destinations. Although this implementation is generally used in conjunction with *Next Hop Resolution Protocol* (NHRP)^[10], the effectiveness and utility of NHRP is questionable and should be tested prior to deployment. It is also noteworthy that NHRP is known to produce steady-state forwarding loops when used to establish shortcuts between routers. In the scenario discussed previously, NHRP is used for establishing shortcuts between routers.

Tunnels, however, do have numerous compelling attractions when used to construct VPNs. The architectural concept is to create VPNs as a collection of tunnels across a common host network. Each point of attachment to the common network is configured as a physical link that uses addressing and routing from the common host network, and one or more associated tunnels. Each tunnel endpoint logically links this point of attachment to other remote points from the same VPN. The technique of tunneling uses a tunnel egress address defined within the address space of the common host network, whereas the packets carried within the tunnel use the address space of the VPN, which in turn constrains the tunnel endpoints to be collocated to those points in the network where the VPN and the host network interconnect.

Pros and Cons

The advantage of this approach is that the routing for the VPN is isolated from the routing of the common host network. The VPNs can reuse the same private address space within multiple VPNs without any cross impact, providing considerable independence of the VPN from the host network. This requirement is key for many VPNs in that private VPNs typically may not use globally unique or coordinated address space, and there is often the consequential requirement to support multiple VPNs which independently use the same address block. Such a configuration is not supportable within a controlled route leakage VPN architecture. The tunnel can also encapsulate numerous different protocol families, so that it is possible for a tunnel-based VPN to mimic much of the functionality of dedicated private networks. Again, the need to support multiple protocols in a format which preserves the functionality of the protocol is a critical requirement for many VPN support architectures. This requirement is one in which an IP common network with controlled route leakage cannot provide such services, whereas a tunneling architecture can segment the VPN-private protocol from the common host network. The other significant advantage of the tunneled VPN is the segregation of the common host routing environment with that of the VPN. To the VPN, the common host network assumes the properties of numerous point-to-point circuits, and the VPN can use a routing protocol across the virtual network which matches the administrative requirements of the VPN. Equally, the common host network can use a routing design which matches the administrative requirements of the host network (or collection of host networks), and is not constrained by the routing protocols used by the VPN client networks.

Although it could be said that these advantages indicate that GRE tunneling is the panacea for VPN design, using GRE tunnels as a mechanism for VPNs does have several drawbacks, mostly with regard to administrative overhead, scaling to large numbers of tunnels, and QoS and performance.

Since GRE tunnels must be manually configured, there is a direct relationship to the number of tunnels that must be configured and the amount of administrative overhead required to configure and maintain them—each time the tunnel endpoints must change, and they must be manually reconfigured. Also, although the amount of processing required to encapsulate a packet for GRE handling may appear to be small, there is a direct relationship to the number of configured tunnels and the total amount of processing overhead required for GRE encapsulation. Of course, tunnels can be structured to be triggered automatically, but such an approach has numerous drawbacks that dictate careful consideration of related routing and performance issues. The worst end state of such automatic tunnel generation is that of a configuration loop where the tunnel passes traffic over itself. It is important, once again, to reiterate the impact of a large number of routing peering adjacencies that result from a complete mesh of tunnels; this scenario can result in a negative effect on routing efficiency.

An additional concern with GRE tunneling is the ability of traffic classification mechanisms to identify traffic with a fine enough level of granularity, and not become a hindrance to forwarding performance. If the traffic classification process used to identify packets (that are to be forwarded across the tunnel) interferes with the router's ability to maintain acceptable packet-per-second forwarding rates, then this becomes a performance liability.

Privacy of the network remains an area of concern because the tunnel is still vulnerable—privacy is not absolute. Packets that use GRE formatting can be injected into the VPN from third-party sources. To ensure a greater degree of integrity of privacy of the VPN, it is necessary to deploy ingress filters that are aligned to the configured tunnel structure.

It is also necessary to ensure that the CPE routers are managed by the VPN service provider, because the configuration of the tunnel endpoints is a critical component of the overall architecture of integrity of privacy. However, most VPN service providers are reluctant to add CPE equipment to their asset inventory and undertake remote management of such CPE equipment, due to the high operational overheads and poor capital efficiencies which are typical of CPE deployment. Arguably, one might suggest that having a dedicated CPE router defeats one of the basic premises of constructing a VPN—the use of shared infrastructure as a way to reduce the overall network cost.

It should be noted that VPNs can be constructed using tunnels without the explicit knowledge of the host network provider, and the VPN can span numerous host networks without any related underlying agreements between the network operators to mutually support the overlay VPN. Such an architecture is little different from provider-operated VPN architecture; the major difference lies in the issue of traffic and

performance engineering, and the administrative boundary of the management of the VPN overlay. Independently configured VPN tunnels can result in injection of routes back into the VPN in a remote location, a scenario that can cause traffic to traverse the same link twice, once in an unencapsulated format and again within a tunnel. This situation can then lead to adverse performance impacts.

It is also true that the overlay VPN model has no control over which path is taken in the common host network, nor the stability of that path. This scenario can then lead to adverse performance impacts on the VPN. Aside from the technology aspects of this approach, the major issue is one of whether the VPN management is outsourced to the network provider, or undertaken within administrative functions of the VPN. One of the more serious considerations in building a VPN on tunneling is that there is virtually no way to determine the cost of the route across a tunnel, because the true path is masked by the cut-through nature of the tunnel. This situation could ultimately result in highly suboptimal routing, meaning that a packet could take a path determined by the cut-through mechanism that is excessively suboptimal, while native per-hop routing protocols might find a much more efficient method to forward the packets to their destinations.

Conclusion

So far in our discussion of VPNs, we have introduced a working definition of the term “Virtual Private Network” and discussed the motivations behind the adoption of such networks. We have outlined a framework for describing the various forms of VPNs, and then examined numerous network-layer VPN structures, in particular, that of controlled route leakage and tunneling techniques.

In Part II we will continue this examination of network-layer VPNs, including virtual private dial networks and network-layer encryption. In addition, we will examine link-layer VPNs that use ATM and Frame Relay substrates, and also look at switching and encryption techniques, and issues concerning QoS and non-IP VPNs.

References

- [1] Steinberg, Steve G., “Hype List—Deflating this month’s overblown memes.” *Wired Magazine*, 6.02, February 1998, p. 80. Ironically, number 1 on the Hype List is virtual private networks, with a life expectancy of 18 months.
- [2] Raymond, Eric S., compiler. *The New Hacker’s Dictionary, Third Edition*. MIT Press, ISBN 0-262-68092-0, 1996. The Jargon File online: <http://www.ccil.org/jargon/>
- [3] *Webster’s Revised Unabridged Dictionary* (1913). Hypertext Webster Gateway: http://work.ucsd.edu:5141/cgi-bin/http_webster

- [4] Aiken, R., R. Carlson, I. Foster, T. Kuhfuss, R. Stevens, and L. Winkler. "Architecture of the Multi-Modal Organizational Research and Production Heterogeneous Network (MORPHnet)," Argonne National Laboratory, ECT and MCS Divisions, January 1997.
<http://www.anl.gov/ECT/Public/research/morphnt2.htm>
- [5] Chandra, R., P. Traina, and T. Li. RFC 1997, "BGP Communities Attribute." August 1996; E. Chen and T. Bates. RFC 1998, "An Application of the BGP Community Attribute in Multi-home Routing." August 1996.
- [6] Hanks, S., T. Li, D. Farinacci, and P. Traina. RFC 1701, "Generic Routing Encapsulation." October 1994; S. Hanks, T. Li, D. Farinacci, and P. Traina. RFC 1702, "Generic Routing Encapsulation over IPv4 networks." October 1994.
- [7] Valencia, A., K. Hamzeh, A. Rubens, T. Kolar, M. Littlewood, W. M. Townsley, J. Taarud, G. S. Pall, B. Palter, and W. Verthein. "Layer Two Tunneling Protocol 'L2TP.'" `draft-ietf-pppext-l2tp-10.txt`, March 1998.
- [8] Hamzeh, K., G. Singh Pall, W. Verthein, J. Taarud, and W. A. Little. "Point-to-Point Tunneling Protocol—PPTP." `draft-ietf-pppext-pptp-02.txt`, July 1997.
See also: <http://www.microsoft.com/backoffice/communications/morepptp.htm>
- [9] Waitzman, D., C. Partridge, and S. Deering. RFC 1075, "Distance Vector Multicast Routing Protocol." November 1988. For historical purposes, see also <ftp://ftp.isi.edu/mbone/faq.txt>
- [10] Luciani, J., D. Katz, D. Piscitello, B. Cole, and N. Doraswamy. "NBMA Next Hop Resolution Protocol (NHRP)," `draft-ietf-rolc-nhrp-15.txt`, February 1998.

PAUL FERGUSON is a consulting engineer at Cisco Systems and an active participant in the Internet Engineering Task Force (IETF). His principal areas of expertise include large-scale network architecture and design, global routing, Quality of Service (QoS) issues, and Internet Service Providers. Prior to his current position at Cisco Systems, he worked in network engineering, analytical, and consulting capacities for Sprint, Computer Sciences Corporation (CSC), and NASA. He is coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, published by John Wiley & Sons, ISBN 0-471-24358-2, a collaboration with Geoff Huston. E-mail: ferguson@cisco.com

GEOFF HUSTON holds a B.Sc and a M.Sc from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the the initial build of the Internet within the Australian academic and research sector. Huston is currently the Chief Technologist in the Internet area for Telstra. He is also an active member of the IETF, and was an inaugural member of the Internet Society Board of Trustees. He is coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, published by John Wiley & Sons, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. E-mail: gih@telstra.net

SSL: Foundation for Web Security

by William Stallings

Virtually all businesses, most government agencies, and many individuals now have Web sites. The number of individuals and companies with Internet access is expanding rapidly, and all of them have graphical Web browsers. As a result, businesses are enthusiastic about setting up facilities on the Web for electronic commerce. But the reality is that the Internet and the Web are extremely vulnerable to compromises of various sorts. As businesses utilize the Internet for more than information dissemination, they will need to use trusted security mechanisms.

An increasingly popular general-purpose solution is to implement security as a protocol that sits between the underlying transport protocol (TCP) and the application. The foremost example of this approach is the *Secure Sockets Layer* (SSL) and the follow-on Internet standard of SSL known as *Transport Layer Security* (TLS). At this level, there are two implementation choices. For full generality, SSL (or TLS) could be provided as part of the underlying protocol suite and therefore be transparent to applications. Alternatively, SSL can be embedded in specific packages. For example, Netscape and Microsoft Explorer browsers come equipped with SSL, and most Web servers have implemented the protocol. Although it is possible to use SSL for applications other than Web transactions, its use at present is typically as part of Web browsers and servers and hence limited to Web traffic. Most of this article deals with the technical details of SSL; the status of TLS is described at the end.

If you have viewed an HTML source document, you have seen that the links are referenced with `HREF=<URL>` within an anchor (A) tag. In most cases, the reference is to another document through the use of the *Hyper Text Transfer Protocol*, or HTTP. For this, the browser initiates one or more sessions to the destination port of TCP/80 (the well-known port for HTTP) on the server. In some cases, a plug-in can be called, and data specific to that plug-in can be transferred to or from the browser. For that, the browser would initiate a session to the well-known TCP port of the plug-in. SSL is called when the reference starts like the following: `HREF="https://..` By calling "https" within the browser, it is mandating that the data be transferred through the use of SSL. By clicking on this hot link, the browser initiates a session to the server on port TCP/443. SSL attempts to negotiate a secure link and transfers the data across it. If the negotiation fails, no data is transferred. The browser usually indicates that a secure connection has been requested. Netscape Navigator version 3 indicates this with a blue border around the page and a highlighted key in the lower left corner. Netscape Communicator version 4 displays this with a closed padlock in a lower status window. Microsoft's Internet Explorer indicates it

with a padlock in a lower information window. Display of these signs indicates that the information within the browser window has been delivered through the security of SSL.

SSL was originated by Netscape. Version 3 of the protocol was designed with public review and input from industry and was published as an Internet Draft document. Subsequently, when a consensus was reached to submit the protocol for Internet standardization, the TLS working group was formed within the *Internet Engineering Task Force* (IETF) to develop a common standard. The current work on TLS is aimed at producing an initial version as an Internet Standard. This first version of TLS can be viewed as essentially an SSLv3.1, and is very close to SSLv3. TLS includes a mechanism by which a TLS entity can back down to the SSLv3.0 protocol; in that sense, TLS is backward compatible with SSL.

SSL Architecture

SSL is designed to make use of TCP to provide a reliable end-to-end secure service. SSL is not a single protocol but rather two layers of protocols.

The SSL Record Protocol provides basic security services to various higher-layer protocols. In particular, the HTTP, which provides the transfer service for Web client/server interaction, can operate on top of SSL. Three higher-layer protocols are defined as part of SSL: the *Handshake Protocol*, the *Change CipherSpec Protocol*, and the *Alert Protocol*. These SSL-specific protocols are used in the management of SSL exchanges.

Two important SSL concepts are the SSL session and the SSL connection, which are defined in the specification as follows:

- **Connection:** A logical client/server link that provides a suitable type of service. For SSL, such connections are peer-to-peer relationships. The connections are transient. Every connection is associated with one session.
- **Session:** An association between a client and a server. Sessions are created by the Handshake Protocol. Sessions define a set of cryptographic security parameters, which can be shared among multiple connections. Sessions are used to avoid the expensive negotiation of new security parameters for each connection.

Between any pair of parties (applications such as HTTP on client and server), there may be multiple secure connections. In theory, there may also be multiple simultaneous sessions between parties, but this feature is not used in practice.

Several states are associated with each session. When a session is established, there is a current operating state for both read and write (that is, receive and send). In addition, during the Handshake Protocol, pending read and write states are created. Upon successful conclusion of the Handshake Protocol, the pending states become the current states. A session state is defined by the following parameters (definitions taken from the SSL specification):

- Session identifier: An arbitrary byte sequence chosen by the server to identify an active or resumable session state.
- Peer certificate: An X509.v3 certificate of the peer. This element of the state may be null.
- Compression method: The algorithm used to compress data prior to encryption.
- CipherSpec: Specifies the bulk data encryption algorithm (such as DES) and a hash algorithm (such as MD5 or SHA-1). It also defines cryptographic attributes such as the hash size.
- Master secret: 48-byte secret shared between the client and server.
- Is resumable: A flag indicating whether the session can be used to initiate new connections.

A connection state is defined by the following parameters:

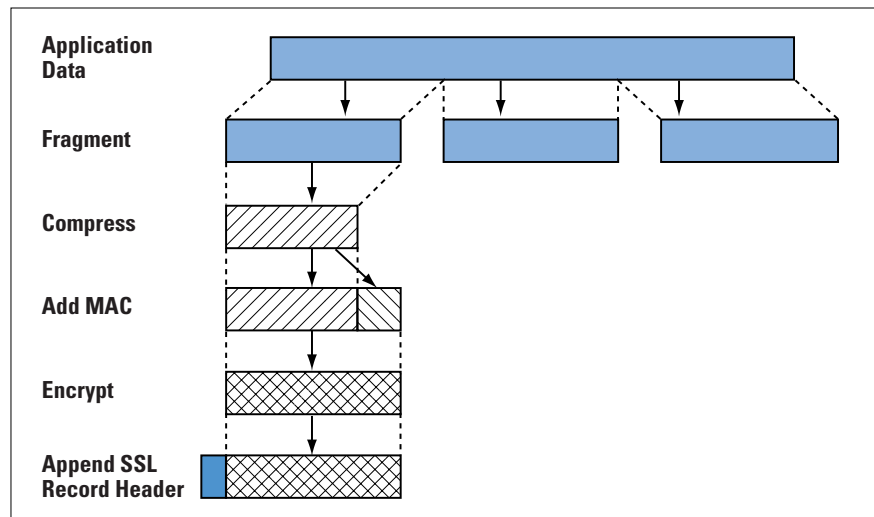
- Server and client random: Byte sequences that are chosen by the server and client for each connection.
- Server write MAC secret: The secret key used in MAC operations on data sent by the server.
- Client write MAC secret: The secret key used in MAC operations on data sent by the client.
- Server write key: The conventional encryption key for data encrypted by the server and decrypted by the client.
- Client write key: The conventional encryption key for data encrypted by the client and decrypted by the server.
- Initialization vectors: When a block cipher in CBC mode is used, an initialization vector (IV) is maintained for each key. This field is first initialized by the SSL Handshake Protocol. Thereafter the final ciphertext block from each record is preserved for use as the IV for the next record.
- Sequence numbers: Each party maintains separate sequence numbers for transmitted and received messages for each connection. When a party sends or receives a change CipherSpec message, the appropriate sequence number is set to zero.

SSL Record Protocol

The SSL Record Protocol provides two services for SSL connections: confidentiality, by encrypting application data; and message integrity, by using a *message authentication code* (MAC). The Record Protocol is a base protocol that can be utilized by some of the upper-layer protocols of SSL. One of these is the handshake protocol which, as described later, is used to exchange the encryption and authentication keys. It is vital that this key exchange be invisible to anyone who may be watching this session.

Figure 1 indicates the overall operation of the SSL Record Protocol. The Record Protocol takes an application message to be transmitted, fragments the data into manageable blocks, optionally compresses the data, applies a MAC, encrypts, adds a header, and transmits the resulting unit in a TCP segment. Received data is decrypted, verified, decompressed, and reassembled and then delivered to the calling application, such as the browser.

Figure 1:
SSL Record Protocol
Operation



The first step is fragmentation. Each upper-layer message is fragmented into blocks of 2^{14} bytes (16,384 bytes) or less. Next, compression is optionally applied. In SLLv3 (as well as the current version of TLS), no compression algorithm is specified, so the default compression algorithm is null. However, specific implementations may include a compression algorithm.

The next step in processing is to compute a message authentication code over the compressed data. For this purpose, a shared secret key is used. In essence, the hash code (for example, MD5) is calculated over a combination of the message, a secret key, and some padding. The receiver performs the same calculation and compares the incoming MAC value with the value it computes. If the two values match, the receiver is assured that the message has not been altered in transit. An attacker would not be able to alter both the message and the MAC, because the attacker does not know the secret key needed to generate the MAC.

Next, the compressed message plus the MAC are encrypted using symmetric encryption. A variety of encryption algorithms may be used, including the Data Encryption Standard (DES) and triple DES.

The final step of SSL Record Protocol processing is to prepend a header, consisting of the following fields:

- **Content Type (8 bits):** The higher-layer protocol used to process the enclosed fragment.
- **Major Version (8 bits):** Indicates major version of SSL in use. For SSLv3, the value is 3.
- **Minor Version (8 bits):** Indicates minor version in use. For SSLv3, the value is 0.
- **Compressed Length (16 bits):** The length in bytes of the plain-text fragment (or compressed fragment if compression is used).

The content types that have been defined are `change_cipher_spec`, `alert`, `handshake`, and `application_data`. The first three are the SSL-specific protocols, mentioned previously. The application-data type refers to the payload from any application that would normally use TCP but is now using SSL, which in turn uses TCP. In particular, the HTTP protocol that is used for Web transactions falls into the application-data category. A message from HTTP is passed down to SSL, which then wraps this message into an SSL record.

Change CipherSpec Protocol

The Change CipherSpec Protocol is one of the three SSL-specific protocols that use the SSL Record Protocol, and it is the simplest. This protocol consists of a single message, which consists of a single byte with the value 1. The sole purpose of this message is to cause the pending state to be copied into the current state, which updates the CipherSuite to be used on this connection. This signal is used as a coordination signal. The client must send it to the server and the server must send it to the client. After each side has received it, all of the following messages are sent using the agreed-upon ciphers and keys.

Alert Protocol

The Alert Protocol is used to convey SSL-related alerts to the peer entity. As with other applications that use SSL, alert messages are compressed and encrypted, as specified by the current state.

Each message in this protocol consists of two bytes. The first byte takes the value “warning” (1) or “fatal”(2) to convey the severity of the message. If the level is fatal, SSL immediately terminates the connection. Other connections on the same session may continue, but no new connections on this session may be established. The second byte contains a code that indicates the specific alert. An

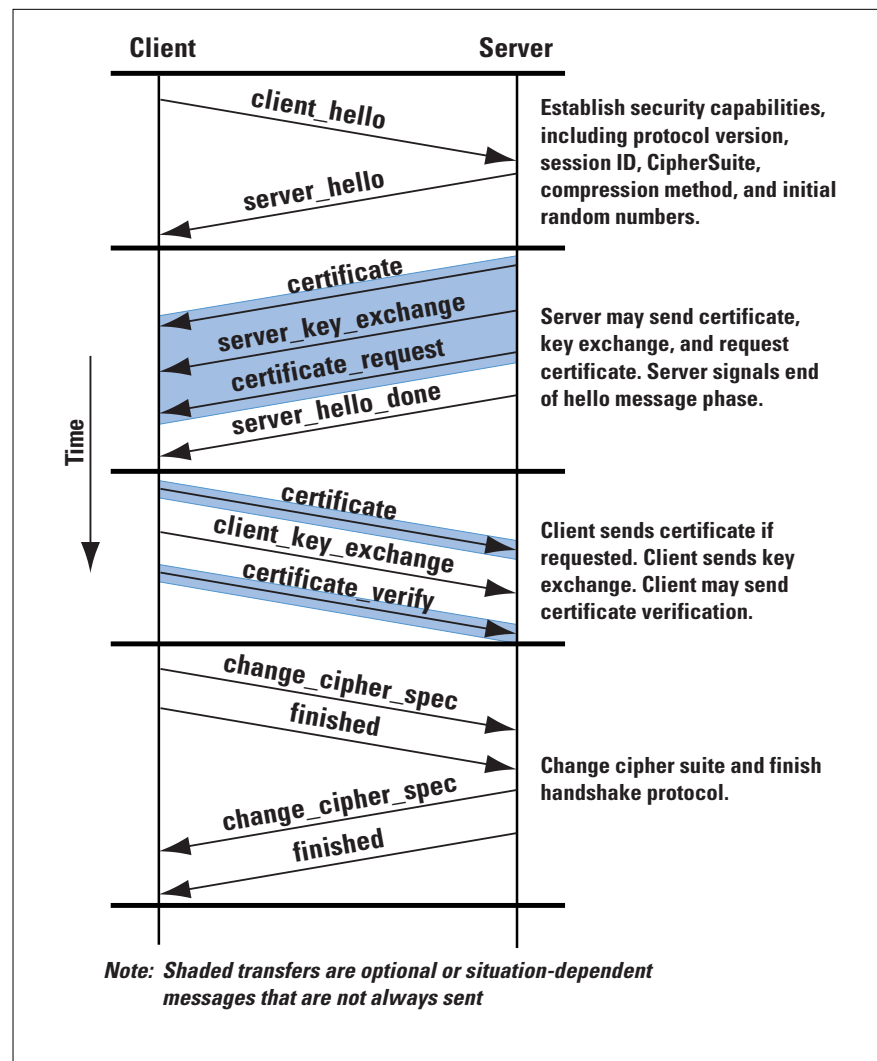
example of a fatal message is `illegal_parameter` (a field in a handshake message was out of range or inconsistent with other fields). An example of a warning message is `close_notify` (notifies the recipient that the sender will not send any more messages on this connection; each party is required to send a `close_notify` alert before closing the write side of a connection).

Handshake Protocol

The most complex part of SSL is the Handshake Protocol. This protocol allows the server and client to authenticate each other and to negotiate an encryption and MAC algorithm and cryptographic keys to be used to protect data sent in an SSL record. The Handshake Protocol is used before any application data is transmitted. The Handshake Protocol consists of a series of messages exchanged by the client and the server.

Figure 2 shows the initial exchange needed to establish a logical connection between the client and the server. The exchange can be viewed as having four phases.

Figure 2:
Handshake Protocol
Action



Phase 1 is used to initiate a logical connection and to establish the security capabilities that will be associated with it. The exchange is initiated by the client, which sends a `client_hello` message with the following parameters:

- **Version:** The highest SSL version understood by the client.
- **Random:** A client-generated random structure, consisting of a 32-bit timestamp and 28 bytes generated by a secure random number generator. These values serve as nonces and are used during key exchange to prevent replay attacks.
- **Session ID:** A variable-length session identifier. A nonzero value indicates that the client wishes to update the parameters of an existing connection or create a new connection on this session. A zero value indicates that the client wishes to establish a new connection on a new session.
- **CipherSuite:** A list that contains the combinations of cryptographic algorithms supported by the client, in decreasing order of preference. Each element of the list (each CipherSuite) defines both a key exchange algorithm and a CipherSpec; these are discussed subsequently.
- **Compression Method:** A list of the compression methods the client supports.

After sending the `client_hello` message, the client waits for the `server_hello` message, which contains the same parameters as the `client_hello` message. For the `server_hello` message, the following conventions apply. The Version field contains the lower of the version suggested by the client and the highest version supported by the server. The Random field is generated by the server and is independent of the client's Random field. If the SessionID field of the client was nonzero, the same value is used by the server; otherwise the server's SessionID field contains the value for a new session. The CipherSuite field contains the single CipherSuite selected by the server from those proposed by the client. The Compression field contains the compression method selected by the server from those proposed by the client.

The first element of the CipherSuite parameter is the key exchange method (that is, the means by which the cryptographic keys for conventional encryption and MAC are exchanged). The following key exchange methods are supported:

- **RSA:** The secret key is encrypted with the receiver's RSA public key. A public-key certificate for the receiver's key must be made available.
- **Fixed Diffie-Hellman:** This is a Diffie-Hellman key exchange in which the server's certificate contains the Diffie-Hellman public parameters

signed by the *certificate authority* (CA). That is, the public-key certificate contains the Diffie-Hellman public-key parameters. The client provides its Diffie-Hellman public key parameters either in a certificate, if client authentication is required, or in a key exchange message. This method results in a fixed secret key between two peers, based on the Diffie-Hellman calculation using the fixed public keys.

- **Ephemeral Diffie-Hellman:** This technique is used to create ephemeral (temporary, one-time) secret keys. In this case, the Diffie-Hellman public keys are exchanged, and signed using the sender's private RSA or DSS key. The receiver can use the corresponding public key to verify the signature. Certificates are used to authenticate the public keys. This option appears to be the most secure of the three Diffie-Hellman options because it results in a temporary, authenticated key.
- **Anonymous Diffie-Hellman:** The base Diffie-Hellman algorithm is used, with no authentication. That is, each side sends its public Diffie-Hellman parameters to the other, with no authentication. This approach is vulnerable to man-in-the-middle attacks, in which the attacker conducts anonymous Diffie-Hellman exchanges with both parties.

Following the definition of a key exchange method is the Cipher-Spec, which indicates the encryption and hash algorithms and other related parameters.

The server begins Phase 2 by sending its certificate, if it needs to be authenticated; the message contains one or a chain of X.509 certificates. The certificate message is required for any agreed-on key exchange method except anonymous Diffie-Hellman. Note that if fixed Diffie-Hellman is used, this certificate message functions as the server's key exchange message because it contains the server's public Diffie-Hellman parameters.

Next, a `server_key_exchange` message may be sent, if it is required. It is not required in two instances: (1) The server has sent a certificate with fixed Diffie-Hellman parameters; or (2) RSA key exchange is to be used.

Next, a nonanonymous server (server not using anonymous Diffie-Hellman) can request a certificate from the client. The `certificate_request` message includes two parameters: `certificate_type` and `certificate_authorities`. The certificate type indicates the type of public-key algorithm. The second parameter in the `certificate_request` message is a list of the distinguished names of acceptable certificate authorities.

The final message in Phase 2, and one that is always required, is the `server_done` message, which is sent by the server to indicate the end of the server hello and associated messages. After sending this message, the server waits for a client response. This message has no parameters.

Upon receipt of the `server_done` message, the client should verify that the server provided a valid certificate, if required, and check that the server hello parameters are acceptable. If all is satisfactory, the client sends one or more messages back to the server in Phase 3. If the server has requested a certificate, the client begins this phase by sending a certificate message. If no suitable certificate is available, the client sends a `no_certificate` alert instead.

Next is the `client_key_exchange` message, which must be sent in this phase. The content of the message depends on the type of key exchange.

Finally, in this phase, the client may send a `certificate_verify` message to provide explicit verification of a client certificate. This message is only sent following any client certificate that has signing capability (that is, all certificates except those containing fixed Diffie-Hellman parameters).

Phase 4 completes the setting up of a secure connection. The client sends a `change_cipher_spec` message and copies the pending CipherSpec into the current CipherSpec. Note that this message is not considered part of the Handshake Protocol but is sent using the Change CipherSpec Protocol. The client then immediately sends the finished message under the new algorithms, keys, and secrets. The finished message verifies that the key exchange and authentication processes were successful.

In response to these two messages, the server sends its own `change_cipher_spec` message, transfers the pending to the current CipherSpec, and sends its finished message. At this point the handshake is complete and the client and server may begin to exchange application layer data.

After the records have been transferred, the TCP session is closed. However, since there is no direct link between TCP and SSL, the state of SSL may be maintained. For further communications between the client and the server, many of the negotiated parameters are retained. This may occur if, in the case of Web traffic, the user clicks on another link that also specifies HTTPs on the same server. If the clients or servers wish to resume the transfer of records, they don't have to again negotiate encryption algorithms or totally new keys. The SSL specifications suggest that the state information be cached for no longer than 24 hours. If no sessions are resumed within that time, all information is deleted and any new sessions have to go through the handshake again. The specifications also recommend that neither the client nor the server have to retain this information, and shouldn't if either of them suspects that the encryption keys have been compromised. If either the client or the server does not agree to resume the session, for any reason, then both will have to go through the full handshake.

Transport Layer Security

TLS is an IETF standardization initiative whose goal is to produce an Internet standard version of SSL. In fact, the charter for the TLS working group states:

“The TLS working group is a focused effort on providing security features at the transport layer, rather than general purpose security and key management mechanisms. The standard track protocol specification will provide methods for implementing privacy, authentication, and integrity above the transport layer.”

This means that TLS can be used to provide security services to any application that uses TCP or the *User Datagram Protocol* (UDP). However, the driving force behind this work is to develop a standardized version of SSL. Microsoft has indicated that TLS will go into the next major version of its browser and Web server products, and Netscape has made a similar commitment. With this kind of support, it is likely that TLS will move quickly along the Internet Standards track.

The current draft version of TLS is very similar to SSLv3. TLS uses slightly different cryptographic algorithms for such things as the MAC function generation of secret keys. TLS also includes more alert codes.

SSL is already widely deployed and, under the name TLS, is moving toward Internet standardization. It is the solution of choice for Web transaction security.

References

- [1] <http://www.phaos.com/sslresource.html>
(has links to vendors, SSL specifications, and FAQs)
- [2] <http://www.netscape.com/newsref/std/SSL.html>
(PostScript versions of the spec are available there)
- [3] <http://www.ietf.org/html.charters/tls-charter.html>
(contains latest RFCs and Internet Drafts for TLS)
- [4] <http://www.imc.org/ietf-tls/mail-archive/>
(mailing list archive)
- [5] <ftp://ftp.ietf.org/internet-drafts/draft-ietf-tls-protocol-05.txt>
- [6] <ftp://ftp.ietf.org/internet-drafts/draft-ietf-tls-https-01.txt>
- [7] <http://www.consensus.com/security/ssl-talk-faq.html>

WILLIAM STALLINGS is a consultant, lecturer, and author of over a dozen books on data communications and computer networking. He has a PhD in computer science from M.I.T. This article is based on material in the author's latest book, *Cryptography and Network Security, Second Edition* (Prentice-Hall, 1998). His home in cyberspace is <http://www.shore.net/~ws> and he can be reached at ws@shore.net

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal will carry tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It will provide readers with technology and standardization updates for all levels of the protocol stack and serve as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and quality of service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

Book Reviews

Groupware *Groupware: Collaborative Strategies for Corporate LANs and Intranets*, by David Coleman, ISBN 0-13-727728-8, Prentice-Hall PTR, 1997, <http://www.prenhall.com>.

Some areas of science provide very poor training for dealing with primarily human processes. One might think that packet switching would be an exception because it lives on the stochastic nature of bursty communications. Because our knowledge of human and group activity is, at best, characterized by statistical assessments, those working in networking should do well in understanding and dealing with the unpredictable and human nature of communication, especially when it involves using networks.

So much for theory. In general, the world of lower-level networking has done little for the upper strata of computer-mediated human communication, except to provide a platform for the work of others. An apparent exception in the world of Internet technology is e-mail, yet it actually serves more as proof of the problem than as an exception. The basic facilities in Internet e-mail are the same today as they were 25 years ago. As nice as they are, the word “basic” is essential when characterizing them. Almost none of the Internet’s standardized e-mail facilities are really targeted at providing automated or structural support for the work of a group.

Groupware Defined

The collection of products and services designed to help people collaborate via computer, by direct interaction, or by information dissemination is called “groupware.” Coleman’s book is a revision of *Groupware: Technology and Applications*. Written only 15 months earlier, the world changed more than enough in that time to require the revision. The first book had relatively little to say about the Internet, whereas this new book tries mightily to factor it into the equation. The result is a bit erratic, but the digressions serve to highlight how rapidly things are changing, rather than to suggest looking elsewhere for a better source on the topic.

The new book has an entirely different subtitle, giving a reasonable sense that the content targets more an understanding of system organization and function than detailed technical explanation. That’s just fine, because the book really is not particularly technical. It covers the requirements and functions for supporting activity by groups.

Downsizing and working remotely are two very strong driving forces for increased use of groupware. This book is essentially an introduction to concepts, functionality, and use of systems that attempt to help staff members work together. Oddly, that does not only mean working together when physically separated, because there is discussion of meeting room assistance, such as with automated sense-of-the-group tallying devices.

Organization

The first two chapters introduce the topic, emphasizing that human and group process concerns dominate the field and are intimately tied to the aggressive efforts that organizations are making to run more productively and, frequently, with fewer people. The third chapter discusses functionality in terms of the World Wide Web. The book reflects the current enthusiasm for the Web, sometimes to the detriment of the appropriate use of messaging technology, although messaging is more prevalent among groupware than other kinds of commercial Internet systems.

The realm of groupware does not have a firm taxonomy. My own synthesis includes: Message (text and document) Exchange, Forms Exchange, Calendaring & Scheduling, Workflow, Presentations and Interactive Meetings, and Document Development and Sharing. The next six chapters cover the functional pieces of this groupware realm.

The next five chapters cover the major vendors of integrated groupware products: Lotus Notes, Novel GroupWise, TeamWARE, Hewlett-Packard, and Oracle Interoffice. HP's chapter discusses "strategy," suggesting the lack of a well-integrated product suite, but one more survey of the terrain is nonetheless useful. And that, perhaps, is the major reason for reading this book: It constantly emphasizes the human and process-oriented aspect of organizational behavior and the need to attend carefully both to the needs of the humans and the nature of the processes. It is easy to understand that an improper travel authorization, will bring an organization to its knees. It is easy to forget that the system is used by humans who well might not want the added complexity or rigidity of the system and who, therefore, must be part of the design and adoption effort. In my opinion, the book takes a rather more negative view about groupware acceptability than is necessary, but then I like such technology, and the average worker in the average organization does not.

The last six chapters of this book intermix case studies and Hahn, of Collabra and Netscape, points the reader to Chapter 17, "Groupware & Reengineering: The Human Side of Change." Although one of the better considerations of these issues in the book, it is far from the only one.

A Useful Survey

If you have little familiarity with these "upper level application" areas of networking, the functionality, products, or use, then this book is a good one to read. You will not learn much about the underlying technology, nor will you be able to qualify as a "certified groupware support engineer," but you will obtain an extremely useful survey of the field, and you will obtain it from the perspective of human and organization use. As the Internet moves into the mass market, that perspective is a good one.

—Dave Crocker
Brandenburg Consulting
dcrocker@brandenburg.com

High-Speed Networks *High-Speed Networks: TCP/IP and ATM Design Principles*,
by William Stallings, ISBN 0-13-525965-7 Prentice-Hall, 1997,
<http://www.shore.net/~ws/HsNet.html>

High-speed networks now dominate both the WAN and LAN markets. In the WAN market, data networks have evolved from packet-switching networks to ATM networks operating at 155 Mbps or more. In the LAN market, the staple 10-Mbps Ethernet is being replaced with 100-Mbps Fast Ethernet, Gigabit Ethernet, and even Asynchronous Transfer Mode (ATM) LANs. This book provides a survey of high-speed networks and the design issues related to them. Much of the book is devoted to the study of various techniques aimed at reducing network congestion.

Organization

The book is divided into seven sections. The first section deals with the fundamentals: TCP/IP principles; packet switching and Frame Relay networks; and internetworking principles. The second section provides an overview of ATM and Fast and Gigabit Ethernet. These two sections can easily be torn out of the book and serve as an excellent primer on today's modern networks. I am going to recommend to my employer that they be made mandatory reading.

In the third section of the book, Stallings focuses on one treatment of queueing theory, namely, how it is applied to modeling network behavior. Stallings has an undeniable gift for taking large complicated subjects and teaching the fundamentals, and then some, without belittling the subject at hand or the reader. This book is witness to this gift, and this chapter but one fine example. But once the reader has an understanding of queueing theory, Stallings throws a wrench in the gears. The chapter on self-similarity explains why traditional queueing models are inadequate when trying to predict the performance of Ethernet traffic and other self-similar streams. While this section is by far the most theoretical, it is at the same time necessary for the reader's understanding of network performance, and while many readers may not care to devote the time necessary to gain a complete understanding of self-similarity, astute students are urged to invest in more than a simple gloss-over of this section.

Having understood the basics of self-similarity, I hoped the fifth section of the book, on network traffic management, would be addressed with greater emphasis on delivering quality of service and the problems related to self-similarity. Instead, the material is based on traditional queueing models.

The fourth section, flow control, is divided into two categories. The first, link control mechanisms, focuses on some of the performance issues related to the use of *Automatic Repeat Request* (ARQ) link control protocols. The second category, transport control mechanisms,

concentrates on the TCP flow control mechanism. I expected to find references to bugs in some TCP implementations exposed by high-volume WWW servers, but didn't. Stallings goes on to present an overview of some of the performance issues of TCP over ATM. As institutions begin upgrading their networks, this issue is sure to receive a great deal of interest. The section concludes with a look at the *Real-Time Transport Protocol*, another area sure to spark attention as the need to move large multimedia data across WANs, in real time, becomes more relevant.

The sixth section of the book covers Internet routing protocols and opens with a primer on graph theory. Four routing protocols (RIP, OSPF, BGP, and IDRP) are covered. The section concludes with a discussion of multicasting as an introduction to RSVP. This section sparked my curiosity enough to call for a visit to the WWW site for RSVP development.

Stallings shies away from directly addressing application-driven improvements aimed at increasing network performance. In today's Web/CGI-driven world, I would expect this to be a topic of interest to many. Perhaps this is a subject for another book. But the topic is not entirely avoided. The last section of the book focuses on various lossless and lossy compression techniques. The quirkiess of material covered makes this section a darling.

Recommended

This book rates an A+. Unlike most books about computers being published today, this book is neither superficial nor is it insulting to the reader. It is intended for both professional and academic audiences. Stallings' desire to truly educate is apparent. This is not a book about promoting the hype, this is a book about serious learning.

—Neophytos Iacovou,
University of Minnesota
Academic & Distributed Computing Services
iacovou@boombox.micro.umn.edu

Fragments

The Fragments page is intended to provide you with updates and pointers to information related to Internet technology developments.

The Future of the Domain Name System (DNS)

For more than a year, a debate has taken place regarding the future of the DNS. In particular, the issue of competitive name registries, possible addition of new *global Top Level Domains* (gTLDs) and the future of the *Internet Assigned Numbers Authority* (IANA) have been discussed. Information regarding the initial proposal can be found at: <http://www.gtld-mou.org/>. The US Government has issued a so-called *Green Paper* entitled “Technical Management of Internet Names and Addresses.” The Green Paper and comments received on this document can be found at:

<http://www.ntia.doc.gov/ntiahome/domainname/>

IETF and Related Links

The *Internet Engineering Task Force* (IETF) is responsible for the development of standards for Internet technology. Membership to the IETF is open and you can participate in person or subscribe to the IETF mailing list. The IETF meets three times per year. For a list of future meetings and other IETF information see: <http://www.ietf.org>. On this website you will also find a number of links to organizations which are related to the IETF in one way or another:

- *The Internet Society* (ISOC) and its annual INET conference.
- *The Internet Architecture Board* (IAB)
- *The Internet Assigned Numbers Authority* (IANA)
- *The Internet Research Task Force* (IRTF)

SIGCOMM

If you want to learn about the latest developments on the research side of networking you should check out SIGCOMM, the Association for Computing Machinery’s Special Interest Group on Communications. You can find out more about the group and their annual conference at: <http://www.acm.org/sigcomm/sigcomm98>

Send Us Your Comments!

We look forward to hearing your comments and suggestions regarding anything you read in this publication. Send e-mail to: ipj@cisco.com.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or noninfringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Engineering
MCI Communications, USA

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Edward R. Kozel, Sr. VP, Corporate Development
Cisco Systems, Inc., USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President,
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the Cisco News
Publications Group, Cisco Systems, Inc.
www.cisco.com*

*Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 1998 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-J4
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

Bulk Rate Mail U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol Journal

September 1998

Volume 1, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

F R O M T H E E D I T O R

In This Issue

From the Editor 1

What Is a VPN?—Part II 2

Reliable Multicast Protocols
and Applications 19

Layer 2 and Layer 3
Switch Evolution 38

Book Review 44

Fragments 47

We begin this issue with Part II of “What Is a VPN?” by Paul Ferguson and Geoff Huston. In Part I they introduced a definition of the term “Virtual Private Network” (VPN) and discussed the motivations behind the adoption of such networks. They outlined a framework for describing the various forms of VPNs, and examined numerous network-layer VPN structures, in particular, that of controlled route leakage and tunneling. In Part II the authors conclude their examination of VPNs by describing virtual private dial networks and network-layer encryption. They also examine link-layer VPNs, switching and encryption techniques, and issues concerning Quality of Service and non-IP VPNs.

IP Multicast is an emerging set of technologies and standards that allow many-to-many transmissions such as conferencing, or one-to-many transmissions such as live broadcasts of audio and video over the Internet. Kenneth Miller describes multicast in general, and reliable multicast protocols and applications in particular. Although multicast applications are primarily used in the research community today, this situation is likely to change as the demand for Internet multimedia applications increases and multicast technologies improve.

Successful deployment of networking technologies requires an understanding of a number of technology options ranging from wiring and transmissions systems via switches, routers, bridges and other pure networking components, to networked applications and services. *The Internet Protocol Journal* (IPJ) is designed to look at all aspects of these “building blocks.” This time, Thayumanavan Sridhar details some of the issues in the evolution of Layer 2 and Layer 3 switches.

Interest in the first issue of IPJ has exceeded our expectations, and hard copies are almost gone. However, you can still view and print the issue in PDF format on our Web site at www.cisco.com/ipj. The current edition is also available on the Web. If you want to receive our next issue, please complete and return the enclosed card.

We welcome your comments, questions and suggestions regarding anything you read in this journal. We are also actively seeking authors for new articles. The Call for Papers and Author Guidelines can be found on our Web page. Please send your comments to ipj@cisco.com

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

Missed the first issue of IPJ?
Download your copy in
PDF format from:
www.cisco.com/ipj

What Is a VPN? — Part II

by Paul Ferguson, Cisco Systems
and Geoff Huston, Telstra

In Part I we introduced a working definition of the term “Virtual Private Network” (VPN), and discussed the motivations behind the adoption of such networks. We outlined a framework for describing the various forms of VPNs, and then examined numerous network-layer VPN structures, in particular, that of controlled route leakage and tunneling techniques. We begin Part II with examining other network-layer VPN techniques, and then look at issues that are concerned with non-IP VPNs and Quality-of-Service (QoS) considerations.

Types of VPNs

This section continues from Part I to look at the various types of VPNs using a taxonomy derived from the layered network architecture model. These types of VPNs segregate the VPN network at the network layer.

Network-Layer VPNs

A network can be segmented at the network layer to create an end-to-end VPN in numerous ways. In Part I we described a controlled route leakage approach that attempts to perform the segregation only at the edge of the network, using route advertisement control to ensure that each connected network received a view of the network (only peer networks). We pick up the description at this point in this second part of the article.

Tunneling

As outlined in Part I, the alternative to a model of segregation at the edge is to attempt segregation throughout the network, maintaining the integrity of the partitioning of the substrate network into VPN components through the network on a hop-by-hop basis. Part I examined numerous tunneling technologies that can achieve this functionality. Tunneling is also useful in servicing VPN requirements for dial access, and we will resume the description of tunnel-based VPNs at this point.

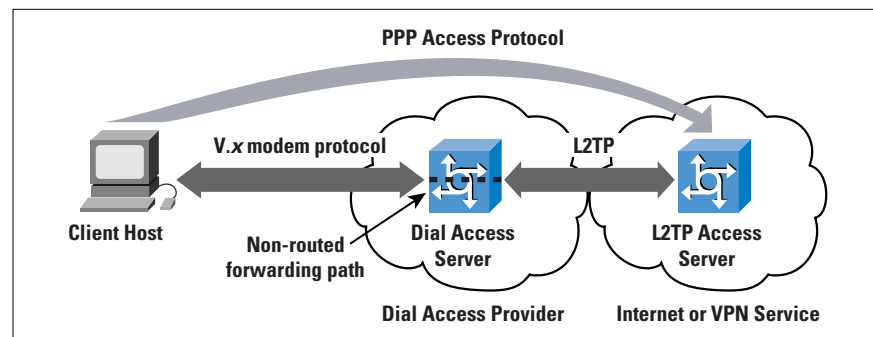
Virtual Private Dial Networks

Although several technologies (vendor-proprietary technologies as well as open, standards-based technologies) are available for constructing a *Virtual Private Dial Network* (VPDN), there are two principal methods of implementing a VPDN that appear to be increasing in popularity—*Layer 2 Tunneling Protocol* (L2TP) and *Point-to-Point Tunneling Protocol* (PPTP) tunnels. From an historical perspective, L2TP is the technical convergence of the earlier Layer 2 Forwarding (L2F)^[1] protocol specification and the PPTP protocol. However, one might suggest that because PPTP is now being bundled into the desktop operating system of many of the world’s personal computers, it stands to be quite popular within the market.

At this point it is worthwhile to distinguish the difference between “client-initiated” tunnels and “NAS-initiated” (Network Access Server, otherwise known as a Dial Access Server) tunnels. The former is commonly referred to as “voluntary” tunneling, whereas the latter is commonly referred to as “compulsory” tunneling. In voluntary tunneling, the tunnel is created at the request of the user for a specific purpose; in compulsory tunneling, the tunnel is created without any action from the user, and without allowing the user any choice in the matter.

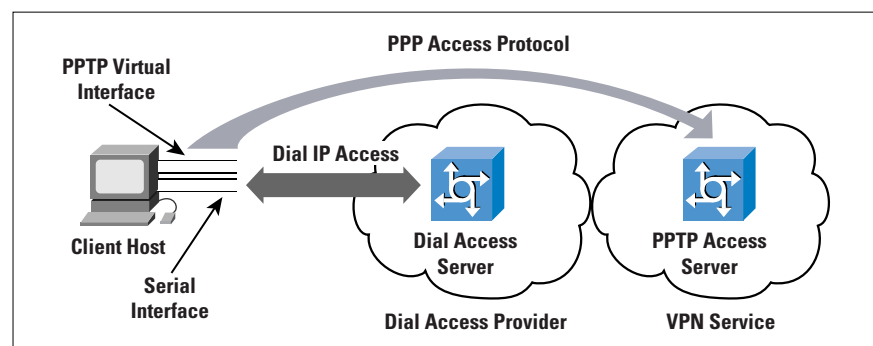
L2TP, as a compulsory tunneling model, is essentially a mechanism to “off-load” a dialup subscriber to another point in the network, or to another network altogether. In this scenario, a subscriber dials into a NAS, and based on a locally configured profile (or a NAS negotiation with a policy server) and successful authentication, a L2TP tunnel is dynamically established to a predetermined endpoint, where the subscriber’s *Point-to-Point Protocol* (PPP) session is terminated (Figure 1).

Figure 1:
PPP Tunnel
Termination Model
of L2TP



PPTP, as a voluntary tunneling model, on the other hand, allows end systems (for example, desktop computers) to configure and establish individual discrete point-to-point tunnels to arbitrarily located PPTP servers, without the intermediate NAS participating in the PPTP negotiation and subsequent tunnel establishment. In this scenario, a subscriber dials into a NAS, but the PPP session is terminated on the NAS, as in the traditional Internet access PPP model. The layered PPTP session is then established between the client end system and any upstream PPTP server that the client desires to connect to. The only caveats on PPTP connectivity are that the client can reach the PPTP server via conventional routing processes, and that the user has been granted the appropriate privileges on the PPTP server (Figure 2).

Figure 2:
PPP Tunnel
Termination Model
of PPTP



Although L2TP and PPTP may sound extraordinarily similar, there are subtle differences that deserve further examination. The applicability of both protocols is very much dependent on what problem is being addressed. It is also about control—who has it, and why it is needed. It also depends heavily on how each protocol implementation is deployed—in either the voluntary or the compulsory tunneling models.

With PPTP in a voluntary tunneling implementation, the dial-in user can choose the PPTP tunnel destination (the PPTP server) after the initial PPP negotiation has completed. This feature is important if the tunnel destination changes frequently, because no modifications are needed to the client's view of the base PPP access when there is a change in the server and the transit path to the server. It is also a significant advantage that the PPTP tunnels are transparent to the service provider, and no advance configuration is required between the NAS operator and the overlay dial access VPN. In such a case, the service provider does not house the PPTP server, and simply passes the PPTP traffic along with the same processing and forwarding policies as all other IP traffic. In fact, this feature should be considered a significant benefit of this approach. The configuration and support of a tunneling mechanism within the service provider network would be one less parameter that the service provider has to operationally manage, and the PPTP tunnel can transparently span multiple service providers without any explicit service provider configuration. However, the economic downside to this feature for the service provider, of course, is that a "VPDN-enabled" network service can be marketed to yield an additional source of revenue. Where the client undertakes the VPDN connection, there is no direct service provider involvement and no consequent value added to the base access service.

From the subscriber's perspective, this is a "win-win" situation, because the user is not reliant on the upstream service provider to deliver the VPDN service—at least no more than any user is reliant for basic IP-level connectivity. The other "win" is that the subscriber does not have to pay a higher subscription fee for a VPN service. Of course, the situation changes when the service provider takes an active role in providing the VPDN, such as housing the PPTP servers, or if the subscriber resides within a subnetwork in which the parent organization wants the service provider's network to make the decision concerning where tunnels are terminated. The major characterization of PPTP-based VPDN is one of a roaming client base, where the clients of the VPDN use a local connection to the public Internet data network, and then overlay a private data tunnel from the client's system to the desired remote service point. Another perspective is to view this approach as "on-demand" VPDN virtual circuits.

With L2TP in a "compulsory" tunneling implementation, the service provider controls where the PPP session is terminated. This setup can be extremely important in situations where the service provider to whom

the subscriber is actually dialing into (let's call it the "modem pool provider" network) must transparently hand off the subscriber's PPP session to another network (let's call this network the "content provider"). To the subscriber, it appears as though the local system is directly attached to the content provider's network, when in fact the access path has been passed transparently through the modem pool provider's network to the subscribed content service. Very large content providers, for instance, may outsource the provisioning and maintenance of thousands of modem ports to a third-party access provider, who in turn agrees to transparently pass the subscribers' access sessions back to the content provider. This setup is generally called "wholesale dial." The major motivation for such L2TP-based wholesale dial lies in the typical architecture of the *Public Switched Telephone Network* (PSTN), where the use of wholesale dial facilities can create a more rational PSTN call load pattern with Internet access PSTN calls terminated in the local Central Office.

Of course, if all subscribers who connect to the modem pool provider's network are destined for the same content provider, then there are certainly easier ways to hand this traffic off to the content provider's network—such as simply aggregating all the traffic in the local Central Office and handing the content provider a "big fat pipe" of the aggregated session traffic streams. However, in situations where the modem pool provider is providing a wholesale dial service for multiple upstream "next-hop" networks, the methods of determining how each subscriber's traffic must be forwarded to his/her respective content provider are somewhat limited. Packet forwarding decisions could be made at the NAS, based on the source address of the dialup subscriber's computer. This scenario would allow for traffic to be forwarded along the appropriate path to its ultimate destination, in turn intrinsically providing a virtual connection. However, the use of assigning static IP addresses to dial-in subscribers is highly discouraged because of the inefficiencies in IP address utilization policies, and the critical success of the *Dynamic Host Configuration Protocol* (DHCP).

There are, however, some serious scaling concerns in deploying a large-scale L2TP network; these concerns revolve around the issue of whether large numbers of tunnels can actually be supported with little or no network performance impact. Since there have been no large-scale deployments of this technology to date, there is no empirical evidence to support or invalidate these concerns.

In some cases, however, appearances are everything—some content providers do not wish for their subscribers to know that when they connect to their service, they have instead been connected to another service provider's network, and then passed along ultimately to the service to which they have subscribed. In other cases, it is merely designed to be a matter of convenience, so that subscribers do not need to log into a device more than once.

Regrettably, the L2TP draft does not detail all possible implementations or deployment scenarios for the protocol. The basic deployment scenario is quite brief when compared to the rest of the document, and is arguably biased toward the compulsory tunneling model. Nonetheless, there are implementations of L2TP that follow the voluntary tunneling model. To the best of our knowledge, there has never been any intent to exclude this model of operation. In addition, at various recent interoperability workshops, several different implementations of a voluntary L2TP client have been modeled. Nothing in the L2F protocol would prohibit deploying it in a voluntary tunneling manner, but to date it has not been widely implemented. Further, PPTP has also been deployed using the compulsory model in a couple of specific vendor implementations.

In summary, consideration of whether PPTP or L2TP is more appropriate for deployment in a VPDN depends on whether control needs to lie with the service provider or with the subscriber. Indeed, the difference can be characterized with respect to the client of the VPN, where the L2TP model is one of a “wholesale” access provider who has numerous configured client service providers who appear as VPNs on the common dial access system, whereas the PPTP model is one of distributed private access where the client is an individual end user and the VPN structure is that of end-to-end tunnels. One might also suggest that the difference is also a matter of economics, because the L2TP model allows service providers to actually provide a “value-added” service, beyond basic IP-level connectivity, and charge their subscribers accordingly for the ability to access it, thus creating new revenue streams. By contrast, the PPTP model enables distributed reach of the VPN at a much more basic level, enabling corporate VPNs to extend access capabilities without the need for explicit service contracts with a multitude of network access providers.

Network-Layer Encryption

Encryption technologies are extremely effective in providing the segmentation and virtualization required for VPN connectivity, and they can be deployed at almost any layer of the protocol stack. The evolving standard for network-layer encryption in the Internet is *IP Security* (IPSec)^[3, 4]. (IPSec is actually an architecture—a collection of protocols, authentication, and encryption mechanisms. The IPSec security architecture is described in detail in [3].)

While the *Internet Engineering Task Force* (IETF) is finalizing the architecture and the associated protocols of IPSec, there is relatively little network-layer encryption being done in the Internet today. However, some vendor proprietary solutions are currently in use.

Whereas IPSec has yet to be deployed in any significant volume, it is worthwhile to review the two methods in which network-layer encryption is predominantly implemented. The most secure method for network-

layer encryption to be implemented is end-to-end, between participating hosts. End-to-end encryption allows for the highest level of security. The alternative is more commonly referred to as “tunnel mode,” in which the encryption is performed only between intermediate devices (routers), and traffic between the end system and the first-hop router is in plaintext. This setup is considerably less secure, because traffic intercepted in transit between the first-hop router and the end system could be compromised.

As a more general observation on this security vulnerability, where a VPN architecture is based on tunnels, the addition of encryption to the tunnel still leaves the tunnel ingress and egress points vulnerable, because these points are logically part of the host network as well as being part of the unencrypted VPN network. Any corruption of the operation, or interception of traffic in the clear, at these points will compromise the privacy of the private network.

In the end-to-end encryption scheme, VPN granularity is to the individual end-system level. In the tunnel mode scheme, the VPN granularity is to the subnetwork level. Traffic that transits the encrypted links between participating routers, however, is considered secure. Network-layer encryption, to include IPSec, is merely a subset of a VPN.

Link-Layer VPNs

One of the most straightforward methods of constructing VPNs is to use the transmission systems and networking platforms for the physical and link-layer connectivity, yet still be able to build discrete networks at the network layer. A link-layer VPN is intended to be a close (or preferably exact) functional analogy to a conventional private data network.

ATM and Frame Relay Virtual Connections

A conventional private data network uses a combination of dedicated circuits from a public carrier, together with an additional private communications infrastructure, to construct a network that is completely self-contained. Where the private data network exists within private premises, the network generally uses a dedicated private wiring plant to carry the VPN. Where the private data network extends outside the private boundary of the dedicated circuits, it is typically provisioned for a larger public communications infrastructure by using some form of time-division or frequency-division multiplexing to create the dedicated circuit. The essential characteristic of such circuits is the synchronization of the data clock, such that the sender and receiver pass data at a clocking rate that is fixed by the capacity of the dedicated circuit.

A link-layer VPN attempts to maintain the critical elements of this self-contained functionality, while achieving economies of scale and operation, by utilizing a common switched public network infrastructure. Thus, a collection of VPNs may share the same infrastructure for connectivity, and share the same switching elements within the interior of

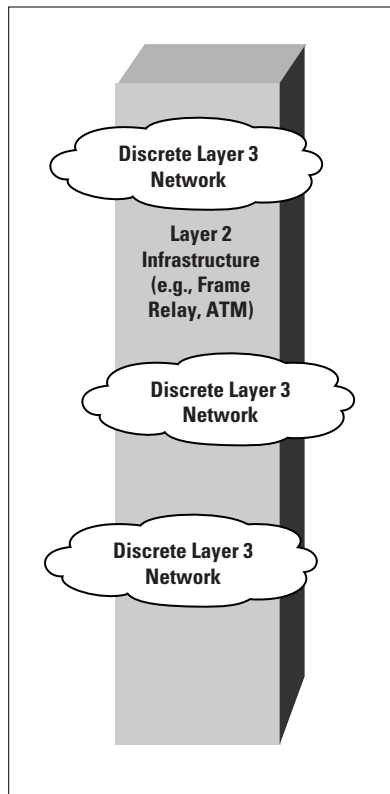


Figure 3:
Conceptualization of
Discrete Layer 3
Networks on a
Common Layer 2
Infrastructure

the network, but explicitly must have no visibility, either direct or inferred, of one another. Generally, these “networks” operate at Layer 3 (the network layer) or higher in the OSI Reference Model, and the “infrastructure” itself commonly consists of either a *Frame Relay* or *Asynchronous Transfer Mode (ATM)* network (Figure 3). The essential difference here between this architecture of virtual circuits and that of dedicated circuits is that there is now no synchronized data clock shared by the sender and receiver, nor necessarily is there a dedicated transmission path that is assigned from the underlying common host network. The sender generally has no a priori knowledge of the available capacity of the virtual circuit, because the capacity varies in response to the total demand placed on it by other simultaneous transmission and switching activity. Instead, the sender and receiver can use adaptive clocking of data, where the sender can adjust the transmission rate to match the requirements of the application and any signaling received from the network and the receiver. It should be noted that a dedicated circuit system using synchronized clocking cannot be oversubscribed, whereas the virtual circuit architecture (where the sender does not have a synchronized end-to-end data clock) can indeed be oversubscribed. It is the behavior of the network when it transitions into this oversubscribed state that is of most interest here.

One of the nice things about a public switched wide-area network that provides virtual circuits is that it can be extraordinarily flexible. Most subscribers to Frame Relay services, for example, have subscribed to the service for economic reasons—it is cheap, and the service provider usually adds a *Service-Level Agreement (SLA)* that “guarantees” some percentage of frame delivery in the Frame Relay network itself.

The remarkable thing about this service offering is that the customer is generally completely unaware of whether the service provider can actually deliver the contracted service at all times and under all possible conditions. The Layer 2 technology is not a synchronized clock blocking technology in which each new service flow is accepted or denied based on the absolute ability to meet the associated resource demands. Each additional service flow is accepted into the network and carried on a best-effort basis. Admission functions provide the network with a simple two-level discard mechanism that allows a graduated response to instances of overload; however, when the point of saturated overload is reached within the network, all services will be affected.

This situation brings up several other important issues: The first concerns the engineering practices of the Frame Relay service provider. If the Frame Relay network is poorly engineered and is constantly congested, then obviously the service quality delivered to the subscribers will be affected. Frame Relay uses a notion of a per-virtual circuit *Committed Information Rate (CIR)*, which is an ingress function associated with Frame Relay that checks the ingress traffic rate against the CIR.

Frames that exceed this base rate are still accepted by the Frame Relay network, but they are marked as *discard eligible* (DE). Because the network can be oversubscribed, the data rate within a switch will at times exceed both the egress transmission rate and the local buffer storage. When this situation occurs, the switch will begin to discard data frames, and will do so initially for frames with the DE marker present. This scenario is essentially a two-level discard precedence architecture. It is an administrative decision by the service provider as to the relative levels of provisioning of core transmission and switching capacity, and the ratio of network ingress capacity used by subscribers. The associated CIRs of the virtual circuits against this core capacity are critical determinants of the resultant deliverable quality of performance of the network and the layered VPNs.

For example, at least one successful (and popular) Frame Relay service provider provides an economically attractive Frame Relay service that permits a zero-rate CIR on PVCs, combined with an SLA that ensures that at least 99.8 percent of all frame-level traffic presented to the Frame Relay network will be delivered successfully. If this SLA is not met, then the subscriber's monthly service fee will be appropriately prorated the following month. The Frame Relay service provider provides frame level statistics to each subscriber every month, culled from the Frame Relay switches, to measure the effectiveness of this SLA "guarantee." This particular Frame Relay service provider is remarkably successful in honoring the SLAs because they conduct ongoing network capacity management on a weekly basis, provisioning new trunks between Frame Relay switches when trunk utilization exceeds 50 percent, and ensuring that trunk utilization never exceeds 75 percent. In this fashion, traffic on PVCs with a zero-rate CIR can generally avoid being discarded in the Frame Relay network.

Having said that, the flexibility of PVCs allows discrete VPNs to be constructed across a single Frame Relay network. And in many instances, this scenario lends itself to situations where the Frame Relay network provider also manages each discrete VPN via a telemetry PVC. Several service providers have *Managed Network Services* (MNS) that provide exactly this type of service.

Whereas the previous example revolves around the use of Frame Relay as a link-layer mechanism, essentially the same type of VPN mechanics hold true for ATM. As with Frame Relay, there is no data clock synchronization between the sender, the host network, and the receiver. In addition, the sender's traffic is passed into the ATM network via an ingress function, which can mark cells with a *Cell Loss Priority* (CLP) indication. And, as with Frame Relay, where a switch experiences congestion, the switch will attempt to discard marked (CLP) cells as the primary load shedding mechanism, but if this step is inadequate, the network must shed other cells that are not so marked. Once again, the quality of the service depends on proper capacity engineering of the network, and there is no guarantee of service quality inherently in the technology itself.

The generic observation is that the engineering of Frame Relay and ATM common carriage data networks is typically very conservative. The inherent capabilities of both of these link-layer architectures do not permit a wide set of selective responses to network overload, so that in order for the network to service the broadest spectrum of potential VPN clients, the network must provide high-quality carriage and very limited instances of any form of overload. In this way, such networks are typically positioned as a high-quality alternative to dedicated circuit private network architectures, which are intended to operate in a very similar manner (and, not surprisingly, are generally priced as a premium VPN offering). Technically, the architecture of link-layer VPNs is almost indistinguishable from the dedicated circuit private data network—the network can support multiple protocols, private addressing, and routing schemes, because the essential difference between a dedicated circuit and a virtual link-layer circuit is the absence of synchronized clocking between the sender and the receiver. In all other aspects, the networks are very similar.

These approaches to constructing VPNs certainly involve scaling concerns, especially with regard to configuration management of provisioning new *Virtual Connections* (VCs) and routing issues. Configuration management still tends to be one of the controversial points in VPN management—adding new subscribers and new VPNs to the network requires VC path construction and provisioning, a tedium that requires ongoing administrative attention by the VPN provider. Also, as already mentioned, full mesh networks encounter scaling problems, in turn resulting in construction of VPNs in which partial meshing is done to avoid certain scaling limitations. The liabilities in these cases need to be examined closely, because partial meshing of the underlying link-layer network may contribute to suboptimal routing (for example, extra hops caused by hub-and-spoke issues, or redirects).

These problems apply to all types of VPNs built on the “overlay” model—not just ATM and Frame Relay. Specifically, the problems also apply to *Generic Routing Encapsulation* (GRE) tunnels.

MPOA and the “Virtual Router” Concept

Another unique model of constructing VPNs is the use of *Multiprotocol over ATM* (MPOA)^[5], which uses RFC 1483 encapsulation^[6]. This VPN approach is similar to other “cut-through” mechanisms in which a particular switched link layer is used to enable all “Layer 3” egress points to be only a single hop away from one another.

In this model, the edge routers determine the forwarding path in the ATM switched network, because they have the ability to determine which egress point packets need to be forwarded to. After a network-layer reachability decision is made, the edge router forwards the packet onto a VC designated for a particular egress router. However, since the egress routers cannot use the *Address Resolution Protocol* (ARP) for destination address across the cloud, they must rely on an external server for address resolution (ATM address to IP address).

The first concern here is a sole reliance on ATM—this particular model does not encompass any other types of data link layer technologies, rendering the technology less than desirable in a hybrid network. Whereas this scenario may have some domain of applicability within a homogenous ATM environment, when looking at a broader VPN environment that may encompass numerous link-layer technologies, this approach offers little benefit to the VPN provider.

Secondly, there are serious scaling concerns regarding full mesh models of connectivity, where suboptimal network-layer routing may result because of cut-through. And the reliance on address resolution servers to support the ARP function within the dynamic circuit framework brings this model to the point of excessive complexity.

The advantage of the MPOA approach is the use of dynamic circuits rather than more cumbersome, statically configured models. The traditional approach to supporting private networks involves extensive manual design and operational support to ensure that the various configurations on each of the bearer switching elements are mutually consistent. The desire within the MPOA environment is to attempt to use MPOA to govern the creation of dynamically controlled, edge-to-edge ATM VCs. Although this setup may offer the carrier operator some advantages in reduced design and operational overhead, it does require the uniform availability of ATM, and in many heterogeneous environments this scenario is not present.

In summary, this model is another overlay model, with some serious concerns regarding the ability of the model to withstand scale.

“Peer” VPN models that allow the egress nodes to maintain separate routing tables have also been introduced—one for each VPN—effectively allowing separate forwarding decisions to be made within each node for each distinctive VPN. Although this is an interesting model, it introduces concerns about approaches in which each edge device runs a separate routing process and maintains a separate *Routing Information Base* (RIB, or routing table) process for each VPN community of interest. It also should be noted that the “virtual router” concept requires some form of packet labeling, either within the header or via some lightweight encapsulation mechanism, in order for the switch to be able to match the packet against the correct VPN routing table. If the label is global, the issue of operational integrity is a relevant concern, whereas if the label is local, the concept of label switching and maintenance of edge-to-edge label switching contexts is also a requirement.

Among the scaling concerns are issues regarding the number of supported VPNs in relation to the computational requirements, and stability of the routing system within each VPN (that is, instability in one VPN affecting the performance of other VPNs served by the same device). The aggregate scaling demands of this model are also significant. Given a change in the underlying physical or link-layer topology, the consequent

requirement to process the routing update on a per-VPN basis becomes a significant challenge. Use of distance vector protocols to manage the routing tables would cause a corresponding sudden surge in traffic load, and the surge grows in direct proportion to the number of supported VPNs. The use of link-state routing protocols would require the consequent link-state calculation to be repeated for each VPN, causing the router to be limited by available CPU capacity.

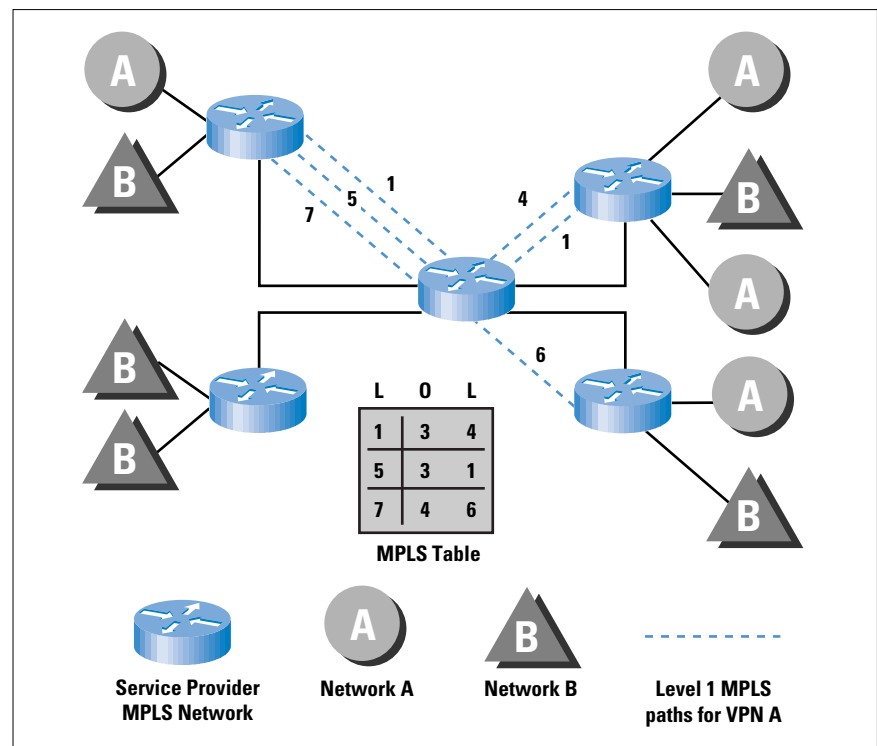
Multiprotocol Label Switching

One method of addressing these scaling issues is to use VPN labels within a single routing environment, in the same way that packet labels are necessary to activate the correct per-VPN routing table. The use of local label switching effectively recreates the architecture of a Multiprotocol Label Switching VPN. It is perhaps no surprise that when presented with two basic approaches to the architecture of the VPN—the use of network-layer routing structures and per-packet switching, and the use of link-layer circuits and per-flow switching—the industry would devise a hybrid architecture that attempts to combine aspects of these two approaches. This hybrid architecture is referred to as *Multiprotocol Label Switching* (MPLS)^[7, 8].

The architectural concepts used by MPLS are generic enough to allow it to operate as a peer VPN model for switching technology for a variety of link-layer technologies, and in heterogeneous Layer 2 transmission and switching environments. MPLS requires protocol-based routing functionality in the intermediate devices, and operates by making the interswitch transport infrastructure visible to the routing. In the case of IP over ATM, each ATM bearer link becomes visible as an IP link, and the ATM switches are augmented with IP routing functionality. IP routing is used to select a transit path across the network, and these transit paths are marked with a sequence of labels that can be thought of as locally defined forwarding path indicators. MPLS itself is performed using a label swapping forwarding structure. Packets entering the MPLS environment are assigned a local label and an outbound interface based on a local forwarding decision. The local label is attached to the packet via a lightweight encapsulation mechanism. At the next MPLS switch, the forwarding decision is based on the incoming label value, where the incoming label determines the next hop interface and next hop label, using a local forwarding table indexed by label. This lookup table is generated by a combination of the locally used IP routing protocol, together with a label distribution protocol, which creates end-to-end transit paths through the network for each IP destination. It is not our intention to discuss the MPLS architecture in detail, apart from noting that each MPLS switch uses a label-indexed forwarding table, where the attached label of an incoming packet determines the next-hop interface and the corresponding outgoing label.

The major observation here is that this lightweight encapsulation, together with the associated notion of boundary-determined transit paths, provides many of the necessary mechanisms for the support of VPN structures^[9]. MPLS VPNs have not one, but three key ingredients: (1) constrained distribution of routing information as a way to form VPNs and control inter-VPN connectivity; (2) the use of VPN-IDs, and specifically the concatenation of VPN-IDs with IP addresses to turn (potentially) nonunique addresses into unique ones; and (3) the use of label switching (MPLS) to provide forwarding along the routes constructed via (1) and (2). The generic architecture of deployment is that of a label-switched common host network and a collection of VPN environments that use label-defined virtual circuits on an edge-to-edge basis across the MPLS environment. An example is indicated in Figure 4, which shows how MPLS virtual circuits are constructed.

Figure 4:
MPLS "Tunnels,"
or VPNs



Numerous approaches are possible to support VPNs within an MPLS environment. In the base MPLS architecture, the label applied to a packet on ingress to the MPLS environment effectively determines the selection of the egress router, as the sequence of label switches defines an edge-to-edge virtual path. The extension to the MPLS local label hop-by-hop architecture is the notion of a per-VPN global identifier (or *Closed User Group* (CUG) identifier, as defined in [5]), which is used effectively within an edge-to-edge context. This global identifier could be assigned on ingress, and is then used as an index into a per-VPN routing table to determine the initial switch label. On egress from the MPLS environment, the CUG identifier would be used again as an index into a per-VPN global identifier table to undertake next-hop selection.

Routing protocols in such an environment need to carry the CUG identifier to trigger per-VPN routing contexts, and a number of suggestions are noted in [5] as to how this could be achieved.

It should be stressed that MPLS itself, as well as the direction of VPN support using MPLS environments, is still within the area of active research, development, and subsequent standardization within the IETF, so this approach to VPN support is still somewhat speculative in nature.

Link-Layer Encryption

As mentioned previously, encryption technologies are extremely effective in providing the segmentation and virtualization required for VPN connectivity, and can be deployed at almost any layer of the protocol stack. Because there are no intrinsically accepted industry standards for link-layer encryption, all link-layer encryption solutions are generally vendor specific and require special encryption hardware.

Although this scenario can avoid the complexities of having to deal with encryption schemes at higher layers of the protocol stack, it can be economically prohibitive, depending on the solution adopted. In vendor proprietary solutions, multivendor interoperability is certainly a genuine concern.

Transport and Application-Layer VPNs

Although VPNs can certainly be implemented at the transport and application layers of the protocol stack, this setup is not very common. The most prevalent method of providing virtualization at these layers is to use encryption services at either layer; for example, encrypted e-mail transactions, or perhaps authenticated *Domain Name System* (DNS) zone transfers between different administrative name servers, as described in DNSSec (*Domain Name System Security*)^[10].

Some interesting, and perhaps extremely significant, work is being done in the IETF to define a *Transport Layer Security* (TLS) protocol^[11], which would provide privacy and data integrity between two communicating applications. The TLS protocol, when finalized and deployed, would allow applications to communicate in a fashion that is designed to prevent eavesdropping, tampering, or message forgery. It is unknown at this time, however, how long it may be before this work is finalized, or if it will be embraced by the networking community as a whole after the protocol specification is completed.

The significance of a “standard” transport-layer security protocol, however, is that when implemented, it could provide a highly granular method for virtualizing communications in TCP/IP networks, thus making VPNs a pervasive commodity, and native to all desktop computing platforms.

Non-IP VPNs

Although this article has focused on TCP/IP and VPNs, it is recognized that multiprotocol networks may also have requirements for VPNs. Most of the same techniques previously discussed can also be applied to multiprotocol networks, with a few obvious exceptions—many of the techniques described herein are solely and specifically tailored for TCP/IP protocols.

Controlled route leaking is not suitable for a heterogeneous VPN protocol environment, in that it is necessary to support all protocols within the common host network. GRE tunnels, on the other hand, are constructed at the network layer in the TCP/IP protocol stack, but most routable multiprotocol traffic can be transported across GRE tunnels (for example, IPX and AppleTalk). Similarly, the VPDN architectures of L2TP and PPTP both provide a PPP end-to-end transport mechanism that can allow per-VPN protocols to be supported, with the caveat that it is a PPP-supported protocol in the first place.

The reverse of heterogeneous VPN protocol support is also a VPN requirement in some cases, where a single VPN is to be layered above a heterogeneous collection of host networks. The most pervasive method of constructing VPNs in multiprotocol networks is to rely upon application-layer encryption, and the resulting VPNs are generally vendor proprietary, although some would contend that one of the most pervasive examples of this approach was the mainstay of the emergent Internet in the 1970s and 1980s—that of the UNIX-to-UNIX Copy Program (UUCP) network, which was (and remains) an open technology.

Quality-of-Service Considerations

In addition to creating a segregated address environment to allow private communications, the expectation that the VPN environment will be in a position to support a set of service levels also exists. Such per-VPN service levels may be specified either in terms of a defined service level that the VPN can rely upon at all times, or in terms of a level of differentiation that the VPN can draw upon the common platform resource with some level of priority of resource allocation.

Using dedicated leased circuits, a private network can establish fixed resource levels available to it under all conditions. Using a shared switched infrastructure, such as Frame Relay virtual circuits or ATM virtual connections, a quantified service level can be provided to the VPN through the characteristics of the virtual circuits used to implement the VPN.

When the VPN is moved away from such a circuit-based switching environment to that of a general Internet platform, is it possible for the Internet Service Provider to offer the VPN a comparable service level that attempts to quantify (and possibly guarantee) the level of resources that the VPN can draw upon from the underlying host Internet?

This area is evolving rapidly, and much of it remains within the realm of speculation rather than a more concrete discussion about the relative merits of various Internet QoS mechanisms. Efforts within the *Integrated Services Working Group* of the IETF have resulted in a set of specifications for the support of guaranteed and controlled load end-to-end traffic profiles using a mechanism that loads per-flow state into the switching elements of the network^[12, 13]. There are numerous caveats regarding the use of these mechanisms, in particular relating to the ability to support the number of flows that will be encountered on the public Internet^[14]. Such caveats tend to suggest that these mechanisms will not be the ones that are ultimately adopted to support service levels for VPNs in very large networking environments.

If the scale of the public Internet environment does not readily support the imposition of per-flow state to support guarantees of service levels for VPN traffic flows, the alternative query is whether this environment could support a more relaxed specification of a differentiated service level for overlay VPN traffic. Here, the story appears to offer more potential, given that differentiated service support does not necessarily imply the requirement for per-flow state, so stateless service differentiation mechanisms can be deployed that offer greater levels of support for scaling the differentiated service^[15]. However, the precise nature of these differentiated service mechanisms, and their capability to be translated to specific service levels to support overlay VPN traffic flows, still remain in the area of future activity and research.

Conclusions

So what is a virtual private network? As we have discussed, a VPN can take several forms. A VPN can be between two end systems, or it can be between two or more networks. A VPN can be built using tunnels or encryption (at essentially any layer of the protocol stack), or both, or alternatively constructed using MPLS or one of the “virtual router” methods. A VPN can consist of networks connected to a service provider’s network by leased lines, Frame Relay, or ATM, or a VPN can consist of dialup subscribers connecting to centralized services or other dialup subscribers.

The pertinent conclusion here is that although a VPN can take many forms, a VPN is built to solve some basic common problems, which can be listed as virtualization of services and segregation of communications to a closed community of interest, while simultaneously exploiting the financial opportunity of economies of scale of the underlying common host communications system.

To borrow a popular networking axiom, “When all you have is a hammer, everything looks like a nail.” Every organization has its own problem that it must solve, and each of the tools mentioned in this article can be used to construct a certain type of VPN to address a particular set of functional objectives. More than a single “hammer” is

available to address these problems, and network engineers should be cognizant of the fact that VPNs are an area in which many people use the term generically—there is a broad problem set with equally as many possible solutions. Each solution has numerous strengths and also numerous weaknesses and vulnerabilities. No single mechanism for VPNs that will supplant all others in the months and years to come exists, but instead a diversity of technology choices in this area of VPN support will continue to emerge.

Acknowledgments

Thanks to Yakov Rekhter, Eric Rosen, and W. Mark Townsley, all of Cisco Systems, for their input and constructive criticism.

References

- [1] Valencia, A., M. Littlewood, and T. Kolar. “Layer Two Forwarding (Protocol) ‘L2F’.” **draft-valencia-l2f-00.txt**, work in progress, October 1997.
- [2] Droms, R. “Dynamic Host Configuration Protocol.” RFC 2131, March 1997.
- [3] Kent, S., and R. Atkinson. “Security Architecture for the Internet Protocol.” **draft-ietf-ipsec-arch-sec-04.txt**, work in progress, March 1998.
- [4] Additional information on IPSec can be found on the IETF IPSec home page, located at <http://www.ietf.org/html.charters/ipsec-charter.html>
- [5] Heinanen, J. “Multiprotocol Encapsulation over ATM Adaptation Layer 5.” RFC 1483, July 1993.
- [6] The ATM Forum. “Multi-Protocol Over ATM Specification v1.0.” **af-mpoa-0087.000**, July 1997.
- [7] Callon, R., P. Doolan, N. Feldman, A. Fredette, G. Swallow, and A. Viswanathan. “A Framework for Multiprotocol Label Switching.” **draft-ietf-mpls-framework-02.txt**, work in progress, November 1997.
- [8] Rosen, E., A. Viswanathan, and R. Callon. “A Proposed Architecture for MPLS.” **draft-ietf-mpls-arch-01.txt**, work in progress, March 1998.
- [9] Heinanen, J. and E. Rosen. “VPN Support for MPLS.” **draft-heinanen-mpls-vpn-01.txt**, work in progress, March 1998.
- [10] Eastlake, D. and C. Kaufman. “Domain Name System Security Extensions.” RFC 2065, January 1997. For further information regarding DNSSec, see: <http://www.ietf.org/html.charters/dnssec-charter.html>

- [11] Dierks, T. and C. Allen. “The TLS Protocol—Version 1.0.” **draft-ietf-tls-protocol-05.txt**, work in progress, November 1997. For more information on the IETF TLS working group, see <http://www.ietf.org/html.charters/tls-charter.html>. See also the article on SSL in the *Internet Protocol Journal*, Volume 1, No. 1, June 1998.
- [12] Wroclawski, J. “Specification of the Controlled-Load Network Element Service.” RFC 2211, September 1997.
- [13] Shenker, S., C. Partridge, and R. Guerin. “Specification of Guaranteed Quality of Service.” RFC 2212, September 1997.
- [14] Mankin, A., F. Baker, S. Bradner, M. O’Dell, A. Romanow, A. Weinrib, and L. Zhang. “Resource ReSerVation Protocol (RSVP) Version 1—Applicability Statement, Some Guidelines on Deployment.” RFC 2208, September 1997.
- [15] “Differentiated Services Operational Model and Definitions.” **draft-nichols-dsopdef-00.txt**, work in progress, K. Nichols and S. Blake (editors), February 1998.

PAUL FERGUSON is a consulting engineer at Cisco Systems and an active participant in the Internet Engineering Task Force (IETF). His principal areas of expertise include large-scale network architecture and design, global routing, Quality of Service (QoS) issues, and Internet Service Providers. Prior to his current position at Cisco Systems, he worked in network engineering, analytical, and consulting capacities for Sprint, Computer Sciences Corporation (CSC), and NASA. He is coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, published by John Wiley & Sons, ISBN 0-471-24358-2, a collaboration with Geoff Huston. E-mail: ferguson@cisco.com

GEOFF HUSTON holds a B.Sc and a M.Sc from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the the initial build of the Internet within the Australian academic and research sector. Huston is currently the Chief Technologist in the Internet area for Telstra. He is also an active member of the IETF, and was an inaugural member of the Internet Society Board of Trustees. He is coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, published by John Wiley & Sons, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. E-mail: gih@telstra.net

Reliable Multicast Protocols and Applications

by C. Kenneth Miller, StarBurst Communications

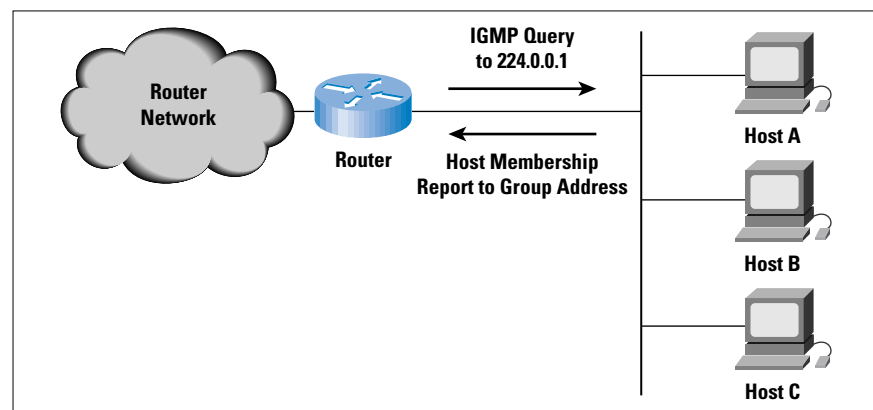
Multicast IP network services offer new opportunities to provide value-added applications that involve many-to-many transmission such as conferencing or network gaming, or one-to-many transmission such as multimedia events, tickertape feeds, and file transfer, where the many could be thousands or even conceivably millions. Multicast IP services use a different kind of IP address, called Class D. In contrast to individual host addresses (Classes A–C), which include a host and a network component and usually are semipermanent, Class D multicast addresses may by design be used only for a particular session, or can be semipermanent, as multicast groups may be set up and torn down relatively quickly, on the order of seconds. The IP address structure is shown in Figure 1.

Figure 1:
IP Address Types

	0	1	2	3	4	8	16	24	31												
Class A	0	netid				hostid															
Class B	1	0	netid								hostid										
Class C	1	1	0	netid											hostid						
Class D	1	1	1	0	multicast address																

Hosts join groups at the receiver's initiation using the *Internet Group Management Protocol* (IGMP). When a host joins a group, it notifies the nearest multicast subnet router of its presence in the group, as shown in Figure 2. First defined in RFC 1112^[1], IGMPv1 is still the version of IGMP most widely supported. IGMPv2 has recently been documented as an official RFC (RFC 2236^[2]). The main feature that IGMPv2 brings is reduced latency for leaving groups. In IGMPv1, the designated multicast router for the subnet polls for multicast group members; no response between polls indicates that all hosts in a particular multicast group have left the group, and that the routers can prune back the multicast routing tree.

Figure 2:
IGMPv1 Dialog

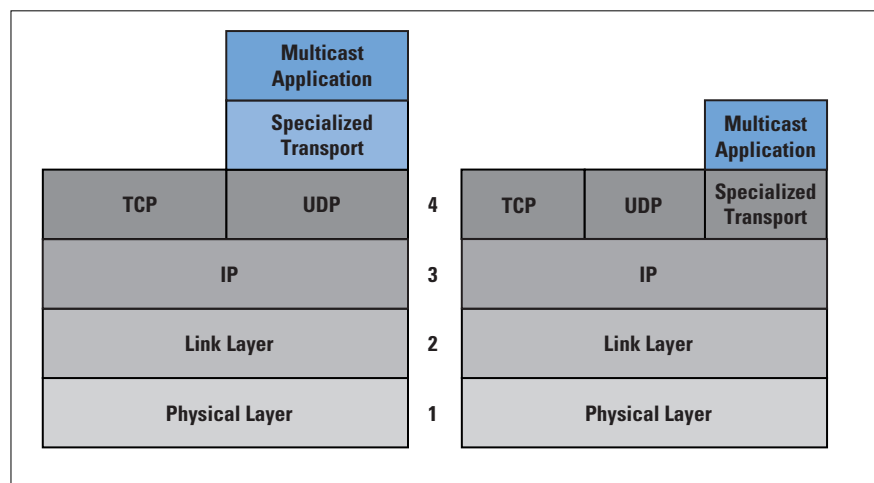


Network infrastructure devices, for example, routers, need to provide a routing protocol to forward multicast packets to group members, in a fashion similar to that performed for unicast routing. Multicast IP packet forwarding is best effort, just as it is with unicast packet forwarding. However, most unicast applications use TCP as a transport layer to provide guaranteed packet ordering and delivery. Some examples of applications that use TCP are the *File Transfer Protocol* (FTP) for file transfer and the *Hypertext Transfer Protocol* (HTTP) for Web access.

However, TCP is a unicast (point-to-point) only transport protocol. Thus, all multicast applications must run on top of the *User Datagram Protocol* (UDP) or alternatively, interface directly to IP via “raw” sockets and provide their own customized transport layer, as shown in Figure 3. UDP provides only minimal transport-layer services, error detection, and port multiplexing. Thus, if any errors or packet loss due to congestion occur, packets are simply lost to the application, and they are not recoverable. Thus, *all multicast applications must have a specific transport-layer service to support that particular application*. When that transport layer operates over UDP, it operates in the application layer with the application. When it interfaces directly to IP using “raw” sockets, the specialized transport layer operates at the transport layer, but is specialized to the particular application that uses it.

It should be noted that TCP supports only data reliability; it is not suited for transport of multimedia streams, which require consistent time delivery at the receiver and only need to be semireliable. Thus, multimedia streaming applications need a specialized transport layer such as the *Real-Time Transport Protocol* (RTP)^[3] for unicast as well as multicast transmissions.

Figure 3:
Specialized Multicast
Transport Protocols
Operate over UDP
or IP



Many equate multicast with multimedia, thinking that the Internet and private intranets will become an alternative entertainment media to television by using multicast IP network services and multimedia streaming technology. However, numerous other multicast applications require reliability rather than timeliness; they are multicast applications that are similar to those unicast applications that operate over TCP, except that delivery is to many recipients rather than just one.

Reliable Multicast Application Categories and Requirements

Reliable multicast applications come in three basic categories with differing requirements, as shown in Figure 4.

Figure 4:
Reliable Multicast
Application
Categories

Application Type	Latency Req.	Reliability	Scalability
Collaborative	Low	Semi/Strict	<100
Message Str.	Low/Medium	Semi/Strict	to Millions
Bulk Data	Not Real Time	Strict	to Millions

Collaborative applications such as data conferences (whiteboarding) and network-based games are many-to-many applications with modest scaling requirements of less than 100 participants. This kind of application requires low latency of less than 400 msec so that responses do not cause discomfort to the human participants. Transmission does not always need strict reliability; for example, refresh of background information for a network game could wait for the next refresh.

Message streaming applications such as tickertape and news feeds also often require low latency. Tickertape feeds to brokerage houses need to be very timely because the information loses value greatly with time. Time is very much money in this application, and there is also a need for strict reliability.

Tickertape feeds to consumers are purposely delayed by minutes because they are usually transmitted without charge, but they cannot be so stale as to be viewed as “old” information. This data does not have a strict reliability requirement because the next trade of a particular security refreshes the data. News feeds likewise have only a moderate latency requirement. If the news feeds are sent in a carousel fashion, that is, each news story is repeated, strict reliability may not be needed because it is refreshed in the next transmission of the same story.

Bulk data delivery has no specific latency requirement. Often there is a desire to schedule delivery during the night, when there is less network traffic. At other times, the desire is to receive the data almost

immediately. However, at all times the entire “file” or piece of data needs to be received to be complete. Strict reliability is the rule; for example, if any bit of a software image is lost, the data is worthless.

Message streaming and bulk data application scaling requirements span the gamut from tens to possibly even millions.

Reliable multicast transport protocols, in contrast to multimedia streaming transport protocols, have not yet been standardized. However, numerous reliable multicast protocols exist; some have been used only for research, while others have been commercialized.

The *Reliable Multicast Research Group* (RMRG) in the *Internet Research Task Force* (IRTF) is now studying reliable multicast. It is chartered to recommend techniques for a working group in the *Internet Engineering Task Force* (IETF) to create a set of reliable multicast standards.

Standardization Effort

The standardization effort has been started in an IRTF research group to study the problems and possible solutions by Internet researchers. This effort was first placed in the hands of researchers because the problems were considered very difficult to solve in the global Internet. Some of the concerns about reliable multicast were discussed in an expired Internet Draft published in November 1996 by the Transport Area Directors of IETF.

These concerns formed the basis for the work of the RMRG, which was formed in early 1997. The concerns from that document follow:

“A particular concern for the IETF (and a dominant concern for the Transport Services Area) is the impact of reliable multicast traffic on other traffic in the Internet in times of congestion (more specifically, the effect of reliable multicast traffic on competing TCP traffic). The success of the Internet relies on the fact that best-effort traffic responds to congestion on a link (as currently indicated by packet drops) by reducing the load presented on that link. Congestion collapse in today’s Internet is prevented only by the congestion control mechanism in TCP.

There are a number of reasons to be particularly attentive to the congestion-related issues raised by reliable multicast proposals. Multicast applications in general have the potential to do more congestion-related damage to the Internet than do unicast applications. This is because a single multicast flow can be distributed along a large, global multicast tree reaching throughout the entire Internet.

Further, reliable multicast applications have the potential to do more congestion-related damage than do unreliable multicast applications. First, unreliable multicast applications such as audio and video are, at the moment, usually accompanied by a person at the receiving end, and people typically unsubscribe from a multicast group if congestion is so heavy that the audio or video stream is unintelligible. Reliable multicast applications such as group file transfer applications, on the other hand, are likely to be between computers, with no humans in attendance monitoring congestion levels.

In addition, reliable multicast applications do not necessarily have the natural time limitations typical of current unreliable multicast applications. For a file transfer application, for example, the data transfer might continue until all of the data is transferred to all of the intended receivers, resulting in a potentially-unlimited duration for an individual flow. Reliable multicast applications also have to contend with a potential explosion of control traffic (e.g., ACKs, NAKs, status messages), and with control traffic issues in general that may be more complex than for unreliable multicast traffic.

The design of congestion control mechanisms for reliable multicast for large multicast groups is currently an area of active research. The challenge to the IETF is to encourage research and implementations of reliable multicast, and to enable the needs of applications for reliable multicast to be met as expeditiously as possible, while at the same time protecting the Internet from the congestion disaster or collapse that could result from the widespread use of applications with inappropriate reliable multicast mechanisms. Because of the setbacks and costs that could result from the widespread deployment of reliable multicast with inadequate congestion control, the IETF must exercise care in the standardization of a reliable multicast protocol that might see widespread use.”

One of the statements in this document is very specious:

“First, unreliable multicast applications such as audio and video are, at the moment, usually accompanied by a person at the receiving end, and people typically unsubscribe from a multicast group if congestion is so heavy that the audio or video stream is unintelligible. Reliable multicast applications such as group file transfer applications, on the other hand, are likely to be between computers, with no humans in attendance monitoring congestion levels.”

This statement is a very weak argument; it is not reliable to depend on a human to turn off a nonfunctioning event. Do we typically turn off the television when we leave the house? Or leave the room to do something else?

In contrast, some of the reliable multicast protocols such as the *Multicast File Transfer Protocol* (MFTP) have the sense of a finite session, and automatically time out and leave a group, even if all group members did not receive all the content.

Essentially what is desired is a reliable multicast protocol that behaves like TCP in that it backs off in the face of congestion approximately the same way as TCP and shares the bandwidth with TCP traffic “fairly.” This feature is of prime importance to Internet researchers who wish to specify protocols that can scale to the global Internet and not cause harm to the traffic already present.

Two additional significant problems need to be solved: scalability and the ability to operate with scalability over many different network infrastructures.

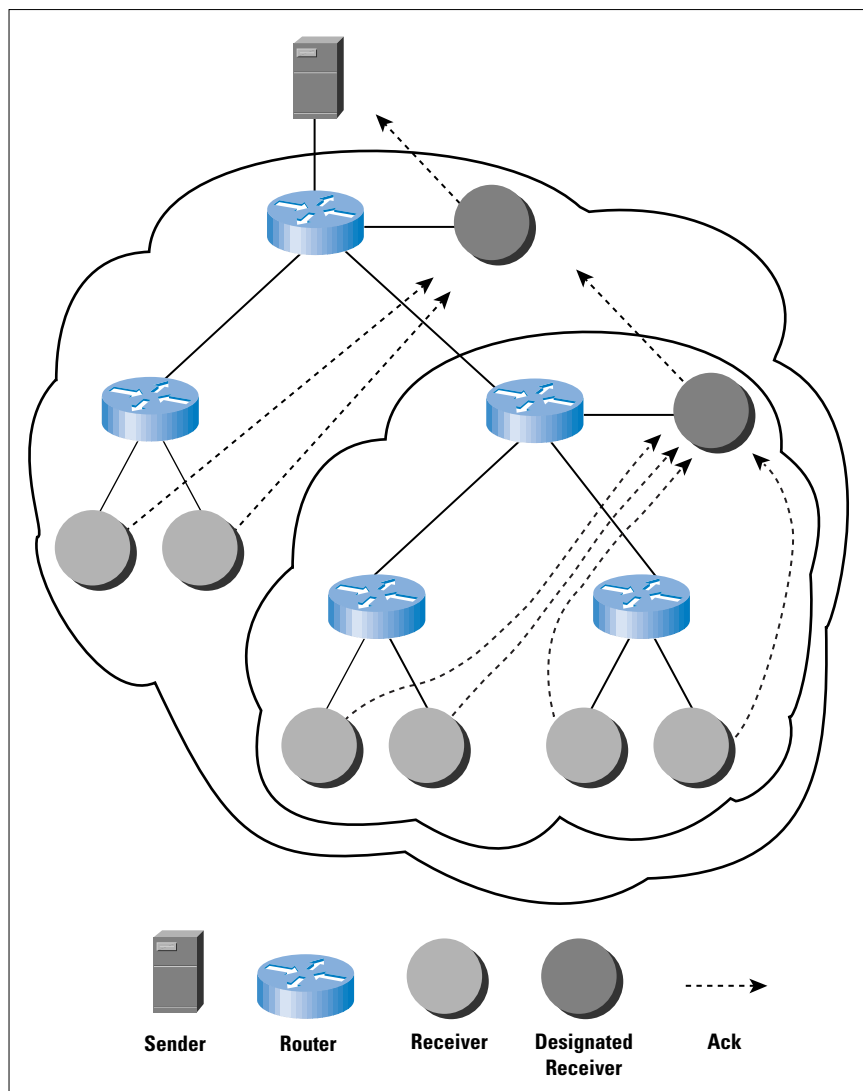
Scaling Issues and How Current Reliable Multicast Protocols Solve Them

Two primary issues are related to scaling, that is, the ability to handle large groups. The first and most significant is widely known as acknowledgment/negative acknowledgment (ACK/NAK) implosion. As the number of receivers grows, the amount of back traffic to the sender eventually overwhelms its capacity to handle them. Additionally, the network at the sender site becomes congested from the cumulative back traffic from the receivers.

The second issue is one of retransmissions (often referred to as “repairs”). If the packet loss is uncorrelated at the receivers, retransmissions grow, so the data may need to be sent multiple times to satisfy all the receivers. Measurements of the *Multicast backbone* (Mbone) have shown that loss consists of both correlated and uncorrelated parts^[4]. Satellite networks will also exhibit mostly uncorrelated loss, unless receivers are geographically close.

Various methods have been used to achieve scaling by reducing the amount of ACK/NAK administrative traffic while still retaining reliability. A straightforward approach is to simply deploy repeaters/aggregators in the network, as shown in Figure 5. This approach is provided by the *Reliable Multicast Transport Protocol* (RMTP)^[5]. RMTP provides for *designated receivers* (DRs) that collect status messages from nodes in a local RMTP domain and provide repairs (retransmissions of missing data), if available. Receivers direct the administrative messages to the DR by unicast. Thus, the DR provides both local recovery and consolidation of control traffic to the next DR in the hierarchy if the data requested is not available.

Figure 5:
RMTP Designated
Receivers



A second approach is to allow any receiver to provide the repair, biasing the request to the nearest receiver that has the requested data. This approach, called *Scalable Reliable Multicast (SRM)*^[6], depends on the concept of repair by any receiver that has the data to gain scalability in reducing administrative back traffic to the source, putting the onus of responsibility on receivers to ensure that they get missed data.

Group members in SRM send low-frequency *session* messages to the group so that their neighbors can learn their status, measure the delay among group members and learn group membership, and detect the last packet in a burst. Session messages are designed to take only about five percent of the traffic in the session.

Receivers with missing data wait a random time period before issuing repair requests, allowing suppression of duplicate requests similar to the mechanism that IGMP uses on its subnet. A similar process occurs for making the actual repairs. The random backoff time for both repair requests made by receivers and repairs made by senders is a

function of “closeness” to the sender and requesting receiver. Thus, those closest to each other time out first and make the repair request or the actual repair in an attempt to keep repairs as local as possible. A receiver that sees the first request and determines that it is the same request that it would have made simply stays silent, reducing potential redundant requests. The requester continues to send repair requests until the repair is received.

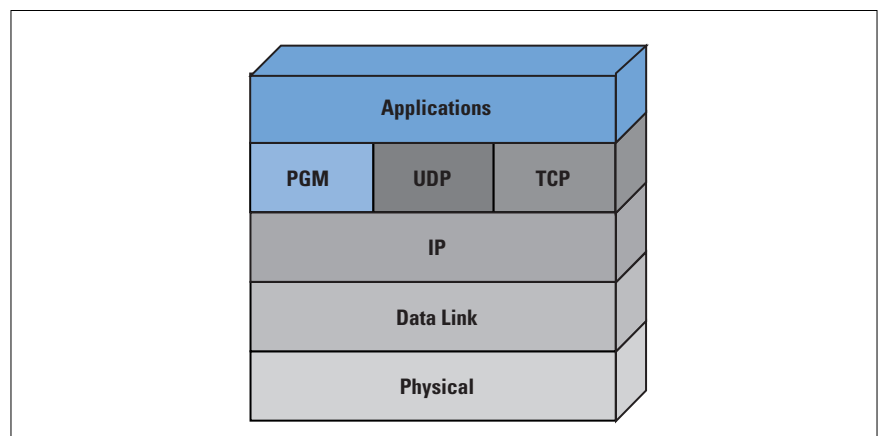
Any receiver may satisfy the repair request, because all receivers are required to cache previously sent data. Any receiver that can satisfy the request is prepared to do so; a random backoff timer is used before a repair is sent, and if it sees the repair being sent by another group member, it stays silent to reduce the probability of sending duplicate repairs.

SRM was first developed to be the reliable multicast protocol to operate with the *wb* whiteboard data conferencing tool developed by Lawrence Berkeley Labs (LBL) researchers, SRM is currently operational over the Mbone, the experimental multicast network of the Internet.

A third approach is to have the network infrastructure, that is, routers, help in providing scaling. This approach, called *Pretty Good Multicast* (PGM)^[7], is a new proposal that was first publicly presented to the RMRG meeting held in February 1998.

One design goal of the creators of PGM was simplicity and the ability to optimally leverage routers in the network to provide scalability. PGM is an example of a protocol that bypasses UDP and interfaces directly to IP via “raw” sockets, as shown in Figure 6.

Figure 6:
PGM Interfaces
Directly to IP



PGM provides no notion of group membership; it simply provides reliability within a source’s transmit window from the time a receiver joins a group until it departs.

PGM has only a few data packets that are defined:

ODATA: original content data

NAK: selective negative acknowledgment

NCF: NAK confirmation

RDATA: retransmission (repair)

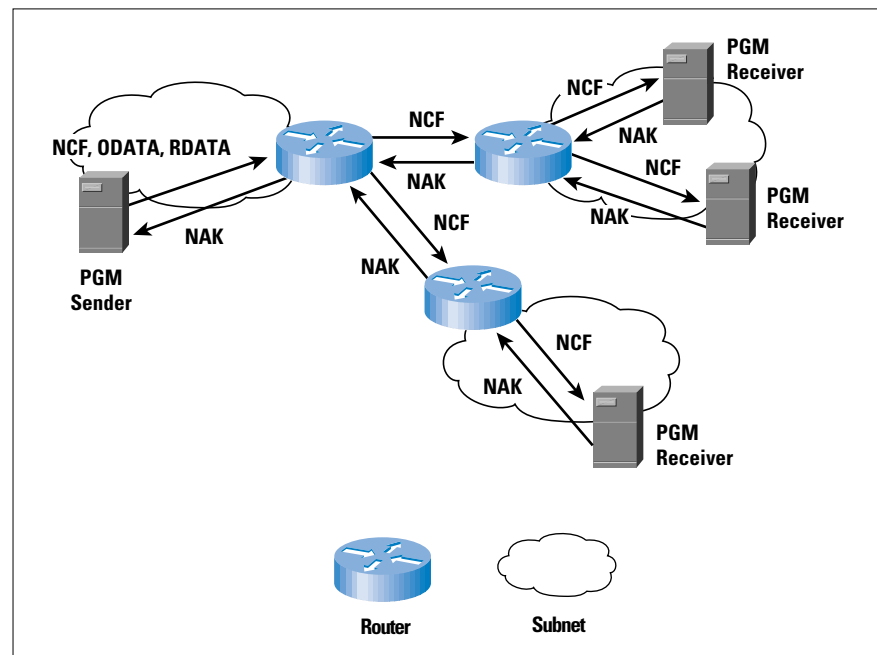
SPM: source path message

Each PGM packet contains a *Transport Session Identifier* (TSI) to identify the session and source of that data, so multiple sessions may be easily identified by PGM-aware routers and receivers. *ODATA*, *NCF*, *RDATA*, and *SPM* packets flow downstream in the distribution tree, and *NAK* packets flow upstream toward the source.

PGM is designed for scalability as well as the ability to serve real-time applications. Thus there is a need for timeliness. This need is handled by the *transmit window*, which defines a sliding window of data such that if no *NAK*s are received by the sender or a designated local retransmitter by the time the window is up, the data is simply not available for repairs.

PGM is totally *NAK* based, so the scaling issue is to reduce the number of *NAK*s sent back to the source, while at the same time protecting against lost *NAK*s. Enter here the router assist, as shown in Figure 7.

Figure 7:
PGM NAK/NCF
Dialog



NAKs are unicast from PGM-router to PGM-router, initiated by the receiver that lost data sending a NAK to its nearest PGM-aware router. Each PGM-aware router keeps forwarding NAKs until it sees an NCF or RDATA, which indicates that a repair is being sent. NAK *suppres-*

sion is provided by a receiver's subnet PGM-aware router, and all PGM-aware routers *eliminate* duplicate NAKs all the way upstream to the source.

The unicast path back to the source must be the same path as the downstream multicast tree. SPMs are sent downstream interleaved with ODATA packets to establish a source path state for a given source and session. PGM-aware routers use this information to determine the unicast path back to the source for forwarding NAKs. SPMs also alert receivers that the oldest data in the transmit window is about to be retired from the window and will thus no longer be available for repairs from the source. SPMs are sent by a source at a rate that is at least the rate at which the transmit window is advanced. This rate provokes “last call” NAKs from receivers and updates the receive window state at receivers.

PGM-aware routers also keep state on where the NAKs come from in the distribution tree so that they may constrain the forwarding of RDATA repairs to only those ports from which NAKs requesting that repair were received. This scenario eliminates the transmission of repair data to parts of the distribution tree where the repair is not needed.

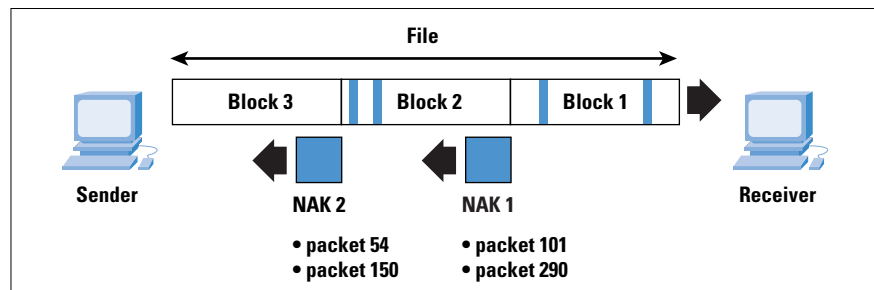
The PGM feature can also optionally redirect NAKs to a *designated local retransmitter* (DLR) rather than the source. A DLR announces its presence to provoke the redirection of NAKs for that session and source.

A fourth approach is to not have a low-latency requirement (that is, only serve “bulk data” delivery applications) and use this feature to advantage to gain scalability. MFTP was first published as an Internet Draft in February 1997, and an update was submitted in April 1998^[8].

MFTP also has a provision for sender-based group creation, with different group models, and the group setup protocol to notify receivers to join the group. Group creation is discussed later in this article.

The basic MFTP protocol breaks the data entity to be sent into maximum size “blocks,” where a block by default consists of thousands or tens of thousands of packets, depending on packet size used. This setup is shown in Figure 8.

Figure 8:
MFTP Blocks



MFTP is a “NAK-only” protocol; that is, if data is received correctly in a block, nothing is sent back to the sender. If one or more packets are in error or missing in a block, receivers respond with a NAK that consists of a bit map of the bad packets in the block. It is thus a *selective reject* mechanism. In this respect, MFTP is similar to RMTP; the main difference is that MFTP explicitly attempts to make the block as large as possible for scaling purposes.

NAKs are normally sent unicast back to the source, unless aggregation to improve scaling using enabled network routers is used. In this case, the NAKs are sent multicast to a special administrative traffic group address.

MFTP does not repair after each block, however; it takes advantage of the non-real time nature of the application for benefit. The data entity, such as a file, is sent initially in its entirety in a *first pass*. The sender collects the NAK packets for a block from all the receivers. One NAK packet from a receiver can represent thousands or even tens of thousands of bad packets, reducing NAK implosion by orders of magnitudes. The collection of NAKs received by the sender from all the receivers is logically OR-ed together to represent the collective need for repairs for the receiving group. These repairs are sent by the sender in a *second pass* to the group. If certain receivers already have the repair, it is simply ignored. This scenario is repeated, if necessary, until all repairs are received by all receivers or until a configurable timeout occurs.

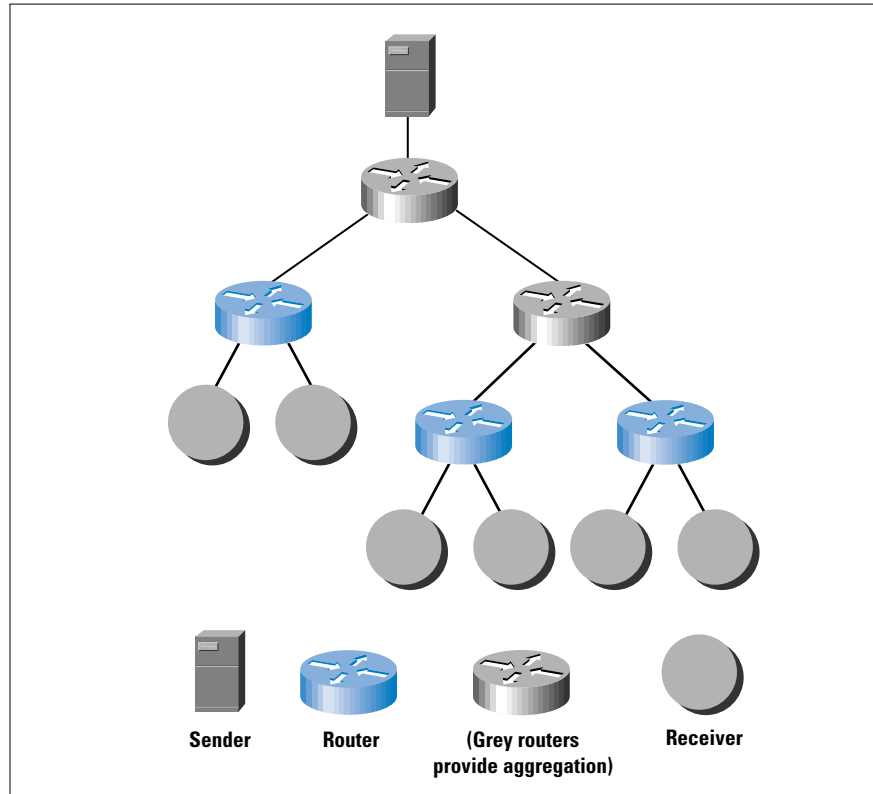
Thus, packet ordering services are not provided, and holes in the data caused by dropped packets or packets in error are filled in as they are received.

The sender is *rate based*; in other words, it transmits at a data rate set by the operator to be less than or equal to what the network can handle. The protocol is thus very efficient with high-latency networks such as satellites, and it is impervious to network asymmetry. It also attempts to be as scalable as possible on one-hop networks such as satellite networks, and it provides for extensions so that network elements may aggregate downstream responses to increase scalability further, depending on the network configuration.

This aggregation capability is shown in Figure 9. The network element, which can be a router, collects MFTP administrative back traffic routers are members. These routers aggregate back traffic from all nodes downstream in the multicast tree from the source, including registrations, NAKs, and dones. Registration and done messages are used by MFTP’s group setup protocol, and they are described later in this article.

Depending on the network configuration, this aggregation capability can further improve the scalability of MFTP by orders of magnitudes.

Figure 9:
Routers as Network
Aggregators



The upper limit to scalability with no network aggregation of administrative traffic is in the tens of thousands of receivers. For example, for a *Maximum Transmission Unit* (MTU) of 1500 bytes (the Ethernet maximum), the default block size is over 11,000 packets. If the number of receivers is 10,000 and each receiver has at least one bad packet per block, then there will be a total of 10,000 NAK packets coming back to the sender from the group about that block, approximately the same number of packets as were sent in the forward direction in that block. MFTP provides for a NAK backoff timer to spread the NAKs out in time to the sender to avoid bursts. If the bandwidth is symmetric at the sender, the sender should be able to handle this maximum NAK. In many situations, the amount of back traffic could exceed forward traffic.

MFTP also has provision for a crude congestion control mechanism. The sender at the beginning of a session sends *announce* messages. These messages are used for many functions, including the setting up of groups. Additionally, it conveys a packet loss parameter to all receivers. This packet loss threshold parameter may be used by receivers to leave the group if the packet loss exceeds the threshold. Leaving the group prunes the distribution tree, relieving the congestion in that section of the tree.

Commercial Usage

The reliable multicast protocols previously discussed are the most prominent ones on the market today. RMTP has been deployed in its message streaming version for a billing record distribution application within a very large telecommunications carrier, but it has had generally limited deployment. It also does not scale over satellite networks, where most of the early multicast deployments reside.

SRM has been used by the research community only over the Mbone, and it is still being refined. Another problem with SRM is that in its current incarnation, it supports neither asymmetric nor satellite networks. Some early Internet Service Provider (ISP) multicast implementations, offer multicast support in only one direction; SRM requires total multicast support.

PGM is new and offers promise, but there is no deployment yet, and it likely will not occur until early 1999. PGM also requires router support in a terrestrial land-line network to gain scaling.

MFTP has the limitation that it supports only bulk transfer applications. However, one trade-off is that it can support all network infrastructures, including satellite infrastructures with scaling. MFTP has also been available commercially in products with the longest application support, dating back to 1995. Thus, MFTP-based products have the largest installed base of any reliable multicast-based product being used over WANs. The largest commercial installation of over 8,500 remote sites in the group is the General Motors^[9] dealer network. Several other commercial installations of MFTP-based applications number over 1,000 group members.

Advanced Research Topics Discussed in Reliable Multicast Research Group

A promising technique to reduce the amount of repair data that needs to be retransmitted is called *erasure correction*. This technique can significantly reduce the amount of repairs that need to be resent if the packet loss is largely uncorrelated at the receivers. It uses a *forward error correction* (FEC) code to generate parity packets to be used for repairs only. This setup provides benefit if errors at receivers are uncorrelated. For example, suppose 16 receivers each have one missing packet, but they are all different. Rather than send all 16 original data packets, one FEC packet could be sent that could correct the one missing packet at all 16 receivers, requiring retransmission of only one packet rather than 16.

If the loss is correlated, then many of the receivers lose the same data, and erasure correction is of no benefit. However, there is also no penalty, except for the need for computing power at both the sender and the receivers to perform the FEC correction calculations. Simulations have shown^[10] that there is a greater than 2:1 reduction in the number of repairs needed to be sent with our example of 10,000 receivers. This benefit will be even larger when group sizes become larger than tens of thousands.

Perhaps a more significant application for FEC is a congestion control technique known as *layering*^[11,12]. With layering, numerous groups are set up by the sender, all with different rates. Receivers that can receive at the highest rate join all the “layer” groups. Those receivers that cannot receive at the highest rate simply leave “layers” until congestion is relieved, and they take longer to receive the data. For this to work without sending data redundantly, the number of parity packets created must be very large compared to the number of data packets.

There are some further issues that have been pointed out by the researchers with the Other issues with the layering approaches have been pointed out by the researchers, however. For layering to be effective, the routing tree should be identical for the different groups; otherwise congestion will not be relieved on a part of the tree. This may not always be the case, especially in sparse mode routing protocols, where selection of the rendezvous point or core is based on group address.

Even if the same distribution tree is used for the different layers, it has been pointed out^[12] that leaves of hosts downstream from a congested link should be coordinated; otherwise the action of less than all of them has no effect on congestion. Additionally, a receiver could cause congestion by adding a layer that another receiver could interpret as congestion, causing it to drop a layer with no effect.

Thus, layering using FEC techniques is an interesting technique that shows promise for use in congestion control. However, there are issues associated with this type of layering that researchers still need to address.

Another technique that has been proposed for congestion control is bulk feedback to the sender^[13]. If the sender receives an excessive number of NAKs from receivers, it drops the sender’s transmission rate with an algorithm that attempts to emulate the behavior of TCP. This approach is an obvious one because it is an extension of the process in which TCP falls back in the face of congestion.

This approach, however, has two basic problems. The first is that there is delay, because the sender needs to get feedback from the multitude of receivers before it acts. This delay can be considerably longer than in the case of TCP, which needs feedback from only one receiver.

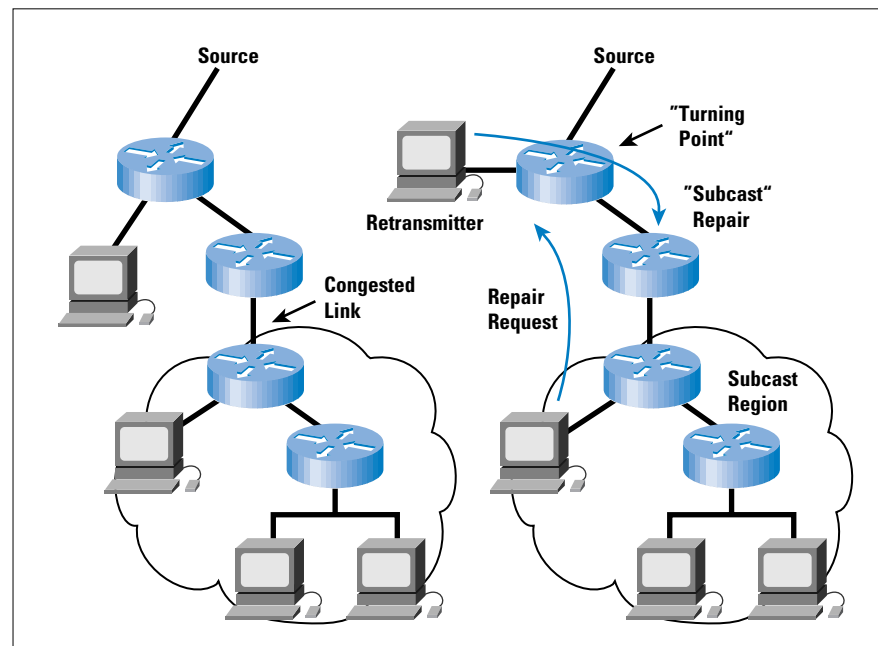
The second flaw is that one errant receiver can effectively penalize the whole group, because the sender reduces the rate to the total group.

This approach is not viewed as a viable solution for these reasons. In fact, the general consensus is that congestion control decision making will be required at the multiple receivers rather than at the sender for both scaling and timeliness reasons.

Another idea that is now receiving intense study by researchers is that of “subcasting”^[14,15,16]. The key idea in subcasting is to optimize local repair to be a retransmitter that may be just above a link congestion point, as shown in Figure 10. The problem is to gain knowledge of the network topology so as to locate a receiving host that is willing to retransmit and that has the repair data.

Then the repairs need to be contained within only the region of the network that lost the original transmission, that is, the “subcast” region.

Figure 10:
Optimized Local
Repair



One proposal is to ask for assistance from the network routers. They know the topology and could be used to find the closest willing retransmitter that has the repair. The router could also direct the repair to only the affected region: the *subcast*.

This technique can be viewed as an extension of concepts originally proposed in SRM to provide local recovery. It assumes that most loss is caused by congested links, and that uncorrelated loss is caused by a series of mildly congested links with few group members. This model is probably the right one for many land-line routed networks; it is problematical with other network infrastructures.

Nevertheless, it is an interesting proposal that merits further research effort. Local repair is destined to be an important tool to meet the goal of improved scalability with minimal traffic overhead.

Group Creation and Destruction

The process of joining a group and leaving a group in IP multicast is left to a potential group member that uses IGMP to notify the nearest multicast router of its membership state. However, mechanisms need to be in place to allow potential members of a group to gain the information needed to decide to join the group.

There are two basic ways to accomplish this scenario for one-to-many sessions. The first and most common is the “broadcast TV” model. The *Multiparty Multimedia Session Control* (MMUSIC) working group of the IETF has developed some protocols that can be used to advertise content. The *Session Announcement Protocol* (SAP)^[17] provides the mechanism to send a stream on a “well-known” multicast address to announce content to any potential listeners who may be interested. It uses the *Session Description Protocol* (SDP)^[18] to describe the contents that are announced. These two protocols together have been used to create a session directory tool that is available on the Mbone. This setup creates essentially the equivalent of a “preview channel” such as is often available on cable television systems.

SDP is also used to post content on Web sites, which advertise that content to anyone who wishes to receive it.

Although these protocols were originally developed primarily to advertise multimedia streaming applications, they are also applicable for data. They provide a useful tool for “push” vendors to advertise multicast “channels” based on content that any consumer can “tune in” to.

Internet researchers describe this model as providing “loosely coupled” sessions, because the sender does not know who is listening, much like radio or TV broadcasters do not know who tunes in to their stations.

MFTP also includes a group setup protocol. The “closed group” option in MFTP provides a mechanism to create a “tightly coupled” session that is very useful to organizations that wish to deliver critical information from a central site to many remote branch offices. The closed group provides a means for the sender to define a group list centrally and direct those members so defined to join the group. This scenario is somewhat similar to e-mail, except more robust.

These instructions are sent in an “announce” message on a special multicast group address that the superset of possible candidate receivers always listens to. Hosts so directed to join the group notify their designated multicast router of their membership directed to join the group notify their designated multicast router of their membership using IGMP and “register” back to the sender of their presence. Thus, the sender knows group membership before transmission commences, and the sender can then also positively confirm delivery.

This approach has proven very desirable for organizations that have many branches where information is desired to be sent at the discretion and time determined by the sender, and usually the information is delivered to a branch office server. Several deployments of applications that use MFTP and the closed group model with group members approaching 10,000 exist.

The MMUSIC group has also created the *Session Invitation Protocol* (SIP)^[19], which is used to invite members to a conference of some sort, including possibly a data conference. This protocol is appropriate for use with whiteboard applications, for example.

Summary and Conclusions

Although multicast has often been viewed as synonymous with multimedia, there is a wide spectrum of reliable multicast applications that involve the transfer of data to multiple group members. Because this wide spectrum of applications has many different requirements, as shown in Figure 4, no one reliable multicast protocol can handle all applications and network infrastructures. The result is that numerous reliable multicast protocols are likely to become standardized, and today numerous reliable multicast protocols are either in commercial products/toolkits or due to be available soon.

The reliable multicast standardization effort now resides in the IRTF, because Internet researchers are concerned about congestion control and fairness to TCP for any protocols that might become standardized for general Internet use. This problem is difficult to solve, given the disparate requirements placed on protocols by the wide variety of applications and different network infrastructures.

Nevertheless, a significant number of reliable multicast-based product deployments have already occurred over private networks. These have been shown to save organizations much money and to help create new business opportunities for them.

Stay tuned; reliable multicast-based applications are ready to be mainstreamed. Together with multimedia multicast applications, multicast applications of all forms will become common soon, first in private intranets and extranets and then in the Internet as a whole.

References

- [1] Deering, S. "Host Extensions for IP Multicasting." RFC 1112, August 1989.
- [2] Fenner, W. "Internet Group Management Protocol, Version 2." RFC 2236, November 1997.
- [3] Schulzrinne, H., Casner, S., Frederick, R., and Jacobson, V. "RTP: A Transport Protocol for Real-Time Applications." RFC 1889, January 1996.
- [4] Handley, M. "An Examination of Mbone Performance." ISI Report, January 10, 1997.
- [5] Paul, S., Sabnani, K. K., Lin, J. C., and Bhattacharyya, S. "Reliable Multicast Transport Protocol (RMTP)." *IEEE Journal on Selected Areas in Communications*, April 1997.
- [6] Floyd, S., Jacobson, V., Liu, C., McCanne, S., and Zhang, L. "A Reliable Multicast Framework for Light-weight Sessions and Application Level Framing." *ACM Transactions on Networking*, November 1996.
- [7] Farinacci, D., Lin, A., Speakman, T., and Tweedly, A. "PGM Reliable Transport Protocol Specification." Work in progress, Internet Draft, **draft-speakman-pgm-spec-01.txt**, January 29, 1998.
- [8] Miller, K., Robertson, K., Tweedly, A., and White, M. "StarBurst Multicast File Transfer Protocol (MFTP) Specification." Work in progress, Internet Draft, **draft-miller-mftp-spec-03.txt**, April 1998.
- [9] Miller, K. "Reliable Multicast Protocols: A Practical View." 22nd Conference on Local Computer Networks, November 1997.
- [10] Kasera, S. K., Kurose, J., Towsley, D., "Scalable Reliable Multicast Using Multiple Multicast Groups." CMPSCI Technical Report TR 96-73, October 1996.
- [11] Nonnenmacher, J., and Biersack, E. W. "Asynchronous Multicast Push: AMP." Proceedings of ICC '97 International Conference on Computer Communications, Cannes, France, November 1997.
- [12] Crowcroft, J., Rizzo, L., and Vicisano, L. "TCP-Like Congestion Control for Layered Multicast Data Transfer." Submitted to INFOCOM '98, August 1997.
- [13] Sano, T., Yamanouchi, N., et al. "Flow and Congestion Control for Bulk Reliable Multicast Protocols—toward coexistence with TCP." Submitted to INFOCOM '98, presented at RMRG meeting in Cannes, France, September 1997.
- [14] Hofmann, M. "Enabling Group Communication in Global Networks." Proceedings of Global Networking '97, June 1997.
- [15] Papadopoulos, C., Parulkar, G., and Varghese, G. "An Error Control Scheme for Large-Scale Multicast Applications." Submitted to INFOCOM '98, presented at RMRG meeting in Cannes, France, September 1997.

- [16] Levine, B. N., Paul, S., and Garcia-Luna-Aceves, J. J. "Deterministic Organization of Multicast Receivers Based on Packet Loss Correlation." Presented at RMRG meeting in Orlando, Fla., February 1998, submitted for publication.
- [17] Handley, M. "SAP: Session Announcement Protocol." Work in progress, Internet Draft, **draft-ietf-mmusic-sap-00.txt**, November 1996.
- [18] Handley, M., and Jacobson, V. "SDP: Session Description Protocol." Work in progress, Internet Draft, **draft-ietf-mmusic-sdp-07.txt**, April 1998.
- [19] Handley, M., Schulzrinne, H., and Schooler, E. "SIP: Session Invitation Protocol." Work in progress, Internet Draft, **draft-ietf-mmusic-sip-04.txt**, November 1997.

(This article is based in part on material in the book *Multicast Networking and Applications* written by C. Kenneth Miller to be published by Addison Wesley Longman, Inc. in 1998. ISBN 0-201-30979-3. Used with permission.)

C. KENNETH MILLER is the founder, Chairman, and Chief Technology Officer of StarBurst Communications. StarBurst Communications provides reliable multicast solutions for commercial applications with such corporate customers as GM, Ford, Chrysler, Toys 'R Us, Thomson Financial, and many others. Miller has been in the data communications industry since 1972. He founded Concord Data Systems in late 1980 and served as its President and CEO until 1986. Concord Data Systems produced high-speed dial modems. He was the author of the IEEE 802.4 LAN standard, which became the lower layer for the Manufacturing Automation Protocol (MAP) factory LAN standard. Miller was a regular columnist in *Data Communications Magazine* from 1992 to 1994. He has also published numerous articles and participated in many panels at trade show and other industry events. He is now writing a book entitled *Multicast Networking and Applications*, to be published in 1998 by Addison-Wesley. Miller received a BEE degree from Rensselaer Polytechnic Institute and a MSEE degree from the University of Pennsylvania, specializing in communications. Miller can be reached at miller@starburstcom.com

Layer 2 and Layer 3 Switch Evolution

by Thayumanavan Sridhar, Future Communications Software

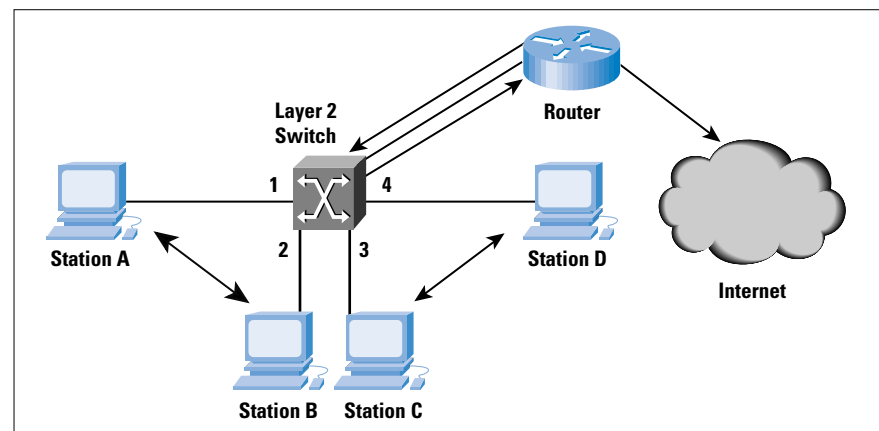
Layer 2 switches are frequently installed in the enterprise for high-speed connectivity between end stations at the data link layer. Layer 3 switches are a relatively new phenomenon, made popular by (among others) the trade press. This article details some of the issues in the evolution of Layer 2 and Layer 3 switches. We hypothesize that the technology is evolutionary and has its origins in earlier products.

Layer 2 Switches

Bridging technology has been around since the 1980s (and maybe even earlier). Bridging involves segmentation of local-area networks (LANs) at the Layer 2 level. A multiport bridge typically learns about the *Media Access Control* (MAC) addresses on each of its ports and transparently passes MAC frames destined to those ports. These bridges also ensure that frames destined for MAC addresses that lie on the same port as the originating station are not forwarded to the other ports. For the sake of this discussion, we consider only Ethernet LANs.

Layer 2 switches effectively provide the same functionality. They are similar to multiport bridges in that they learn and forward frames on each port. The major difference is the involvement of hardware that ensures that multiple switching *paths* inside the switch can be active at the same time. For example, consider Figure 1, which details a four-port switch with stations A on port 1, B on port 2, C on port 3 and D on port 4. Assume that A desires to communicate with B, and C desires to communicate with D. In a single CPU bridge, this forwarding would typically be done in software, where the CPU would pick up frames from each of the ports sequentially and forward them to appropriate output ports. This process is highly inefficient in a scenario like the one indicated previously, where the traffic between A and B has no relation to the traffic between C and D.

Figure 1:
Layer 2 switch with External Router
for Inter-VLAN traffic and connecting
to the Internet



Enter hardware-based Layer 2 switching. Layer 2 switches with their hardware support are able to forward such frames in parallel so that A and B and C and D can have simultaneous conversations. The parallelism has many advantages. Assume that A and B are NetBIOS stations, while C and D are Internet Protocol (IP) stations. There may be no reason for the communication between A and C and A and D. Layer 2 switching allows this coexistence without sacrificing efficiency.

Virtual LANs

In reality, however, LANs are rarely so *clean*. Assume a situation where A,B,C, and D are all IP stations. A and B belong to the same IP subnet, while C and D belong to a different subnet. Layer 2 switching is fine, as long as only A and B or C and D communicate. If A and C, which are on two different IP subnets, need to communicate, Layer 2 switching is inadequate—the communication requires an IP router. A corollary of this is that A and B and C and D belong to different broadcast domains—that is, A and B should not “see” the MAC layer broadcasts from C and D, and vice versa. However, a Layer 2 switch cannot distinguish between these broadcasts—bridging technology involves forwarding broadcasts to all other ports, and it cannot tell when a broadcast is restricted to the same IP subnet.

Virtual LANs (VLANs) apply in this situation. In short, Layer 2 VLANs are Layer 2 broadcast domains. MAC broadcasts are restricted to the VLANs that stations are configured into. How can the Layer 2 switch make this distinction? By configuration. VLANs involve configuration of ports or MAC addresses. Port-based VLANs indicate that all frames that originate from a port belong to the same VLAN, while MAC address-based VLANs use MAC addresses to determine VLAN membership. In Figure 1, ports 1 and 2 belong to the same VLAN, while ports 3 and 4 belong to a different VLAN. Note that there is an implicit relationship between the VLANs and the IP subnets—however, configuration of Layer 2 VLANs does not involve specifying Layer 3 parameters.

We indicated earlier that stations on two different VLANs can communicate only via a router. The router is typically connected to one of the switch ports (Figure 1). This router is sometimes referred to as a *one-armed router* since it receives and forwards traffic on to the same port. In reality, of course, such routers connect to other switches or to wide-area networks (WANs). Some Layer 2 switches provide this Layer 3 routing functionality within the same box to avoid an external router and to free another switch port. This scenario is reminiscent of the large multiprotocol routers of the early '90s, which offered routing and bridging functions.

A popular classification of Layer 2 switches is “cut-through” versus “store-and-forward.” Cut-through switches make the forwarding decision as the frame is being received by just looking at the header of the frame. Store-and-forward switches receive the entire Layer 2 frame

before making the forwarding decision. Hybrid adaptable switches which adapt from cut-through to store-and-forward based on the error rate in the MAC frames are very popular.

Characteristics

Layer 2 switches themselves act as IP end nodes for *Simple Network Management Protocol* (SNMP) management, Telnet, and Web based management. Such management functionality involves the presence of an IP stack on the router along with *User Datagram Protocol* (UDP), *Transmission Control Protocol* (TCP), Telnet, and SNMP functions. The switches themselves have a MAC address so that they can be addressed as a Layer 2 end node while also providing transparent switch functions. Layer 2 switching does not, in general, involve changing the MAC frame. However, there are situations when switches change the MAC frame. The IEEE 802.1Q Committee is working on a VLAN standard that involves “tagging” a MAC frame with the VLAN it belongs to; this tagging process involves changing the MAC frame. Bridging technology also involves the *Spanning-Tree Protocol*. This is required in a multibridge network to avoid loops. The same principles also apply towards Layer 2 switches, and most commercial Layer 2 switches support the Spanning-Tree Protocol.

The previous discussion provides an outline of Layer 2 switching functions. Layer 2 switching is MAC frame based, does not involve altering the MAC frame, in general, and provides transparent switching in parallel with MAC frames. Since these switches operate at Layer 2, they are protocol independent. However, Layer 2 switching does not scale well because of broadcasts. Although VLANs alleviate this problem to some extent, there is definitely a need for machines on different VLANs to communicate. One example is the situation where an organization has multiple intranet servers on separate subnets (and hence VLANs), causing a lot of intersubnet traffic. In such cases, use of a router is unavoidable; Layer 3 switches enter at this point.

Layer 3 Switches

Layer 3 switching is a relatively new term, which has been “extended” by a numerous vendors to describe their products. For example, one school uses this term to describe fast IP routing via hardware, while another school uses it to describe *Multi Protocol Over ATM* (MPOA). For the purpose of this discussion, Layer 3 switches are superfast routers that do Layer 3 forwarding in hardware. In this article, we will mainly discuss Layer 3 switching in the context of fast IP routing, with a brief discussion of the other areas of application.

Evolution

Consider the Layer 2 switching context shown in Figure 1. Layer 2 switches operate well when there is very little traffic between VLANs. Such VLAN traffic would entail a router—either “hanging off” one of the ports as a one-armed router or present internally within the switch. To augment Layer 2 functionality, we need a router—which

leads to loss of performance since routers are typically slower than switches. This scenario leads to the question: Why not implement a router in the switch itself, as discussed in the previous section, and do the forwarding in hardware?

Although this setup is possible, it has one limitation: Layer 2 switches need to operate only on the Ethernet MAC frame. This scenario in turn leads to a well-defined forwarding algorithm which can be implemented in hardware. The algorithm cannot be extended easily to Layer 3 protocols because there are multiple Layer 3 routable protocols such as IP, IPX, AppleTalk, and so on; and second, the forwarding decision in such protocols is typically more complicated than Layer 2 forwarding decisions.

What is the engineering compromise? Because IP is the most common among all Layer 3 protocols today, most of the Layer 3 switches today perform IP switching at the hardware level and forward the other protocols at Layer 2 (that is, bridge them). The second issue of complicated Layer 3 forwarding decisions is best illustrated by IP option processing, which typically causes the length of the IP header to vary, complicating the building of a hardware forwarding engine. However, a large number of IP packets do not include IP options—so, it may be overkill to design this processing into silicon. The compromise is that the most common (fast path) forwarding decision is designed into silicon, whereas the others are handled typically by a CPU on the Layer 3 switch.

To summarize, Layer 3 switches are routers with fast forwarding done via hardware. IP forwarding typically involves a route lookup, decrementing the *Time To Live* (TTL) count and recalculating the checksum, and forwarding the frame with the appropriate MAC header to the correct output port. Lookups can be done in hardware, as can the decrementing of the TTL and the recalculation of the checksum. The routers run routing protocols such as *Open Shortest Path First* (OSPF) or *Routing Information Protocol* (RIP) to communicate with other Layer 3 switches or routers and build their routing tables. These routing tables are looked up to determine the route for an incoming packet.

Combined Layer 2/Layer 3 Switches

We have implicitly assumed that Layer 3 switches also provide Layer 2 switching functionality, but this assumption does not always hold true. Layer 3 switches can act like traditional routers hanging off multiple Layer 2 switches and provide inter-VLAN connectivity. In such cases, there is no Layer 2 functionality required in these switches. This concept can be illustrated by extending the topology in Figure 1—consider placing a pure Layer 3 switch between the Layer 2 Switch and the router. The Layer 3 Switch would off-load the router from inter-VLAN processing.

Figure 2:
Combined Layer2/
Layer3 Switch
connecting directly
to the Internet

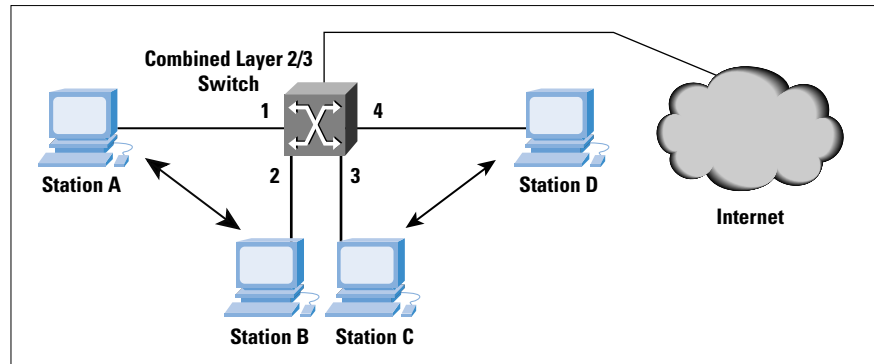


Figure 2 illustrates the combined Layer 2/Layer 3 switching functionality. The combined Layer 2/Layer 3 switch replaces the traditional router also. A and B belong to IP subnet 1, while C and D belong to IP subnet 2. Since the switch in consideration is a Layer 2 switch also, it switches traffic between A and B at Layer 2. Now consider the situation when A wishes to communicate with C. A sends the IP packet addressed to the MAC address of the Layer 3 switch, but with an IP destination address equal to C's IP address. The Layer 3 switch strips out the MAC header and switches the frame to C after performing the lookup, decrementing the TTL, recalculating the checksum and inserting C's MAC address in the destination MAC address field. All of these steps are done in hardware at very high speeds.

Now how does the switch know that C's IP destination address is Port 3? When it performs learning at Layer 2, it only knows C's MAC address. There are multiple ways to solve this problem. The switch can perform an *Address Resolution Protocol* (ARP) lookup on all the IP subnet 2 ports for C's MAC address and determine C's IP-to-MAC mapping and the port on which C lies. The other method is for the switch to determine C's IP-to-MAC mapping by snooping into the IP header on reception of a MAC frame.

Characteristics

Configuration of the Layer 3 switches is an important issue. When the Layer 3 switches also perform Layer 2 switching, they learn the MAC addresses on the ports—the only configuration required is the VLAN configuration. For Layer 3 switching, the switches can be configured with the ports corresponding to each of the subnets or they can perform IP address learning. This process involves snooping into the IP header of the MAC frames and determining the subnet on that port from the source IP address. When the Layer 3 switch acts like a one-armed router for a Layer 2 switch, the same port may consist of multiple IP subnets.

Management of the Layer 3 switches is typically done via SNMP. Layer 3 switches also have MAC addresses for their ports—this setup can be one per port, or all ports can use the same MAC address. The Layer 3 switches typically use this MAC address for SNMP, Telnet, and Web management communication.

Conceptually, the ATM Forum's *LAN Emulation* (LANE) specification is closer to the Layer 2 switching model, while MPOA is closer to the Layer 3 switching model. Numerous Layer 2 switches are equipped with ATM interfaces and provide a LANE client function on that ATM interface. This scenario allows the bridging of MAC frames across an ATM network from switch to switch. The MPOA is closer to combined Layer2/Layer 3 switching, though the MPOA client does not have any routing protocols running on it. (Routing is left to the MPOA server under the Virtual Router model.)

Do Layer 3 switches completely eliminate need for the traditional router ? No, routers are still needed, especially where connections to the wide area are required. Layer 3 switches may still connect to such routers to learn their tables and route packets to them when these packets need to be sent over the WAN. The switches will be very effective on the workgroup and the backbone within an enterprise, but most likely will not replace the router at the edge of the WAN (read Internet in many cases). Routers perform numerous other functions like filtering with access lists, inter-Autonomous System (AS) routing with protocols such as the *Border Gateway Protocol* (BGP), and so on. Some Layer 3 switches may completely replace the need for a router if they can provide all these functions (see Figure 2).

References

- [1] *Computer Networks*, 3rd Edition, Andrew S. Tanenbaum, ISBN 0-13-349945-6, Prentice-Hall, 1996.
- [2] *Interconnections: Bridges and Routers*, Radia Perlman, ISBN 0-201-56332-0, Addison-Wesley, 1992.
- [3] "MAC Bridges," ISO/IEC 10038, ANSI/IEEE Standard 802.1 D-1993.
- [4] "Draft Standard for Virtual Bridged Local Area Networks," IEEE P802.1Q/D6, May 1997.
- [5] "Internet Protocol," Jon Postel, RFC 791, 1981.
- [6] "Requirements for IP Version 4 Routers," Fred Baker, RFC 1812, June 1995.
- [7] "LAN Emulation over ATM Version 1.0," **af-lane-0021.000**, The ATM Forum, January 1995.
- [8] "Multiprotocol over ATM (MPOA) Specification Version 1.0" **af-mpoa-0087.000**, The ATM Forum, July 1997.

THAYUMANAVAN SRIDHAR is Director of Engineering at Future Communications Software in Santa Clara, CA. He received his BE in Electronics and Communications Engineering from the College of Engineering, Guindy, Anna University, Madras, India, his Master of Science in Electrical and Computer Engineering from the University of Texas at Austin. He can be reached at sridhar@futsoft.com

Book Review

Gigabit Ethernet *Gigabit Ethernet: Technology and Applications for High-Speed LANs*, by Rich Seifert, ISBN 0-201-18553-9, Addison-Wesley, 1998, <http://www.awl.com/cseng/titles/0-201-18553-9>.

Gigabit Ethernet is storming its way onto the high-speed LAN scene. From a concept in 1984 to an emerging commercial reality in 1998, Gigabit Ethernet promises to give other high-speed LAN technologies, especially ATM, a serious run for their money. Capitalizing on the basic ease of use and deployment that has made other forms of Ethernet the most popular LAN technology of all, Gigabit Ethernet promises to add major bandwidth to such networks in a straightforward, completely compatible, and relatively affordable way. This book performs an excellent survey of the technologies, algorithms, and design principles that make Gigabit Ethernet possible, and also explains where the tremendous appeal of Gigabit Ethernet really lies. Much of the book is devoted to explaining Ethernet principles and operation in general, as well as exploring recent developments that have enabled gigabit technologies to emerge.

Organization

The book is divided into three parts. Part I explores the foundations that underpin Gigabit Ethernet, starting with a brief but cogent exploration of Ethernet before gigabit versions loomed on the horizon. The rest of Part I covers the trends in LAN usage in general, and Ethernet in particular, that laid the groundwork for Gigabit Ethernet. These trends include the move from shared media to dedicated media on many LANs, and likewise from shared LANs to dedicated LANs, and the concomitant deployment of full-duplex technologies to support bidirectional, high-bandwidth communications. Seifert, an original member of the DIX (Digital-Intel-Xerox) team that developed Ethernet, writes clearly and compellingly about complex issues, such as flow control, medium independence, and automatic configuration, as he explains what made Gigabit Ethernet possible, if not inevitable.

In Part II, Seifert turns his focus onto Gigabit Ethernet itself, beginning with an overview. In the rest of Part II, he explains how *Media Access Control* (MAC) works for half-duplex and full-duplex versions of Gigabit Ethernet, and makes a strong case for the essential irrelevancy of shared-media and half-duplex operation for Gigabit Ethernet. Along the way, Seifert also covers how Gigabit Ethernet networking devices, such as repeaters and switching and routing hubs, must be designed and how they work, and covers the behavior and operation of the physical layer at gigabit speeds.

He concludes this section of the book with a brief overview of the current IEEE Draft 802.3z specification that governs current Gigabit Ethernet operations, and mentions ongoing work in the 802.3ab subcommittee to define a workable implementation for Gigabit Ethernet on twisted-pair media (1000BaseT, as it will probably be known).

In Part III, Seifert tackles some of the most interesting material in this book. He begins with a discussion of how LANs and computers change roles over time in acting as the bottleneck for network use. The point here is that because of its extremely high bandwidth relative to the demands of most applications and end-user requirements, Gigabit Ethernet is likely to remain a backbone or clustering technology for the foreseeable future. He also explores the performance considerations for both networks and applications involved when extreme speeds or excessive bandwidths are available, to point out how bandwidth aggregation is presently Gigabit's most immediate and compelling contribution to networking.

Finally, he explores how Gigabit Ethernet compares to other high-speed networking technologies, including Fast Ethernet, *Fiber Distributed Data Interface* (FDDI), *High-Performance Parallel Interface* (HIPPI), *Fibre Channel*, and ATM. His discussion of why both ATM and Gigabit Ethernet are necessary, and why neither can fully supplant the other, represents a humorous and insightful analysis of why connection-oriented and connectionless communications and applications are both good, and why the two can never truly converge.

An Outstanding Contribution

A rundown of Seifert's layout and content, however, fails to do complete justice to this book. For one thing, Seifert's work includes the funniest and most ingenious footnotes I've seen in recent publications, including some truly horrendous puns and some downright howlers. For example, when discussing how repeaters work, he comments that "A jabbering station causes carrier sense to be continuously asserted and blocks all use of a shared LAN. A repeater looks for this condition and isolates the offending station." To this last sentence, he appends the following footnote: "Research is underway to determine if this mechanism can be extended for use on politicians and university lecturers." And this is just one of dozens of such gems that help to relieve the dryness that deeply technical material can sometimes manifest.

This book is also masterful simply because the author understands his material so well, and does such an outstanding job of explaining and exploring even the most abstruse networking concepts. Although I've been working with Ethernet for 15 years, I learned a great deal of new material from Part I of the book because old concepts were explained in new ways that improved my understanding. I suspect other readers will have one or two "Aha!" experiences from this tome as well.

But it's when making the case for full-duplex Gigabit Ethernet and exploring the requirements for switching and routing behaviors in Gigabit Ethernet networking devices that this material really shines.

Without a doubt, this book is among the very best of any of the literature available on high-speed networking today. I give it an A+ rating, not only because of the breadth and depth of its technical coverage and its compilation of essential concepts and information, but also because the author's deep understanding of networking protocols and communications needs enlivens all of his discussions of matters technical, business, and political. If you want to understand Gigabit Ethernet, this book is the obvious place to begin (and for many, to end) your search for enlightenment.

But even if all you want is a good read about expensive, exotic, and high-performance technology, Seifert's book offers the opportunity for outright enjoyment of the prose, and shared delight at untangling the technical dilemmas that any good design engineer must unravel on the road between a set of requirements and working implementation thereof.

—Ed Tittel
LANWrights, Inc.
etittel@lanw.com

More Book Reviews

We have more book reviews awaiting publication:

- *Internet Cryptography*, by Richard E. Smith, ISBN 0-201-92480-3, Addison-Wesley, 1998. Reviewed by Fred Avolio.
- *Web Security: A Step-by-Step Reference Guide*, by Lincoln D. Stein, ISBN 0-201-63489-9, Addison-Wesley, December 1997. Reviewed by Richard Perlman
- *IP Multicasting: The Complete Guide to Interactive Corporate Networks*, by Dave Kosiur ISBN 0-471-24359-0, Wiley Computer Publishing, 1998. Reviewed by Neophytos Iacovou.

So, make sure you receive the next issue of *The Internet Protocol Journal* due out in December 1998.

Fragments

More on The Future of the Domain Name System (DNS)

Shortly after our first issue went to press, the US Government issued a so-called White Paper as a follow on to the Green Paper. The White Paper, entitled “Management of Internet Names and Addresses,” can be found at:

<http://www.ntia.doc.gov/ntiahome/domainname/domainhome.htm>

In early July, *The International Forum on The White Paper* (IFWP) was formed. The IFWP is “an ad hoc coalition of professional, trade and educational associations representing a diversity of Internet stakeholder groups.” The IFWP held a series of meetings in Reston, Brussels, Geneva, Singapore and Buenos Aires to discuss the White Paper, specifically the incorporation of the *Internet Assigned Numbers Authority* (IANA). For more information on the IFWP process, see: <http://www.ifwp.org>

The IANA has posted draft bylaws for its incorporation on the IANA web site at: <http://www.iana.org>, and asked for community input. By the time you read this, the incorporation should already have taken place. We will provide an update in our next issue.

IETF Wins Award

The *Computer Professionals for Social Responsibility* (CPSR) has chosen the *Internet Engineering Task Force* (IETF) to be honored with the Norbert Wiener award for the group’s influential role in the evolution of the Internet. In its 12-year history, this is only the second time the CPSR has recognized an organization rather than an individual. The IETF will accept the award at CPSR’s annual conference, on Saturday evening, October 10, 1998, in Boston. The IETF is noted for its highly open and democratic processes that have affected the development of the Internet. The CPSR believes that such open processes are both extremely important and seriously threatened, and have accordingly made Internet governance the focus of its 1998 program year. The Norbert Wiener award was established in 1987 by the CPSR in memory of the originator of the field of cybernetics, whose pioneering work was one of the pillars on which the computer technology was created. See: <http://www.cpsr.org> and <http://www.ietf.org>

Send us your comments!

We look forward to hearing your comments and suggestions regarding anything you read in this publication. Send us e-mail at: ipj@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or noninfringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Engineering
MCI Communications, USA

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Edward R. Kozel, Sr. VP, Corporate Development
Cisco Systems, Inc., USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President,
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the Cisco News
Publications Group, Cisco Systems, Inc.
www.cisco.com*

*Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 1998 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-J4
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

Bulk Rate Mail
U.S. Postage
PAID
Cisco Systems, Inc.

The Internet Protocol Journal

December 1998

Volume 1, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
SNMPv3	2
CATV Internet Technology ...	13
Digital TV	27
I Remember IANA.....	38
Book Reviews	40
Call for Papers	46
Fragments	47

FROM THE EDITOR

The *Simple Network Management Protocol* (SNMP) was first standardized in 1988. It quickly became a de facto management standard, not only for Internet technologies, but for a wide range of applications. Like many early Internet protocols, the first two versions of SNMP did not include provisions for security. In 1996, two different proposals for security enhancements to SNMPv2 were put forward, with strong proponents behind each. Everyone agreed that the industry needed just *one* solution, and therefore work proceeded to incorporate the best features of the two security proposals for SNMPv2. The result is SNMPv3, and it is described in this issue by William Stallings.

As the Internet continues to grow, demand for high-speed access for residential users is increasing. Alternatives to traditional dialup service include *Digital Subscriber Line* (DSL) services, wireless solutions, and various television technologies. In this issue, we examine two aspects of Internet access using TV technologies. First, Mark Laubach gives an overview of cable modem technologies and standards, and discusses some deployment issues. In the second article, George Abe looks at the emerging digital television standards and how they could be used to provide Internet access.

The Internet lost one of its most respected pioneers when Jon Postel passed away on October 16, 1998. Jon was well-known as the Director of the *Internet Assigned Numbers Authority* (IANA) and as the editor of the *Request for Comments* (RFC) document series. Included in this issue is "I Remember IANA," a tribute to Jon Postel written by his longtime friend Vint Cerf. The remembrance has also been published as RFC 2468.

With that we have come to the end of 1998 and the end of Volume 1 of *The Internet Protocol Journal*. We wish you a pleasant holiday season and will be back with Volume 2, Number 1 in March 1999. In the meantime, please visit our Web site at www.cisco.com/ipj. There you will find back issues in PDF format, our Call for Papers and guidelines for authors of IPJ articles.

You can download
previous issues of IPJ in
PDF format from:
www.cisco.com/ipj

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

Security Comes to SNMP: The New SNMPv3 Proposed Internet Standards

by William Stallings

Data networks typically include bridges, routers, links into WANs, and end-user equipment from multiple vendors. Users need automated tools to help manage such configurations that are easy to install, easy to use, and don't place a great burden on the network.

This accounts for the popularity of the *Simple Network Management Protocol* (SNMP). Introduced in 1988 to provide management capability for TCP/IP-based networks, SNMP rapidly became the most widely used standardized network management tool. Virtually all vendors of network-based equipment provide SNMP.

The appeal of SNMP has indeed been its simplicity because SNMP provides a bare-bones set of functions, and it is indeed easy to implement, install, and use. And, used sensibly, it will not place undue burden on the network. Moreover, because of its simplicity, achievement of interoperability is a relatively straightforward task: SNMP modules from different vendors can be made to work together with minimal effort.

SNMP—Strengths and Weaknesses

SNMP is based on three concepts: *managers*, *agents*, and the *Management Information Base* (MIB). In any configuration, at least one manager node runs SNMP management software. Network devices to be managed, such as bridges, routers, servers, and workstations, are equipped with an agent software module. The agent is responsible for providing access to a local MIB of objects that reflects the resources and activity at its node. The agent also responds to manager commands to retrieve values from the MIB and to set values in the MIB. An example of an object that can be retrieved is a counter that keeps track of the number of packets sent and received over a link into the node; the manager can track this value to monitor the load at that point in the network. An example of an object that can be set is one that represents the state of a link; the manager could disable the link by setting the value of the corresponding object to the disabled state.

Such capabilities are fine for implementing a basic network-management system. To enhance this basic functionality, a new version of SNMP was introduced in 1993 and revised in 1996. SNMPv2 added bulk transfer capability and other functional extensions. However, neither SNMPv1 nor SNMPv2 offers security features. Specifically, SNMPv1/v2 can neither authenticate the source of a management message nor provide encryption. Without authentication, it is possible for nonauthorized users to exercise SNMP network management functions. It is also possible for nonauthorized users to eavesdrop on

management information as it passes from managed systems to the management system. Because of these deficiencies, many SNMPv1/v2 implementations are limited to simply a read-only capability, reducing their utility to that of a network monitor; no network control applications can be supported.

Enter SNMPv3

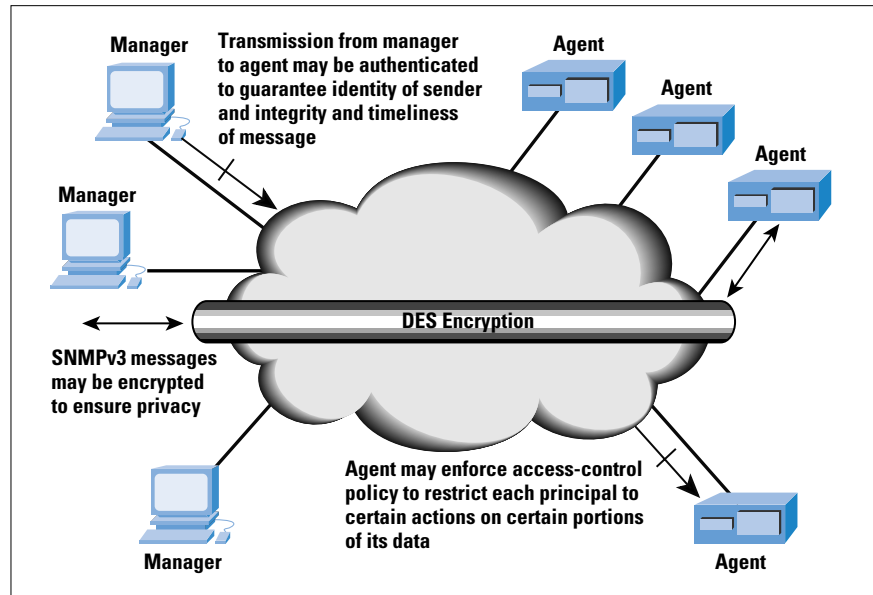
To correct the security deficiencies of SNMPv1/v2, SNMPv3 was issued as a set of Proposed Standards in January 1998 (Table 1). This set of documents does not provide a complete SNMP capability but rather defines an overall SNMP architecture and a set of security capabilities. These are intended to be used with the existing SNMPv2. As one of the SNMPv3 working documents puts it, “SNMPv3 is SNMPv2 plus administration and security.”

Table 1: SNMPv3 RFCs

RFC Number	Title
2271	An Architecture for Describing SNMP Management Frameworks
2272	Message Processing and Dispatching for the Simple Network Management Protocol (SNMP)
2273	SNMPv3 Applications
2274	User-Based Security Model for SNMPv3
2275	View-Based Access Control Model (VACM) for SNMP

SNMPv3 includes three important services: *authentication*, *privacy*, and *access control* (Figure 1). To deliver these services in a flexible and efficient manner, SNMPv3 introduces the concept of a *principal*, which is the entity on whose behalf services are provided or processing takes place. A principal can be an individual acting in a particular role; a set of individuals, with each acting in a particular role; an application or set of applications; or combinations thereof. In essence, a principal operates from a management station and issues SNMP commands to agent systems. The identity of the principal and the target agent together determine the security features that will be invoked, including authentication, privacy, and access control. The use of principals allows security policies to be tailored to the specific principal, agent, and information exchange, and gives human security managers considerable flexibility in assigning network authorization to users.

Figure 1:
SNMPv3 Security
Features



SNMPv3 is defined in a modular fashion, as shown in Figure 2. Each SNMP entity includes a single SNMP *engine*. An SNMP engine implements functions for sending and receiving messages, authenticating and encrypting/decrypting messages, and controlling access to managed objects. These functions are provided as services to one or more applications that are configured with the SNMP engine to form an SNMP *entity*. This modular architecture provides several advantages. First, the role of an SNMP entity is determined by the modules that are implemented in that entity. For example, a certain set of modules is required for an SNMP agent, whereas a different (though overlapping) set of modules is required for an SNMP manager. Second, the modular structure of the specification lends itself to defining different versions of each module. This, in turn, makes it possible to (1) define alternative or enhanced capabilities for certain aspects of SNMP without needing to go to a new version of the entire standard (for example, SNMPv4), and (2) clearly specify coexistence and transition strategies.

Figure 2:
SNMP Entity
(RFC 2271)

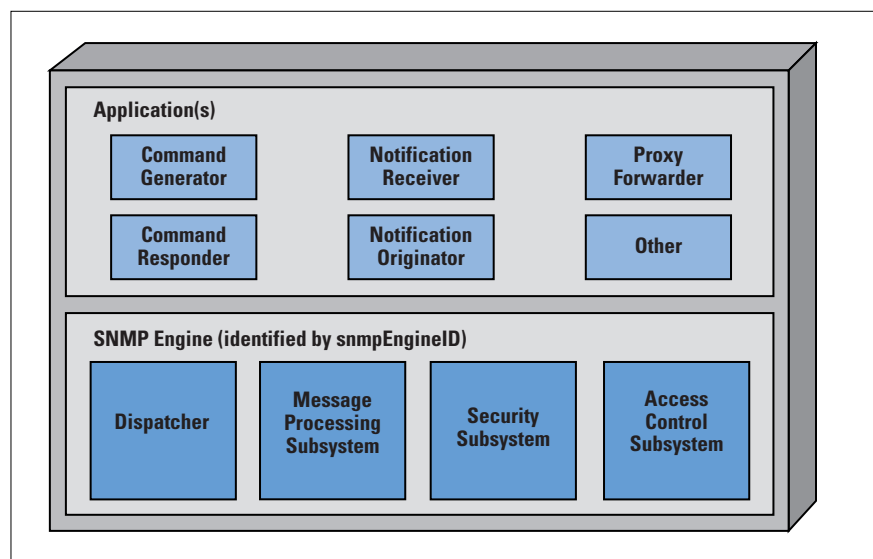


Table 2 provides a brief definition of each module.

Table 2: Components of an SNMP Entity (RFC 2271 and 2273)

Dispatcher	Allows for concurrent support of multiple versions of SNMP messages in the SNMP engine. It is responsible for (1) accepting protocol data units (PDUs) from applications for transmission over the network and delivering incoming PDUs to applications; (2) passing outgoing PDUs to the Message Processing Subsystem to prepare as messages, and passing incoming messages to the Message Processing Subsystem to extract the incoming PDUs; and (3) sending and receiving SNMP messages over the network.
Message Processing Subsystem	Responsible for preparing messages for sending and for extracting data from received messages.
Security Subsystem	Provides security services such as the authentication and privacy of messages. This subsystem potentially contains multiple Security Models.
Access Control Subsystem	Provides a set of authorization services that an application can use for checking access rights. Access control can be invoked for retrieval or modification request operations and for notification generation operations.
Command Generator	Initiates SNMP Get, GetNext, GetBulk, or Set request PDUs and processes the response to a request that it has generated.
Command Responder	Receives SNMP Get, GetNext, GetBulk, or Set request PDUs destined for the local system as indicated by the fact that the contextEngineID in the received request is equal to that of the local engine through which the request was received. The command responder application performs the appropriate protocol operation, using access control, and generates a response message to be sent to the originator of the request.
Notification Originator	Monitors a system for particular events or conditions, and generates Trap or Inform messages based on these events or conditions. A notification originator must have a mechanism for determining where to send messages, and which SNMP version and security parameters to use when sending messages.
Notification Receiver	Listens for notification messages, and generates response messages when a message containing an Inform PDU is received.
Proxy Forwarder	Forwards SNMP messages. Implementation of a proxy forwarder application is optional.

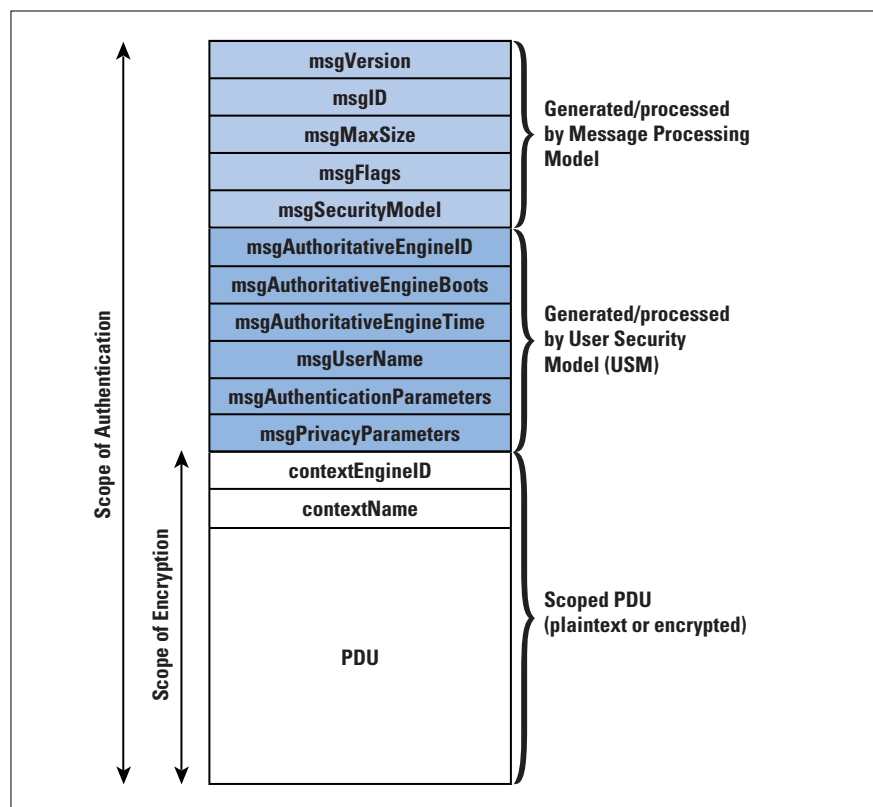
SNMPv3 Message Processing

SNMPv3 relies on the *User Datagram Protocol* (UDP) or some other transport-layer protocol to convey SNMP information. Above the UDP layer, SNMP functionality is organized into two application-level layers: a PDU processing layer and a message processing layer.

The topmost layer is the PDU processing layer. At this layer, management commands (such as Get, Set, Trap, Inform) are realized in a PDU that includes an indication of the command type and a list of variables (management objects) to which the command refers. This PDU is then passed down to the message processing layer, which adds a message header. The message header contains security-related information that may be used for authentication and privacy operations.

Figure 3 illustrates the message structure. The first five fields are generated by the message processing model on outgoing messages and processed by the message processing model on incoming messages. The next six fields show security parameters used by the security model, which is invoked by the message processing model to provide security services. Finally, the PDU, together with the contextEngineID and contextName, constitute a scoped PDU, used for PDU processing.

Figure 3:
SNMPv3 Message
Format with
User-Based
Security Model



The first five fields follow:

- *msgVersion*: Set to snmpv3(3).
- *msgID*: A unique identifier used between two SNMP entities to coordinate request and response messages, and by the message processor to coordinate the processing of the message by different subsystem models within the architecture. The range of this ID is 0 through $2^{31}-1$.

- *msgMaxSize*: Conveys the maximum size of a message in octets supported by the sender of the message, with a range of 484 through $2^{31}-1$. This is the maximum segment size that the sender can accept from another SNMP engine (whether a response or some other message type).
- *msgFlags*: An octet string containing three flags in the least significant three bits: reportableFlag, privFlag, authFlag. If reportableFlag = 1, then a Report PDU must be returned to the sender under those conditions that can cause the generation of a Report PDU; when the flag is zero, a Report PDU may not be sent. The reportableFlag is set to 1 by the sender in all messages containing a request (Get, Set) or an Inform, and set to 0 for messages containing a Response, a Trap, or a Report PDU. The reportableFlag is a secondary aid in determining when to send a Report. It is used only in cases in which the PDU portion of the message cannot be decoded (for example, when decryption fails because of incorrect key). The privFlag and authFlag are set by the sender to indicate the security level that was applied to the message. For privFlag = 1, encryption was applied and for privFlag = 0, authentication was applied. All combinations are allowed except (privFlag = 1 AND authFlag = 0); that is, encryption without authentication is not allowed.
- *msgSecurityModel*: An identifier in the range of 0 through $2^{31}-1$ that indicates which security model was used by the sender to prepare this message and, therefore, which security model must be used by the receiver to process this message. Reserved values include 1 for SNMPv1, 2 for SNMPv2c, and 3 for SNMPv3.

User-Based Security Model

The *User-Based Security Model* (USM) uses the concept of an authoritative engine. In any message transmission, one of the two entities, transmitter or receiver, is designated as the authoritative SNMP engine, according to the following rules:

- When an SNMP message contains a payload that expects a response (for example, a Get, GetNext, GetBulk, Set, or Inform PDU), then the receiver of such messages is authoritative.
- When an SNMP message contains a payload that does not expect a response (for example, an SNMPv2-Trap, Response, or Report PDU), then the sender of such a message is authoritative.

Thus, for messages sent on behalf of a Command Generator and for Inform messages from a Notification Originator, the receiver is authoritative. For messages sent on behalf of a Command Responder or for Trap messages from a Notification Originator, the sender is authoritative. This designation serves two purposes:

- The timeliness of a message is determined with respect to a clock maintained by the authoritative engine. When an authoritative engine sends a message (Trap, Response, Report), it contains the current value of its clock, so that the nonauthoritative recipient can synchronize on that clock. When a nonauthoritative engine sends a message (Get, GetNext, GetBulk, Set, Inform), it includes its current estimate of the time value at the destination, allowing the destination to assess the timeliness of the message.
- A key localization process, described later, enables a single principal to own keys stored in multiple engines; these keys are localized to the authoritative engine in such a way that the principal is responsible for a single key but avoids the security risk of storing multiple copies of the same key in a distributed network.

When an outgoing message is passed to the USM by the Message Processor, the USM fills in the security-related parameters in the message header. When an incoming message is passed to the USM by the Message Processor, the USM processes the values contained in those fields. The security-related parameters include the following:

- *msgAuthoritativeEngineID*: The `snmpEngineID` of the authoritative SNMP engine involved in the exchange of this message. Thus, this value refers to the source for a Trap, Response, or Report, and to the destination for a Get, GetNext, GetBulk, Set, or Inform.
- *msgAuthoritativeEngineBoots*: The `snmpEngineBoots` value of the authoritative SNMP engine involved in the exchange of this message. The object `snmpEngineBoots` is an integer in the range 0 through $2^{31}-1$ that represents the number of times that this SNMP engine has initialized or reinitialized itself since its initial configuration.
- *msgAuthoritativeEngineTime*: The `snmpEngineTime` value of the authoritative SNMP engine involved in the exchange of this message. The object `snmpEngineTime` is an integer in the 0 through $2^{31}-1$ range that represents the number of seconds since this authoritative SNMP engine last incremented the `snmpEngineBoots` object. Each authoritative SNMP engine is responsible for incrementing its own `snmpEngineTime` value once per second. A non-authoritative engine is responsible for incrementing its notion of `snmpEngineTime` for each remote authoritative engine with which it communicates.
- *msgUserName*: The user (principal) on whose behalf the message is being exchanged.
- *msgAuthenticationParameters*: Null if authentication is not being used for this exchange; otherwise, this is an authentication parameter. For the current definition of USM, the authentication parameter is a message authentication code generated using an algorithm referred to as HMAC.

- *msgPrivacyParameters*: Null if privacy is not being used for this exchange; otherwise, this is a privacy parameter. For the current definition of USM, the privacy parameter is a parameter used in the encryption algorithm DES.

Secret-Key Authentication

The authentication mechanism in SNMPv3 assures that a received message was, in fact, transmitted by the principal whose identifier appears as the source in the message header. In addition, this mechanism assures that the message was not altered in transit and that it was not artificially delayed or replayed.

To achieve authentication, each pair of principal and remote SNMP engines that wishes to communicate must share a secret authentication key. The sending entity provides authentication by including a message authentication code with the SNMPv3 message it is sending. This code is a function of the contents of the message, the identity of the principal and engine, the time of transmission, and a secret key that should be known only to the sender and the receiver. The secret key must initially be set up outside of SNMPv3 as a configuration function. That is, the configuration manager or network manager is responsible for distributing initial secret keys to be loaded into the databases of the various SNMP managers and agents. This can be done manually or by using some form of secure data transfer outside of SNMPv3. When the receiving entity gets the message, it uses the same secret key to calculate the message authentication code again. If the receiver's version of the code matches the value appended to the incoming message, then the receiver knows that the message can only have originated from the authorized manager, and that the message was not altered in transit. The shared secret key between sending and receiving parties must be preconfigured.

Another aspect of USM authentication is timeliness verification. USM is responsible for assuring that messages arrive within a reasonable time window to protect against message delay and replay attacks. Two functions support this service: synchronization and time-window checking.

Each authoritative engine maintains two values, *snmpEngineBoots* and *snmpEngineTime*, that keep track of the number of boots since initialization and the number of seconds since the last boot. These values are placed in outgoing messages in the fields *msgAuthoritativeEngineBoots* and *msgAuthoritativeEngineTime*. A nonauthoritative engine maintains synchronization with an authoritative engine by maintaining local copies of *snmpEngineBoots* and *snmpEngineTime* for each remote authoritative engine with which it communicates. These values are updated on receipt of an authentic message from the remote authoritative engine. Between these message updates, the nonauthoritative

engine increments the value of `snmpEngineTime` for the remote authoritative engine to maintain loose synchronization. These values are inserted in outgoing messages intended for that authoritative engine.

When an authoritative engine receives a message, it compares the incoming boot and time values with its own boot and time values. If the boot values match and if the incoming time value is within 150 seconds of the actual time value, then the message is declared to be within the time window and, therefore, to be a timely message.

Privacy Using Conventional Encryption

The SNMPv3 USM privacy facility enables managers and agents to encrypt messages to prevent eavesdropping by third parties. Again, manager entity and agent entity must share a secret key. When privacy is invoked between a principal and a remote engine, all traffic between them is encrypted using the *Data Encryption Standard* (DES). The sending entity encrypts the entire message using the DES algorithm and its secret key, and sends the message to the receiving entity, which decrypts it using the DES algorithm and the same secret key. Again, the two parties must be configured with the shared key.

The *cipher-block-chaining* (CBC) mode of DES is used by USM. This mode requires that an initial value (IV) be used to start the encryption process. The `msgPrivacyParameters` field in the message header contains a value from which the IV can be derived by both sender and receiver.

View-Based Access Control Model (VACM)

The access control facility makes it possible to configure agents to provide different levels of access to the agent's MIB to different managers. An agent entity can restrict access to its MIB for a particular manager entity in two ways. First, it can restrict access to a certain portion of its MIB. For example, an agent may restrict most manager principals to viewing performance-related statistics and allow only a single designated manager principal to view and update configuration parameters. Second, the agent can limit the operations that a principal can use on that portion of the MIB. For example, a particular manager principal could be limited to read-only access to a portion of an agent's MIB. The access control policy to be used by an agent for each manager must be preconfigured; it essentially consists of a table that details the access privileges of the various authorized managers. Unlike authentication, which is done by user, access control is done by group, where a group may be a set of multiple users.

Figure 4:
VACM Flowchart

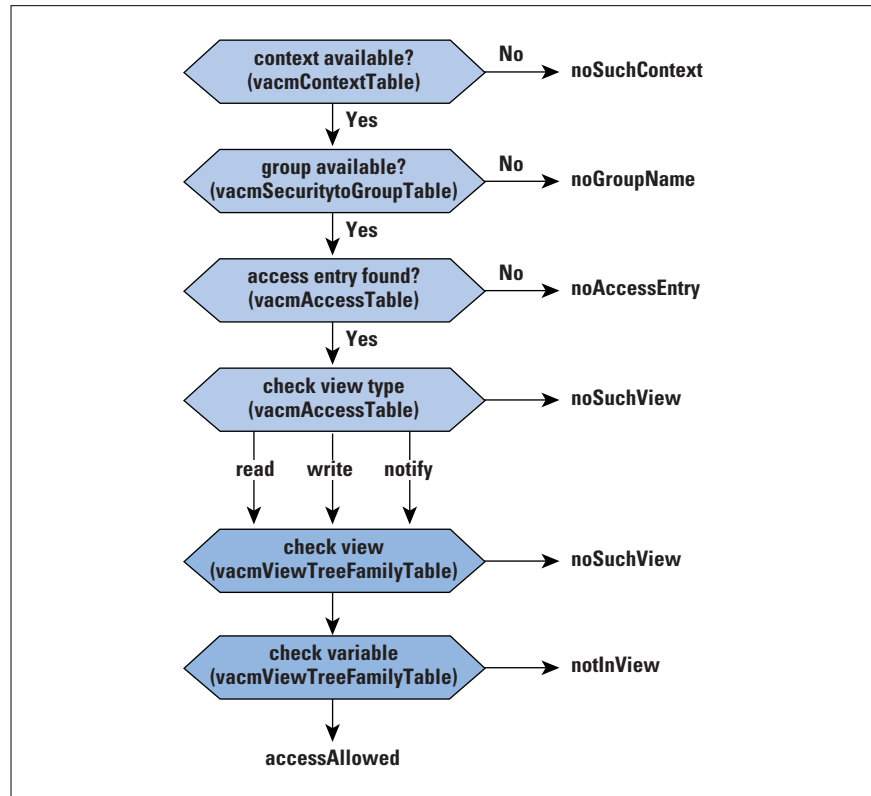


Figure 4 illustrates the overall VACM logic, which proceeds in the following steps:

1. The context name refers to a named subset of the MIB objects at an agent. VACM checks to see if there is an entry in `vacmContextTable` for the requested `contextName`. If so, then this context is known to this SNMP engine. If not, then an `errorIndication` of `noSuchContext` is returned.
2. Each principal operating under a given security model is assigned to at most one group, and access privileges are configured on a group basis. VACM checks `vacmSecurityToGroupTable` to determine if there is a group assigned to the requested `<securityModel, securityName>` pair. If so, then this principal, operating under this `securityModel`, is a member of a group configured at this SNMP engine. If not, then an `errorIndication` of `noGroupName` is returned.
3. VACM next consults the `vacmAccessTable` with `groupName`, `contextName`, `securityModel`, and `securityLevel` (indicates authentication, authentication plus privacy, or neither) as indices. If an entry is found, then an access control policy has been defined for this `groupName`, operating under this `securityModel`, at this `securityLevel`, for access to this `contextName`. If not, then an `errorIndication` of `noAccessEntry` is returned.

4. A MIB view is a structure subset of a context; it is essentially a set of managed object instances viewed as a set for access control purposes. VACM determines whether the selected vacmAccessTable entry includes reference to a MIB view of viewType (read, write, notify). If so, then this entry contains a viewName for this combination of groupName, contextName, securityModel, securityLevel, and viewType. If not, then an errorIndication of noSuchView is returned.
5. The viewName from Step 4 is used as an index into vacmViewTreeFamilyTable. If a MIB view is found, then a MIB view has been configured for this viewName. If not, then an errorIndication of noSuchView is returned.
6. VACM checks the variableName against the selected MIB view. If this variable is included in the view, then a statusInformation of accessAllowed is returned. If not, then an errorIndication of notInView is returned.

References

- [0] The SNMPv3 RFCs, see Table 1 above.
- [1] J. D. Case, M. Fedor, M. L. Schoffstall, and C. Davin, "Simple Network Management Protocol," RFC 1157, May 1990.
- [2] Rose, M. T., *The Simple Book: An Introduction to Networking Management*, Revised Second Edition, Prentice-Hall, ISBN 0-13-451659-1, 1996.
- [3] Waters, G., Editor, "User-based Security Model for SNMPv2," RFC 1910, February 1996.
- [4] *ConneXions—The Interoperability Report*, Volume 10, No. 5, May 1996—Special Issue: "Network Management Today."

WILLIAM STALLINGS is a consultant, lecturer, and author of over a dozen books on data communications and computer networking. He has a PhD in computer science from M.I.T. This article is based on material in the author's latest book: *SNMP, SNMPv2, SNMPv3, and RMON 1 and 2*, Second Edition (Addison Wesley, 1998). His home in cyberspace is <http://www.shore.net/~ws> and he can be reached at ws@shore.net

Residential Area CATV Broadband Internet Technology: Current Status

by Mark Laubach, Com21, Inc.

Cable modem technology has entered commonplace discussion and is in the early stages of widespread deployment throughout the world. The capabilities provided by cable modems promise data bandwidth speeds far in excess of those provided by traditional telephone modem services. In North America the race is on between cable operators deploying services based on standardized cable modems and telephone companies deploying *Digital Subscriber Line* (DSL) services. Internet Service Providers (ISPs) are taking position to promote any method of delivering Internet services to and from the home and are helping to fuel the race. Initially these services will only provide higher-speed Internet access and improved access to major information services (for example, AOL). Cable modem service offerings promise higher than DSL speed to the subscriber and a promise that packet voice services will be available in 1999.

As an introduction to some of the issues surrounding cable modem technology, this article summarizes two of the standardization efforts: the IEEE 802.14 Cable TV Media Access Control and Physical Protocol working group and the North American Data Over Cable Service Interface Specification. Delivering a viable Internet service to a cable TV reached subscriber community has its own set of deployment issues that are briefly reviewed and summarized.

Background

Networks based on packet technology were first presented in 1964^[1]. Since then, and through numerous evolutionary steps, the Internet as we know it today was brought into existence. Today, packets are transmitted over most any media. The next economic and technical frontier is the mass deployment of moving packets over cable television (CATV) networks for serving the Internet to every home. There are several link layer approaches for delivering IP datagrams via cable modems. The always present debate of whether to use fixed or variable length packets continues in the cable modem world. This article presents overviews of two variations of cable modem protocols: first, the concept of sending small, fixed-sized packets over the CATV plant using 53-octet *Asynchronous Transfer Mode* (ATM) cells^[2], as is being defined in the public standards process of the IEEE 802.14 working group; and secondly, by sending variable-length packets (IP over Ethernet) as defined by the *Multimedia Cable Network System* (MCNS) *Data Over Cable Service Interface Specification* (DOCSIS) for the North American cable industry^[3]. As widely accepted standards normally motivate industrial focus and subsequent cost reduction due to vendor competitive pressures, there is an additional drive provided by North American cable operators to get the cost of the cable modem off their books and into retail channels.

The IEEE 802.14 Cable TV MAC and PHY Protocol working group is chartered with providing a single *Media Access Control* (MAC) and multiple physical sublayer (PHY) standard for cable TV networks. The efforts of 802.14 must support IEEE 802 layer services (including Ethernet) and must also be ATM compatible.

The DOCSIS specifications are managed by CableLabs on behalf of its cable television system operator members. The project was initiated by an organization called *Multimedia Cable Network System* (MCNS) Partners, L. P., which consists of Comcast Cable Communications, Cox Communications, Tele-Communications, Inc., and Time Warner Cable. In addition to MCNS, Rogers Cablesystems Limited, MediaOne, and CableLabs have all contributed to the DOCSIS documents, as have several networking and telecommunications vendors. DOCSIS documents describe the internal and external network interfaces for a system that allows bidirectional transfer of IP traffic, between the cable system head-end and customer premises, over a cable television system^[4].

The customer network interface in common use today is Ethernet 10BaseT. There is a mandate for a 10 Mbps Ethernet interface in the home. Subscriber access equipment can be a personal computer, X-Terminal, or any such device that supports the TCP/IP protocol suite. Future home interfaces from the cable modem will include the *Universal Serial Bus* (USB) and IEEE 1394 (*FireWire*).

IP Over CATV System Challenges

From an IP perspective, a CATV system almost appears to be another data link layer. However, experience gained thus far has demonstrated that the marriage of IP over CATV radio frequency (RF) channels is not as straightforward as IP over any other high-speed serial point-to-point link.

In the CATV space, the downstream channels in a cable plant (cable head-end to subscribers) is a point-to-multipoint channel. This does have very similar characteristics to transmitting over an Ethernet segment where one transmitter is being listened to by many receivers. The major difference is that baseband modulation has been replaced by a more densely modulated RF carrier with very sophisticated adaptive signal processing and *forward error correction* (FEC).

In the upstream direction (subscriber cable modems transmitting towards the head-end) the environment is many transmitters and one receiver. This introduces the need for precise scheduling of packet transmissions to achieve high utilization and precise power control so as to not overdrive the receiver or other amplifier electronics in the cable system. Since the upstream direction is like a single receiver with many antennas, the channels are much much more susceptible to interfering noise products^[5, 6]. In the cable industry, we generally

call this *ingress noise*. As ingress noise is an inherent part of CATV plants, the observable impact is an unfortunate rise in the average noise floor in the upstream channel. To overcome this noise jungle, upstream modulation is not as dense as in the downstream and we have to use more effective FEC as used in the downstream. There is a further complication that there are many upstream “ports” on a fully deployed *Hybrid Fiber-Coaxial* (HFC) plant that requires matching head-end equipment ports for high-speed data^[7].

To further the rub on the upstream channel use, the arcane regulations of the FCC from back in the mid 1980s mandated that upstream frequency spectrum be reserved on all cable plants, regardless of whether it was actually used. This was typically the 5–42 MHz region, leaving above 50 MHz for downstream transmissions. (Note that there are other regions available for upstream, but the overwhelming majority of cable plants only use 5–42 MHz.) This leaves precious little spectral bandwidth for upstream communications.

The existing environment for high-speed data protocols therefore provides for relatively clean bandwidth in the downstream direction, allowing for higher-speed data rate channels, while in the upstream, individual channels are of lesser data rate. However, multiple upstream channels can be used per downstream channel to get effective symmetric aggregate bandwidth. Typically, we speak of cable modem systems as providing asymmetric services (higher downstream data rate than upstream). Note though that this asymmetry closely matches what we expect initially for residential high-speed data services. That is, many more subscribers at home pulling things off the Internet via web services, than pushing data back in.

Modern modulation techniques provide for a range of data carrying capability (“baud rate”). A low order modulation rate called *Quadrature Phase Shift Keying* (QPSK) provides for two data bits per symbol encoding. *Quadrature Amplitude Modulation* (QAM) provides a lower order modulation of 16 QAM (four bits per symbol) through higher order rates of 64 QAM (six bits per symbol) and 256 QAM (eight bits per symbol). Low order modulations are more robust in higher average noise environments. Higher order modulations are least robust. Therefore, high order modulations are suitable for downstream channels due to the low noise performance, while the order of upstream channel modulation is heavily effected by noise. Typically, cable modem systems will see QPSK used for upstream channels. When the plant is very clean, noise-wise, 16 QAM may be used.

One additional challenge is that the speed of RF signals in fiber and coaxial cable is much less than the speed of light. For system deployments to be effective, the cable modem protocols must support cable modems out to a wire distance of 50 miles (80 km).

At these distances the round trip propagation delay will be on the order of 800 microseconds; which is several times the length of time it takes to transmit a 64-byte packet on the upstream channel. The IEEE and DOCSIS cable modem protocols have been engineered to overcome these propagation delays in order to increase channel utilization; that is demand-based scheduling of a slotted upstream channel coupled with precise station ranging and timing.

Another challenge is in using an IP-over Ethernet approach to providing a reliable public switched packet service to an abundance of subscribers. Traditional Ethernet networking has always relied on all the Ethernet stations being within the same administrative walls with all users sharing the same common interests. Not so with metropolitan area public access networks. Data communications must now be encrypted such that the privacy of user communications is not invaded by promiscuous neighbors. In addition, users are paying for access in this cable modem world, and any abusive behavior of users must be contained so as to not affect other users. This calls for sophisticated fairness scheduling in the head-end systems and the use of comprehensive cryptological and packet filtering techniques. It is all very complicated both to create, and to manage. Each standard has its own approach for dealing with these issues.

Where IP over CATV appeared to be fundamentally similar to Ethernet when the industry first started out, in reality it is not. High-speed cable data networking, as demonstrated by the work output from various standards activities, is fundamentally a new approach to what at first appeared to be similar old problems. It's not ALOHA anymore^[8], nor is it your grandfather's Ethernet^[9, 10].

IEEE 802.14 Cable TV MAC and PHY Protocol Working Group

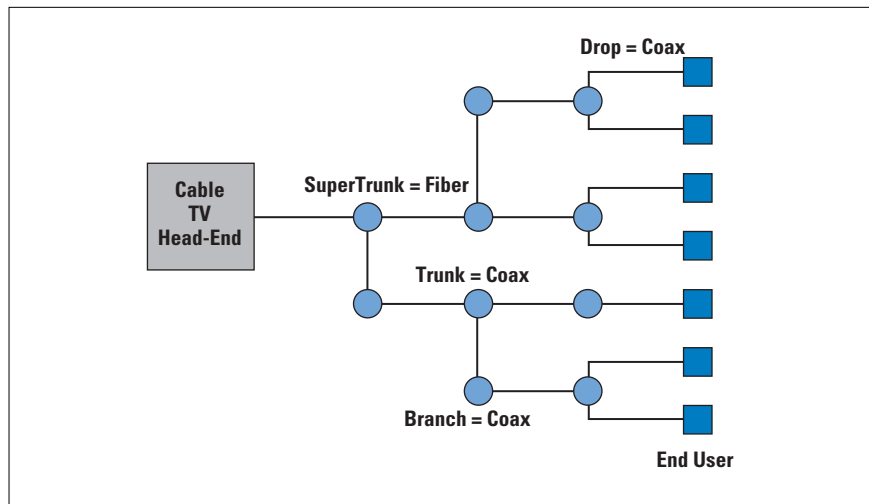
Let's briefly examine the first comprehensive standard activity created to address the current emerging world of high-speed cable data systems. In November 1994, the IEEE 802.14 CATV MAC and PHY Protocol working group met for the first time as an approved project within the 802 standards committee. Previous work had been done in 1993 through 1994 in the 802.catv study group in preparation for formal IEEE 802 project approval. The *Project Authorization Request* (PAR) charter of the group specifies that it will standardize a single MAC layer protocol and multiple PHY layer protocols for two-way HFC networks. Consistent with the IEEE LAN/MAN 802 Reference Model^[11], 802.14 is producing a solution that supports the 802 protocol stack while at the same time supporting ATM in an ATM-compatible manner.

The general 802.14 requirements include:

- Communications support for all coaxial and hybrid fiber-coaxial cable TV network tree and branch topologies. (See Figure 1)

- Support of symmetrical and asymmetrical rates
- Support of *Operation, Administrations, and Maintenance* (OAM) functions
- Support of one-way delays on the order of 400 microseconds (round-trip delays to 800 microseconds)
- Support of a large number of users
- Support for moving data from an originating subnetwork to a destination subnetwork, which may be the same or a different one

Figure 1:
CATV Tree and
Branch Network



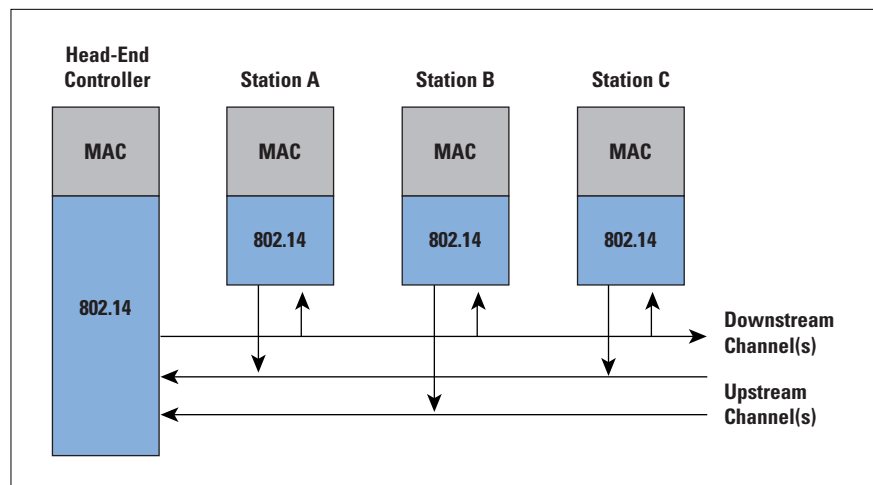
The working group completed a first-release revision of a functional requirements document back in 1995^[12], which detailed the 802.14 cable topology model; defined key assumptions, constraints, and parameters; defined key performance metrics and criteria for the selection of multiple PHY protocols and a MAC protocol; and defined the support of *Quality-of-Service* (QoS) parameters. The working group's work plan called for the close of formal proposals in November 1995, with the recommended protocol defined in July 1996. Seventeen MAC protocol proposals were submitted to the working group. Needless to say, it took awhile for the working group to sort through all the issues and opinions. After much consideration, debate, and wrangling of both solutions and personalities, IEEE 802.14 stabilized on a working group draft in September 1998. This working group draft is now being submitted through the IEEE 802 standard approval process.

The 802.14 MAC and PHY specification includes:

- Definition and operational specifications for cable system Head-End Controller and cable modem Stations. (See Figure 2)
- Support of both connectionless and connection-oriented services

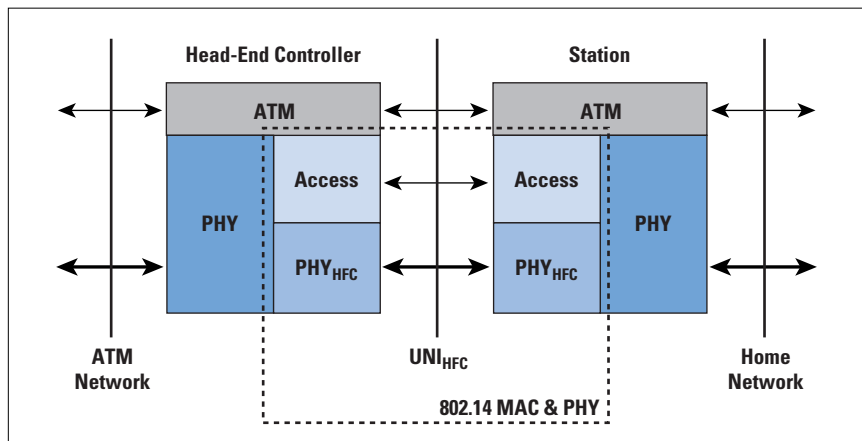
- Support of a formal QoS for connections; support for dynamically allocated bandwidth for different types of traffic, including *Constant Bit Rate* (CBR), *Variable Bit Rate* (VBR), and *Available Bit Rate* (ABR)
- Support for unicast, multicast, and broadcast services; interoperability with ATM
- Predictable low-average access delay without sacrificing network throughput
- Fair arbitration for shared access to the network within any level of service
- Downstream channel support for 64 QAM or 256 QAM modulation
- Compatibility for both international and North American downstream digital video standards
- Upstream channel support for QPSK or 16 QAM modulation

Figure 2:
IEEE 802.14 General
Model



The selection of ATM cells as the data link layer protocol data unit for IEEE 802.14 networks has the advantage that it provides a suitable integrated multiplexing platform capable of supporting a mix of guaranteed (predictive) traffic flows with best-effort (reactive) traffic flows. See Figure 3. Cable operators can deploy IEEE 802.14 based ATM systems as part of an evolutionary path to a fully integrated multimedia bearer service offering. A residential ATM bearer service easily supports Internet access to the home via the Classical IP over ATM standards of the Internet Engineering Task Force^[13] or by providing an IP over Ethernet adaptation overlay service^[14]. The development of QoS scheduling support in the Head-End Controller is left for vendors to implement^[15, 16, 17].

Figure 3:
IEEE 802.14 ATM
Protocol Model



IEEE 802.14 Status

At the time of this writing, the IEEE 802.14 working group just finalized a working group draft suitable to introduction into the IEEE standards process. The entire IEEE process takes about a year from acceptance of the working group letter ballot to producing a published standard.

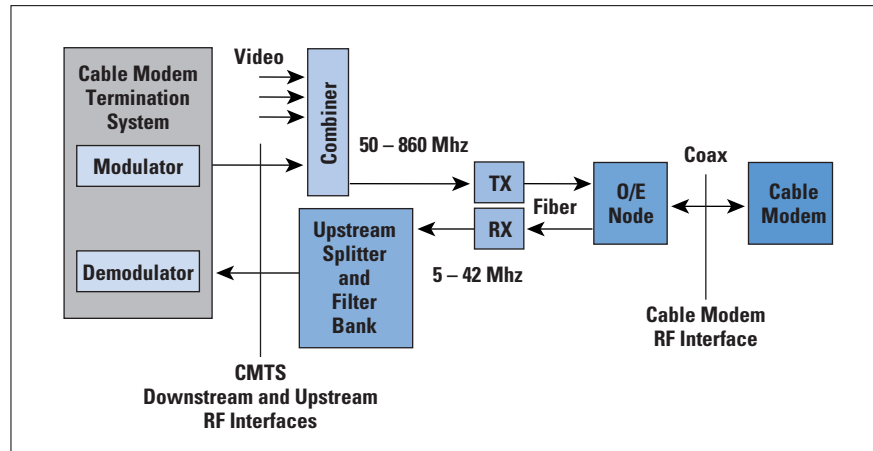
MCNS DOCSIS

The DOCSIS project is an activity of major cable companies and selected vendors to rapidly develop, on behalf of the North American cable industry, the necessary set of communications and operations support interface specifications for cable modems and associated equipment. The activity was triggered by John Malone in December 1995, in response to competition, vendor postures, and unfortunate lack of progress in the public standards process (that is, IEEE 802.14). The target for the specification was to produce a residential, “low-cost,” off-the-shelf, Internet access service, with wide-scale vendor interoperability for base functions with sufficient hooks and room for vendor differentiation.

MCNS specifications are intended to be non-vendor specific, allowing cross-manufacturer compatibility for high-speed data communications services over two-way HFC cable television systems. MCNS met its specification release deadline and published versions of the *DOCSIS Radio Frequency (RF) Interface Specification V1.0*. The first draft specification was published in December 1996. The latest specification was published in July 1998^[3]. The DOCSIS RFI protocol is based on the original LANCity symmetric 10 Mbps protocol, evolved to an asymmetric system, with multiple upstream and high-speed downstream (for example, 30 Mbps) channel support.

The MCNS system model is very similar to the IEEE 802.14 general model and includes many interfaces to a cable modem system, as shown in Figure 4. The goal of the DOCSIS project is to produce specifications for the CATV RF interfaces, including behavior of the *Cable Modem Termination System (CMTS)* and *Cable Modem (CM)* with respect to delivery of the residential IP over Ethernet service.

Figure 4:
Data-Over-Cable RFI
Reference
Architecture



The DOCSIS RFI system is asymmetric, with one to several downstream channels operating asymmetrically with one to several upstream channels. Specific features of MCNS DOCSIS RFI Version 1.0 include:

- Switched Ethernet service for Internet transport via a variable length MAC packet protocol
- Best-effort service
- Downstream data channel rates from 20 Mbps (16 QAM) to 40 Mbps (256 QAM) with a typical configuration of 30 Mbps (64 QAM) in 6 MHz channels
- Compatibility for North American downstream digital video standards. (See article starting on page 27.)
- Downstream data channel rates selected from 320 Kbps (QPSK) through 10.24 Mbps (16 QAM). Channel spectral widths from 200 KHz to 3.2 MHz
- Software flexibility: ability to download new software to change/update CM behavior
- Many filters and features for controlling packet flow and classification
- Comprehensive MIB specifications for control of the cable modem and cable modem termination system
- A single large LAN segment

Due to the time-to-market push for DOCSIS RFI V1.0 interoperable modems, little to no attention was been given for QoS needs however, vendors will likely include some QoS support in their offerings. (Upstream packet fragmentation was removed from the December 1996 draft release.)

CMs and the CMTSs have basically the same protocol stack: downstream and upstream PHY, the DOCSIS RFI MAC, Ethernet and an Ethernet switching layer with substantial filtering, IP/*Address Resolution Protocol* (ARP), *User Datagram Protocol* (UDP), and *Simple Network Management Protocol/Dynamic Host Configuration Protocol/Trivial File Transfer Protocol* (SNMP/DHCP/TFTP).

The DOCSIS RFI includes upstream and downstream optional packet encryption using the *Data Encryption Standard* (DES) to provide link privacy. RSA public key exchange is used between the CM and CMTS.

DOCSIS RFI Status

CableLabs is actively driving multiple vendor interoperability with the goal of having “silicon interoperability” as soon as possible for DOCSIS “certified” CMs and CMTSs. CableLabs runs a variety of test and certification laboratories in their facility. Numerous vendors are participating. It was the expectation to have many cable modem vendors certified by the cable industry major trade show, the Western Cable Show, in December, 1998. However, as interoperability does take time to work out, the process is taking longer than expected. There will likely be some certified vendors by December 1998, with many more in first quarter 1999. It is now expected that the first widespread deployments of DOCSIS cable modems will start in late first quarter 1999.

The DOCSIS project is currently updating the RFI Version 1.0 specification to include better support for bandwidth management and QoS support. The changes being studied include support for multiple *Service Identifiers* (SIDs), filters to perform the classification of IP packets to different SIDs for differentiated services (QoS), and the signaling support for dynamic SID creations and deletion. A scheme for packet fragmentation will be included which will give substantially better support for managing jitter for delay sensitive traffic, such as packet voice. The primary motivation for adding these extensions to DOCSIS RFI V1.0 is to provide for better support of packet voice and video over DOCSIS IP services. A major focus of the North American cable industry is to support “near toll quality” voice and video services via DOCSIS systems. The cable industry effort writing specification for packet voice and video is called *PacketCable*^[18]. It is expected that the DOCSIS RFI V1.1 and initial PacketCable specifications will appear in December 1998.

DOCSIS RFI Version 1.0 was adopted by the *Society of Cable Television Engineers* (SCTE) Data Standards Subcommittee in July 1997 as the North American residential cable modem system standard.

Substantial work is in progress in the IETF *IP over Cable Data Networks* (ipcdn) working group to standardize the DOCSIS MIBs^[19, 20] and to standardize IP over DOCSIS^[21].

An IP over Cable Modem Example

This section presents a brief overview of a hypothetical IP over HFC system. It is meant to be an informative example to illustrate the application of the IP technology and some of the issues that surround provision of the service over a residential cable TV network. Moving IP datagrams in and out of the home over the cable plant is the important issue. The specific technology and protocols used by the cable modem vendor are important only in their ability to provide required IP service support.

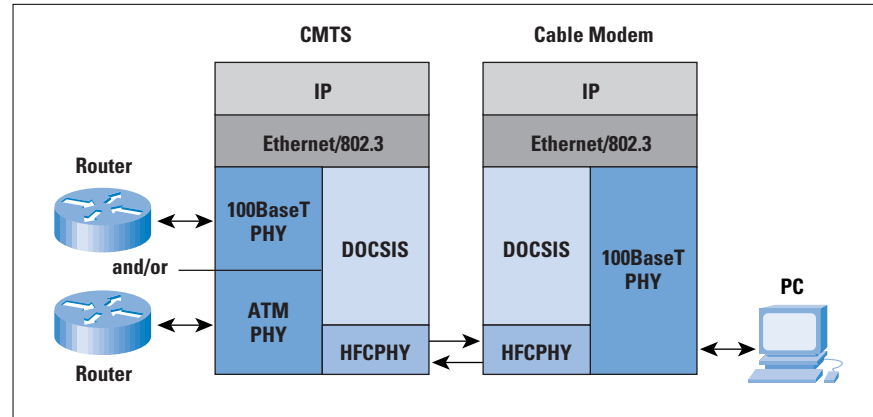
For this example, consider a system that has the following design goals and requirements:

- One-to-many service will be supported in the downstream direction; that is, many cable modems are reachable via the downstream channel
- Many-to-one service will be supported in the upstream direction; that is, the upstream channel bandwidth will be shared. There may be up to several upstream channels
- The protocol used between the Head-End Controller and the head-ends is not significant as long as it meets the needs of the IP service
- The head-end owns the upstream bandwidth and allocates resources to cable modems
- IP over Ethernet 10BaseT is the required interface in the home
- IP over Ethernet or IP over ATM is the required interface at the head-end

This example will rely on the DOCSIS RFI information presented previously in this article. The CMTS can transmit packets to any cable modem on the channel in any order or rate appropriate to the scheduling information it has and controls. The CMTS also participates in the IP multicast group membership (*Internet Group Management Protocol* [IGMP]) and *IP Resource Reservation Protocol* (RVSP) and makes changes in the cable modem resource assignments and allocations as needed. The home cable modem is permitted to use only the upstream channel under direction of the CMTS. Guaranteed and best-effort bandwidth allocations are dynamically assignable by the CMTS. It is assumed that the cable modem protocol has a bandwidth request facility that allows a CM to ask the CMTS for bandwidth. The function of the bandwidth management process is to sort these requests for service and give fair access to the requesting cable modems.

The method for implementation of an Ethernet and 802.3 bridging function over DOCSIS essentially permits the RF channels to act as a serial connection between a half-bridge function in each cable modem with a master in the CMTS. Figure 5 illustrates the protocol stack for this solution. The system presents an Ethernet-like segment to the cable operator. It is well-known how to put together such segments to construct larger internetworks.

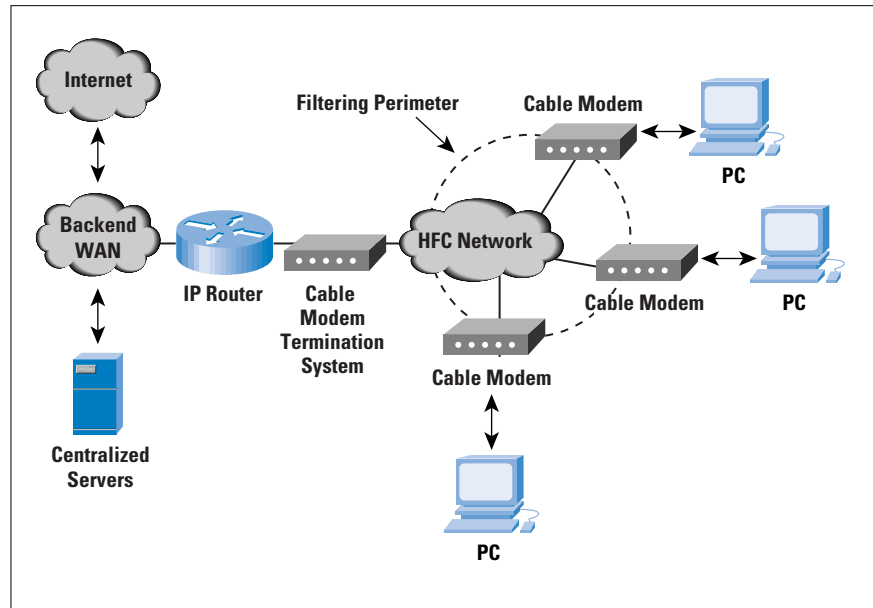
Figure 5:
Bridged Ethernet via
DOCSIS Example



Cable modems provide demarcation between the Internet Service Provider's network and each home network. To help the Internet Service Provider offer fair access service to its residential customers, the cable modem will require sufficient dynamic functionality for multilayer protocol filtering and various forms of rate management (see Figure 6). The goal of this filter is to create a defense perimeter at the first point of entry to the cable network; this perimeter will protect the upstream channel from being saturated or abused by misbehaving home networks. Some examples of this filtering functionality include, but are not limited to:

- Filtering on Ethertype for permitting only certain protocols to pass upstream; for example, IP and ARP only
- Filtering on IP source or destination address to permit/deny access from the home network
- IP and Ethernet broadcast rate limiting; that is, keep any home network broadcast storms confined to the home network
- IP Multicast group address filtering; that is, explicitly permit participation of the home network in an IP multicast group

Figure 6:
Internet Services via
Cable Modem
Deployment Model



It should be noted that these filtering functions are under consideration by numerous cable modem manufacturers, and they are being discussed in the IETF ipcdn working group.

A brief overview of IP over cable TV networks has been presented. From an engineering and deployment viewpoint, making the Internet move over cable modems is deceptively straightforward. Many issues are beyond the scope of this article: address allocation methods, back-end network design, configuration services, server placement, home customer support services, installation, firewalls, and troubleshooting.

Summary

This article has presented an overview of the work in progress of the IEEE 802.14 Cable TV MAC and PHY Protocol Standards working group and the MCNS DOCSIS effort. Initial review of these works is positive; indications are that data over HFC systems are viable. The IEEE 802.14 effort began as a study group in late 1993 and has yet to produce a standard. The MCNS DOCSIS process started in early 1996, moved rapidly, and has produced an accepted international standard specification for North American cable operators for residential cable modem service. The IEEE 802.14 standard appears to be destined for some international use and in systems where ATM over CATV is preferred by cable operators.

The cable network environment will provide a very usable and scaleable bandwidth platform for delivering Internet services to and from the home^[22]. A hypothetical example was provided that illustrates a general equipment deployment model. Actual deployment of Internet to the home will occur in many areas of North America in 1998 with increasing and substantial deployment in 1999.

For More Information

Information on the IEEE's 802.14 working group can be found on the World Wide Web at: <http://www.walkingdog.com/>

Information the Internet Engineering Task Force's IP over Cable Data Networks working group can be found at: <http://www.ietf.org/>

Information on the North American MCNS DOCSIS effort can be found at: <http://www.cablemodem.com/>

Information on the North American PacketCable effort can be found at: <http://www.packetcable.com/>

Information on the SCTE Data Standards Subcommittee can be found at: http://www.cablenet.org/scte/scte_dcs.html

References

- [1] Baran, Paul, "On Distributed Communication Networks." *IEEE Transactions on Communication Systems*, Vol. CS-12, pp. 1-9, March 1964.
- [2] ATM Forum, "ATM User-Network Interface Signaling 4.0," Specification number af-sig-0061.000, www.atmforum.com, July, 1996.
- [3] MCNS, "Data-Over-Cable Service Interface Specification—Radio Frequency Interface." SP-RFI-102-981008, www.cablemodem.com, July, 1998.
- [4] MCNS, www.cablemodem.com, main page, April 1998.
- [5] Kim, Albert. "Two-Way Plant Characterization." Technical Session 23, National Cable Television Association Show and Conference, Dallas, Texas, May 9, 1995.
- [6] Chelehemal, M., Prodan, R., et al., "Field Evaluation of Reverse-Band Channel Impairments." Society of Cable Telecommunications Engineers, Emerging Technologies Conference, San Francisco, California, January 9-12, 1996.
- [7] Laubach, Mark, "Avoiding Gridlock on the Data Infobahn: Port Mismatches Pose Challenges." *CED Magazine*, March 1998
- [8] Abramson, Norman, "Development of the ALOHANET." *IEEE Transactions on Information Theory*, Vol. IT-31, pp. 119-123, March 1985.
- [9] XEROX, "The Ethernet, A Local Area Network: Data Link Layer and Physical Layer Specification." X3T51/80-50, Xerox Corporation, Stamford, Connecticut, October 1980.
- [10] IEEE, "Carrier Sense Multiple Access with Collision Detection (CSMA/CD) Access Method and Physical Layer Specifications." Standard 802.3-1985 (ISO DIS 8802/3), IEEE, New York, ISBN 0-471-82749-5, 1985.
- [11] IEEE, "IEEE Standards for Local Area Networks: Logical Link Control, ANSI/IEEE Std 802.2-1985." Fifth printing, February 1988.

- [12] IEEE 802.14 Working Group, "Cable-TV Functional Requirements and Evaluation Criteria." Work in progress, IEEE802.14/94-002R2, IEEE 802 Committee, February 1995.
- [13] Laubach, Mark. "Classical IP and ARP over ATM." RFC 1577, January 1994.
- [14] Laubach, Mark, "Logical IP Subnetworks over IEEE 802.14 Services." Work in progress, **draft-ietf-ipcdn-ipover-802d14-01.txt**, November 1997.
- [15] Laubach, Mark, "Serving Up Quality of Service." *CED Magazine*, April 1997.
- [16] Laubach, Mark, "Deploying ATM Residential Broadband Networks." NCTA Cable 96 Conference, Los Angeles, California, April 30, 1996.
- [17] Nichols, Kathleen, and Laubach, Mark, "On Quality of Service in an ATM-based HFC Architecture." IEEE ATM Workshop 96, San Francisco, California, August 27, 1996.
- [18] PacketCable, "What is PacketCable?" <http://www.packetcable.com>, April 1998.
- [19] Roeck, Guenter, "Cable Device Management Information Base for MCNS compliant Cable Modems and Cable Modem Termination Systems." Work in progress, **draft-ietf-ipcdn-cable-device-mib-05.txt**, October 1998.
- [20] Roeck, Guenter, "Radio Frequency (RF) Interface Management Information Base for MCNS compliant RF interfaces." Work in progress, **draft-ietf-ipcdn-rf-interface-mib-05.txt**, October 1998.
- [21] White, Gerry, "Logical IP Subnetworks over MCNS Data Link Services." Work in progress, **draft-ietf-ipcdn-ip-over-mcns-00.txt**, August 1997.
- [22] Lucien Rhodes, "The Race for More Bandwidth." (Interview with Milo Medin of @Home), *Wired Magazine*, Vol. 4.01, January 1996

Internet Drafts are *works in progress* and can be retrieved from:

ftp://ds.internic.net/internet-drafts

MARK LAUBACH holds a B.E.E. and M.Sc. from the University of Delaware. He is Vice President and Chief Technical Officer at Com21, Inc. in Milpitas, California, and is responsible for the end-to-end systems architecture and ATM over HFC protocol specification of the Com21 product family. Prior to Com21, he was with the Hewlett-Packard Company for 14.5 years. Laubach is a member of the IETF, and is past chair of the IP over ATM working group. He is the author of RFC 1577, "Classical IP and ARP over ATM." He regularly attends IETF, IEEE, and SCTE working group meetings. He is a Senior member of the IEEE and a member of the SCTE. E-mail: **laubach@com21.com**

Digital Television: A New Venue for the Internet

by George Abe, Cisco Systems

The digitization of television is of interest to the Internet community in that it opens the possibility of a new mode of delivering IP packets to the home. IP services can be delivered over television broadcast distribution networks, whether over the air, cable, or satellite. This article introduces the basic concepts of *digital television* (DTV) and provides a point of departure for further reading.

Why Is Digital TV Happening?

The original motivation for the research into advanced TV (we avoid the term DTV for a moment) was to prop up sagging TV sales. It was mostly vendor push.

By the late 1970s, Japan and Korea had achieved domination in the production of TV sets worldwide. They were so successful that the market had become saturated, particularly in the developed world. Everyone had one or, more likely, three or four TVs at home. Further, a TV lasts over 10 years, so the replacement market is low. TV production had ceased to be a growth market. Margins were and are poor and few innovations were on the horizon.

So in the early 1980s Japan had begun research into new high-definition televisions that would stimulate new demand and enable them to keep their market leadership. Their system is called *Multiple Subnyquist* (MUSE). MUSE was an analog system, but it had better-quality pictures.

Not to be outdone, the U.S. decided it needed to try to recapture the TV market, so began its own development, under the aegis of the Federal Government. A partnership called the *Grand Alliance* was formed, and it began working in 1984. Pioneering work was done by the partnership members, particularly Zenith, MIT, and General Instruments. They created a digital specification after more than a decade of research and development. Along the way, the computer industry made contributions (or some would say interferences) of its own until the FCC announced a final specification in December 1996. The basic elements are found at www.atsc.org and referenced later in this article.

Benefits of DTV

The movement toward widespread DTV gained momentum among government officials, broadcasters, and hardware vendors when some of the benefits became clear.

First, because of improvements in technology, it is possible to transmit pictures and sound of significantly higher quality in the same 6 MHz spectrum that analog TV occupies. The 6 MHz spectrum is wasteful of bandwidth, and the government would like to recover the excess so it can be auctioned or used to support other public services (police, fire,

deep space probes, and so on), which could operate at the relatively low frequencies of VHF TV.

Second, digitally encoded TV could provide new services, such as Web access via TV or interactive TV. These have long been dreams of the consumer electronics (CE) industry, but hope springs eternal.

Third, digital TV offers greater security to the programmer and the network. There is a cottage industry in hacking analog set-top boxes. Digital techniques, such as the *Data Encryption Standard* (DES), double DES, and triple DES give operators hope that they can secure their pay-per-view content.

Finally and most interestingly, since digital TV occupies less bandwidth per program, broadcasters, satellite operators, and cable operators have the opportunity to offer more channels. Instead of a mere 10–13 channels available over the air in a single metropolitan area, it is possible to have perhaps 60 or more over the air channels. Cable operators, with their greater bandwidth underground, could have many more channels. Although technically cable could offer 500 channels, it is hard to imagine where the scripts would come from.

What Is DTV?

By our definition, digital television is the capture, production, distribution, and broadcast of programming in a digitally encoded format. Whereas today's analog TV transmits in amplitude modulation, DTV would use *Quadrature Phase Shift Keying* (QPSK), *Quadrature Amplitude Modulation* (QAM), or *Vestigial Side Band* (VSB) modulation techniques. We won't detail these techniques here except to mention that they are mutually incompatible.

When DTV standards were discussed in the 1980s, the industry could not agree on a single display. The deliberations became more protracted with the entry of the computer industry into the discussions, long after the broadcasters and consumer electronics people began their work. Would there be interlaced or progressive scanning? Would there be the existing aspect ratio or would there be a wide-screen display? Square pixels or not? How many lines of resolution would be displayed?

With the broadcasters and consumer electronics vendors arguing for interlacing, oval pixels, and wide screens and the computer people arguing for progressive scanning, square pixels, and a more square display, the disagreements could not be bridged.

Therefore, the FCC had no choice but to declare that the “market should decide” which display format would prevail. Accordingly, the FCC announced in December 1996 that 18 different display formats would be permissible for over-the-air digital TV. A broadcaster could elect to transmit in any of the approved formats. The approved formats are shown in Tables 1 and 2.

Table 1: Progressive Video Scanning Formats for Digital TV

Vertical Lines	Horizontal Pixels	Aspect Ratio	Frame Rate per Second
1080	1920	16:9	24, 30
720	1280	16:9	24, 30, 60
480	704	16:9	24, 30, 60
480	704	4:3	24, 30, 60
480	640	4:3	24, 30, 60

Table 2: Interlaced Video Scanning Formats for Digital TV

Vertical Lines	Horizontal Pixels	Aspect Ratio	Frame Rate per Second
1080	1920	16:9	30
480	704	16:9	30
480	704	4:3	30
480	640	4:3	30

The vernacular to describe the formats typically indicates the number of vertical lines and the scanning format. For example, “1080i” refers to 1080 lines, interlaced scanning; “720p” refers to 720 lines in progressive format.

In practice, only a few of the 18 approved formats are under consideration by the nation’s broadcasters. NBC and CBS have declared they will support 1080i. ABC is opting for 720p, and Fox has opted for 480p.

Apart from the controversy over display, most of the other elements were quickly resolved. Modulation scheme, transport multiplexing, compression, timing, and an overall systems and testing procedure were agreed to. The apparatus for DTV was in place, almost. The time was January 1997.

High Definition or Standard Definition

Some view DTV as synonymous with high-definition television. It is not. DTV encompasses both *High-Definition TV* (HDTV) and *Standard-Definition TV* (SDTV). Hence HDTV is a proper subset of DTV. The difference between HD and SDTV is not standardized, but our definition of HD includes the display formats that have 720 or 1080 lines. Formats with fewer lines are standard definition.

The key point of difference between HD and SD is that with HD and current compression techniques (MPEG-2), only one program is accommodated in one 6-MHz channel. With SD, it is possible for the broadcaster to transmit two or more programs simultaneously, in a single 6-MHz chunk of bandwidth.

This has tremendous implications. If broadcasters can transmit multiple channels at once, it would be possible (technically) for Disney to broadcast ABC, the Disney Channel, ESPN, and A&E over the air in the same bandwidth they use to show ABC today. (Of course they won't do this for commercial and contractual reasons, but the technology makes it doable).

For Internet Service Providers, a broadcast could transmit SD programming simultaneously with datacasting, and go into the push-mode data service business. For example, Disney/ABC could download software updates for Disney Interactive, or perhaps contract with Microsoft to deliver Windows updates. Whereas most Internet folk view MPEG being transported inside IP packets on the Internet, broadcasters intend to insert IP packets into MPEG-2 transport streams. The consumer's digital set-top box would tune to the data "channel," extract the data from its MPEG capsule, and divert the data packet to an Ethernet or ATM port on the set-top.

There are nearly 1,600 broadcasters in the U.S. Each could, in theory, transmit 19.3 megabits per second. Of course, most of these bits will be used for television, but certainly 1 or 2 megabits can be accommodated by each broadcaster for data service.

Given the dearth of programming to fill multiple SD channels, broadcasters are strongly motivated to consider data services and compete for a slice of the Internet service market.

Digital TV—End to End

Whereas one easily thinks of DTV as a distribution and display technology, in fact there are major changes required to capture, edit, and distribute digital content. Thus there is the need for new cameras, post-production editors, sound mixers, and the like.

Digital TV can be transmitted over the air, through cable networks, or via *Direct Broadcast Satellite* (DBS). Today, only DBS has achieved large-scale distribution of digital TV, with over 7 million subscribers in the U.S. and 15 million worldwide.

Content is created either through a digital camera or by converting existing analog content, such as 35mm film, into digital format. Within the production environment, editing changes are made, typically using *Nonlinear Editors* (NLEs) that connect to a local-area network.

Original production is normally done in the high definition. The highest form of resolution is 1.492 Gbps. (See Table 3.) Equipment to do this is not widely available, but it will be eventually. Panasonic is shipping a digital camera capable of 1.5-Gbps output, but rumor has it they cost almost \$500,000, if you can even get one. Nonetheless, 41 stations began HD programming in November, highlighted by an NFL game on CBS between the Buffalo Bills and the New York Jets on November 8.

Some compression is applied within the postproduction and editing environment. The TV industry, through the *Society of Motion Picture and TV Engineers* (www.smpte.org), developed a series of digital transmission standards. Chief among these is SMPTE 305M, which defines a protocol called *Serial Data Transport Interface* (SDTI), which calls for a 270- or 360-Mbps service to link various pieces of production equipment such as NLEs in a postproduction facility. SMPTE 305M is a networking scheme complete with an addressing specification.

(Interesting point about 305M: It is the first and only protocol known to this author that specifies use of IPv6 addressing.)

Another important protocol is SMPTE 259M, which is a link-layer protocol underneath 305M.

A competing protocol to SDTI is the *Digital Video Broadcasters Asynchronous Serial Interface* (DVB-ASI). Information on DVB-ASI is found at www.dvb.org.

From the editing environment, content is distributed via satellite or land lines to local affiliates (for local over-the-air broadcast), cable head-ends (for cable TV distribution) and satellite hubs (for direct-to-home satellite service). The distribution from national feeds to local facilities is normally at T3/E3 speeds because of the availability of T3/E3 services by telephone companies and satellite transponders for affiliate and direct-to-home distribution.

Cable providers, local broadcasters, and satellite services add their own content and make certain changes to the national feeds. Among these changes are assignment of the programming to specific frequencies or channels, insertion of local advertising, local programming, and emergency broadcasts.

After adding their own content, the local services distribute the final programming to consumers. Over-the-air broadcasters will transmit 19.3 Mbps per 6 MHz, cable will transmit 27 Mbps per 6 MHz, and satellite uses variable channelization, kept closely under wraps.

So there is the progression downward from 1492 Mbps of original encoding, to 270 Mbps for editing, to 34/45 Mbps for affiliate distribution, to 27 Mbps or less for distribution to the end user.

Table 3: Bit Rate Requirements for Various Display Formats

Format	Pixels per Line	Lines per Frame	Pixels per Frame	Frames per Second	Millions of Pixels per Second	Bits per Pixel	Mbps
SVGA	800	600	480,000	72	34.6	8	276.5
NTSC	640	480	307,200	30	9.2	24	221.2
PAL	580	575	333,500	50	16.7	24	400.2
SECAM	580	575	333,500	50	16.7	24	400.2
HDTV	1920	1080	2,073,600	30	62.2	24	1492.8
Film	2000	1700	3,400,000	24	81.6	32	2611.2

Note: Film display formats vary, depending on content and directorial prerogative.

Over the Air and Cable

All the huffing and puffing by the FCC, the consumer electronics industry, the computer industry, and the broadcasters pertains to over-the-air transmission. However, about two-thirds of the American viewing public views TV through cable. So if most Americans are to receive DTV, they must receive it through cable.

This raises important technical and regulatory questions. The technical question is: How are the digital signals produced by the broadcasters and their affiliates to be sent through wires, and what is the allocation of functions between the digital set-top and the digital receiver? This question seems simple but it is not, as we shall see.

The regulatory question pertains to whether the cable operators are to be compelled to carry DTV from broadcasters. This problem is referred to as the digital *Must Carry Problem*, now under consideration by the FCC. It certainly will be litigated, whatever the outcome of the FCC's decision.

Technical Question

Among the key provisions agreed to by the Grand Alliance is the use of a modulation technique called 8-VSB for over-the-air digital transmission. The particulars of 8-VSB are not significant here, but we will mention that this particular decision was arrived at in the mid-1980s, before the cable industry had much impact on the viewing public or on the broadcasting industry.

When the cable industry began to think about digital, in the mid-1990s, they settled on a modulation scheme called 64 QAM. 64 QAM is able to produce 27 Mbps in 6 MHz, whereas 8-VSB produces about 19.3 Mbps. The difference occurs because over-the-air broadcasting requires a more robust encoding scheme to combat the more hostile nature of over-the-air transmission, as opposed to the safer environment of coaxial cables. Thus the cable modulation technique can be more aggressive than over-the-air techniques.

(We should add that satellites use an even more robust modulation technique called QPSK, which gets fewer bits per Hertz than VSB or QAM. But robustness is needed because satellite signals must travel far greater distances than cable or local broadcast.)

Thus for cable to carry a digital over-the-air broadcast, some conversion of 8-VSB encoding to 64 QAM encoding is necessary. This necessity does not present a major technical problem, but agreement is needed on where the conversion is done and at what cost. For example, Broadcom and Sony are collaborating on the development of a chip, to be embedded in a TV, that can decode VSB and QAM. It sounds simple, but the cable industry is not interested. They want to carry QAM and QAM only on their networks.

One option is to convert the format of the digital bitstream coming out of the cable box to the IEEE 1394 *FireWire* format. Since DTVs are likely to have FireWire input, this conversion can provide a ubiquitous connection. However, this scenario raises the problem of copy protection, a sore point in Hollywood. Since digital copies are pristine, the content providers (studios and record companies) are firm in their resolve that unless there is strong copy protection, none of their content will be available over FireWire.

Another option is to build a set-top box that takes baseband signals and modulates them to look like 8-VSB broadcast signals on channel 3, similar to how VCRs work in the analog world now. This scenario is clearly rather ugly, but understood by consumers.

Finally, it could be up to the cable operators to transmodulate the 8-VSB into QAM at the cable head-end. Better yet, they can accept broadcasters' feeds in baseband, and then QAM-modulate the baseband signals for their consumers. The cable set-top box would be sending bit maps to a dumb digital monitor, like a computer monitor, which doesn't know or care that it is receiving QAM or VSB programming.

Apart from modulation, there is the issue of display format. NBC and CBS have declared they will transmit in 1080i. ABC has chosen 720p and Fox has chosen 480p, with some vague pledge for higher definition later. After all, it does not seem necessary to show *The Simpsons* in HD.

On the other hand, John Malone, Chairman of TCI, went public in May 1998 with his declaration that TCI would not voluntarily carry 1080i because it (1080i) was wasteful of bandwidth. Implied in his comment is the fact that cable operators do need to be restricted to 6-MHz channelization for digital. In fact, the entire DTV spectrum on cable could be considered a gigantic pool of bandwidth that the cable operator could allocate to individual channels, much as direct satellite does. This setup gives the cable operators incentive to downconvert the broadcasters' DTV signals. For example, when NBC sends 1080i, the cable operator may elect to transmit 720p, or less, to its customers.

Should the cable operators be required to carry the HDTV pictures from the broadcasters in the broadcasters' chosen format? Would they be allowed to downconvert the HD into standard definition? What happens when a broadcaster, say NBC, elects to transmit in SDTV and thereby has the capability of multiplexing several channels onto a single chunk of 6 MHz? What is the duty of the cable operator to carry Internet datacasting offered by the broadcasters over the cable network, in competition with services such as @Home and Roadrunner?

The complexities of multiplexing go further. Let's say ABC elects to broadcast SD. If one of the subprograms in the multiplex is a pay-per-view channel, should the authentication procedures of the cable operator be superceded? Should the electronic program guide of the cable operator be superceded?

Questions like these have technical and regulatory aspects and are being worked in industry, the FCC, and state regulatory agencies. It is possible that Congress will get involved as well. When John Malone made his statement, both sides of the aisle in Congress were not amused. They want DTV to happen so that spectrum can be freed. If the cable operators stand in the way, the conversion to digital is stopped dead in its tracks.

The Open Cable Initiative

The cable industry does not want to be a bottleneck to broadcasters. On the other hand, it needs to make quick progress into DTV to compete against satellite. Therefore, the industry has embarked on a process called *Open Cable*, which seeks to define a digital set-top box that can be available at retail. Available at retail means a nonproprietary, open design. Open Cable strives to make the DTV set-top box independent of processor platform (that is, not an Intel Pentium necessarily) and operating system independent (that is, not a Microsoft Windows CE necessarily).

The Open Cable set-top box will allow for data services through a specification written by the *Digital Audio Visual Council* (DAVIC—www.davic.org) and therefore, is not compatible with the current *Data-over-Cable Service Interface* (DOCSIS) specification supported by the U.S. cable industry. (See article starting on page 13.) However, it is possible for DOCSIS capabilities to be added on to an Open Cable set-top box. We mention Open Cable because it will be the key customer premises device for cable and digital TV and much hinges on its interoperability with broadcasters transmissions.

Digital TV via Satellite

In addition to over-the-air and cable, DTV can be received by satellite. As of this writing, it is the only way to receive DTV. The digital satellite industry has nearly 7 million subscribers who received DTV today. Its role in all the discussions of HD vs. SD and the provision of data services is relatively low key because it is believed that satellite will continue to be a niche provider because of its technical and legal problems in distributing locally originated TV stations.

But satellites bear watching because if they are able to deliver local channels and obtain 15–20 million homes in the U.S., then the financial consequences on cable and over the air could be crucial.

The New Digital Studio

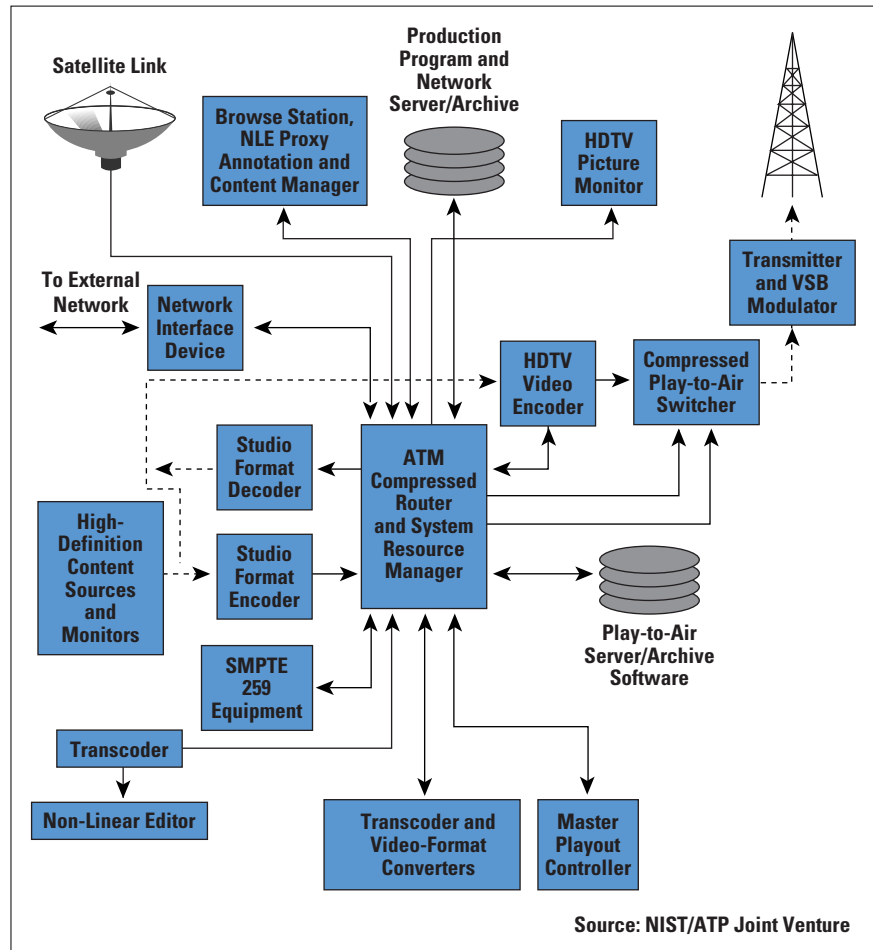
The figure shows a schematic of the elements of a DTV broadcast studio described recently by the U.S. *National Institute of Standard and Technology* (NIST). At the heart of the studio is an ATM switch with new interfaces that connect to DVB or ATSC infrastructures via DVB-ASI or SDTI interfaces.

Connection for wide-area distribution will likely be over ATM. Converters exist for DVB-ASI to ATM. For example, Cellware (www.cellware.de) in Germany markets such a converter, but there is no SDTI-to-ATM interface known to this author at this time.

The digital studio provides a new a marketing opportunity for the LAN industry. Broadcast digital production demands higher speeds than most other LAN applications.

Thus vendors of data communications equipment have two opportunities: to provide equipment to broadcasters who want to enter the Internet service business and to production houses that use ATM or other LANs to support editing and production applications.

Figure 1:
Prototype of HDTV
Broadcast Studio



Web Sites

www.atsc.org: *Advanced TV Standards Committee.* S13 and S16 are subgroups working on datacasting; S13 focuses primarily on the downstream path, whereas S16 focuses primarily on the reverse communication from the receiver. Since over-the-air is one way, this work is limited to the communications between the S13 forward channels and a telephone or Internet return path.

www.dvb.org: The *Digital Video Broadcasting Project (DVB)* has taken the lead in defining DTV specifications as well as defining datacasting interfaces over DTV infrastructures.

www.smpte.org: *Society of Motion Picture and Television Engineers.*

www.sbe.org: *Society of Broadcast Engineers.*

www.scte.org: *Society of Cable TV Engineers.*

www.mpeg.org: *Motion Picture Experts Group.* The word on MPEG compression, controls, and transmission.

References

- [1] ISO/IEC IS 13818-1, International Standard, MPEG-2 Systems.
- [2] ISO/IEC IS 13818-2, International Standard, MPEG-2 Video.
- [3] ISO/IEC 13818-6, International Standard, Digital Storage Media Command and Control (DSM-CC).
- [4] ATSC Standard A/52 (1995), Digital Audio Compression (AC-3).
- [5] ATSC Standard A/53 (1995), ATSC Digital Television Standard.
- [6] ATSC Standard A/55 (1996), Program Guide for Digital Television.
- [7] ATSC Standard A/56 (1996), System Information for Digital Television.
- [8] ATSC Standard A/57 (1996), Program/Episode/Version Identification.
- [9] ATSC Standard A/63 (1997), Standard for coding 25/50-Hz Video.
- [10] ATSC Standard A/64 (1997), Transmission Measurement and Compliance For Digital Television.
- [11] ATSC Standard A/65 (1998), Program and System Information Protocol for Terrestrial Broadcast and Cable.
- [12] ATSC T3/S13 Doc. 010 DVS-yyy Rev z Draft, ATSC Data Broadcast Specification for Terrestrial Broadcast and Cable.
- [13] ETR XXX: Digital Video Broadcasting (DVB); Guidelines for the Use of the DVB Specification: Network Independent Protocols for Interactive Services (ETS 300 802).
- [14] SCTE DVS-nn: SCTE Digital Video Subcommittee (DVS) standard for Cable Headend and Distribution Systems (spec not released—under development)

GEORGE ABE holds an A.B. in Mathematics and an M.S. in Operations Research from UCLA. He currently is a Consulting Engineer at Cisco Systems, where he has dabbled in various areas of residential broadband networking since 1994. He is the author of *Residential Broadband*, Cisco Press (imprint of Macmillan Press). He expects to be a early adopter of digital TV and, when not watching TV, he can be reached at georgea@acm.org

I Remember IANA

by Vint Cerf, MCI WorldCom
October 17, 1998



Photo: Chris Pizzello, New York Times Pictures

A long time ago, in a network, far far away, a great adventure took place! Out of the chaos of new ideas for communication, the experiments, the tentative designs, and crucible of testing, there emerged a cornucopia of networks. Beginning with the ARPANET, an endless stream of networks evolved, and ultimately were interlinked to become the Internet. Someone had to keep track of all the protocols, the identifiers, networks and addresses and ultimately the names of all the things in the networked universe. And someone had to keep track of all the information that erupted with volcanic force from the intensity of the debates and discussions and endless invention that has continued unabated for 30 years. That someone was Jonathan B. Postel, our *Internet Assigned Numbers Authority* (IANA), friend, engineer, confidant, leader, icon, and now, first of the giants to depart from our midst.

Jon, our beloved IANA, is gone. Even as I write these words I cannot quite grasp this stark fact. We had almost lost him once before in 1991. Surely we knew he was at risk as are we all. But he had been our rock, the foundation on which our every Web search and e-mail was built, always there to mediate the random dispute, to remind us when our documentation did not do justice to its subject, to make difficult decisions with apparent ease, and to consult when careful consideration was needed. We will survive our loss and we will remember. He has left a monumental legacy for all Internauts to contemplate. Steadfast service for decades, moving when others seemed paralyzed, always finding the right course in a complex minefield of technical and sometimes political obstacles.

Jon and I went to the same high school, Van Nuys High, in the San Fernando Valley north of Los Angeles. But we were in different classes and I really didn't know him then. Our real meeting came at UCLA when we became a part of a group of graduate students working for Professor Leonard Kleinrock on the ARPANET project. Steve Crocker was another of the Van Nuys crowd who was part of the team and led the development of the first host-to-host protocols for the ARPANET. When Steve invented the idea of the *Request for Comments* (RFC) series, Jon became the instant editor. When we needed to keep track of all the hosts and protocol identifiers, Jon volunteered to be the Numbers Czar and later the IANA once the Internet was in place. Jon was a founding member of the *Internet Architecture Board* (IAB) and served continuously from its founding to the present. He was the *first* individual member of the Internet Society—I know, because he and Steve Wolff raced to see who could fill out the application forms and make payment first and Jon won. He served as a trustee of the Internet Society.

He was the custodian of the .us domain, a founder of the Los Nettos Internet service, and, by the way, managed the networking research division of USC Information Sciences Institute.

Jon loved the outdoors. I know he used to enjoy backpacking in the high Sierras around Yosemite. Bearded and sandaled, Jon was our resident hippie-patriarch at UCLA. He was a private person but fully capable of engaging photon torpedoes and going to battle stations in a good engineering argument. And he could be stubborn beyond all expectation. He could have outwaited the Sphinx in a staring contest, I think.

Jon inspired loyalty and steadfast devotion among his friends and his colleagues. For me, he personified the words “selfless service.” For nearly 30 years, Jon has served us all, taken little in return, indeed sometimes receiving abuse when he should have received our deepest appreciation. It was particularly gratifying at the last Internet Society meeting in Geneva to see Jon receive the Silver Medal of the International Telecommunications Union. It is an award generally reserved for Heads of State, but I can think of no one more deserving of global recognition for his contributions.

While it seems almost impossible to avoid feeling an enormous sense of loss, as if a yawning gap in our networked universe had opened up and swallowed our friend, I must tell you that I am comforted as I contemplate what Jon has wrought. He leaves a legacy of edited documents that tell our collective Internet story, including not only the technical but also the poetic and whimsical as well. He completed the incorporation of a successor to his service as IANA and leaves a lasting legacy of service to the community in that role. His memory is rich and vibrant and will not fade from our collective consciousness. “What would Jon have done?” we will think, as we wrestle in the days ahead with the problems Jon kept so well tamed for so many years.

There will almost surely be many memorials to Jon’s monumental service to the Internet Community. As current chairman of the Internet Society, I pledge to establish an award in Jon’s name to recognize long-standing service to the community, the *Jonathan B. Postel Service Award*, which will be awarded to Jon posthumously as its first recipient.

If Jon were here, I am sure he would urge us not to mourn his passing but to celebrate his life and his contributions. He would remind us that there is still much work to be done and that we now have the responsibility and the opportunity to do our part. I doubt that anyone could possibly duplicate his record, but it stands as a measure of one man’s astonishing contribution to a community he knew and loved.

VINTON G. CERF is senior vice president of Internet Architecture and Technology for MCI WorldCom. Widely known as a “Father of the Internet,” he is the co-designer of the TCP/IP protocol. Cerf served as founding president of the Internet Society from 1992–1995 and is currently chairman of the Board. Cerf holds a Bachelor of Science degree in Mathematics from Stanford University and Master of Science and Ph.D. degrees in Computer Science from UCLA. E-mail: vc erf@mci.net

Book Reviews

Internet Messaging *Internet Messaging: From the Desktop to the Enterprise*, by Marshall T. Rose and David Strom ISBN 0-13-978610-4, Prentice-Hall PTR, 1998, <http://www.prenhall.com>

Very few Internet voices hold a status equivalent to E.F. Hutton's advertising campaign: "When they speak, we should listen." Marshall Rose and David Strom are two such voices, making any product of their combined efforts a serious matter, indeed. Rose has typically written about basic technology, Strom about the pragmatics of use, especially trials and tribulations of fitting networked pieces together. *Internet Messaging* is in the latter category, with a strong added introduction of e-mail and security technology. Anyone who has professional contact with e-mail should get a copy of this book. If commercial use of Internet mail were more advanced and stable, we probably would not need an effort like this. However, e-mail professionals must constantly deal with problems in using interesting functions and in troubleshooting interoperability. *Internet Messaging* helps with the planning, use and debugging of complex, or otherwise "interesting," e-mail services.

Updated Information

The book provides a superb survey of the relevant technology, the popular user mail software, and the rather interesting range of mail and messaging operations issues, including styles of use by organizations. The comparisons of different mail systems leave the reader with a solid understanding of functional and usage requirements for modern systems, as well as the choices available at the time of publication. Mary Houten-Kemp's Web site at <http://www.everythingemail.net> is being used to provide updated information.

E-mail includes a wide range of technical and operations issues, and *Internet Messaging* touches all of them. Its introductions cover user environment, mail transfer, mailing list services, unsolicited bulk e-mail ("spam"), encryption-based security, remote user access, virtual private networks, and directory services. Providing a single discussion, which integrates the use of these disparate technologies, is enough to justify the book.

Organization

Internet Messaging attempts very regular organization and states that the goal is to permit use as a problem/solution reference work. It primarily distinguishes between sending and receiving functions and between desktop and enterprise requirements. This creates a two-by-two matrix, defining the core four chapters. The other chapters include philosophical opening and closing discussions, a separate, very informative chapter on security, and another on general enterprise operations issues.

Most of the chapters are organized into Introduction, Problems, Standards, and Solutions. Unfortunately that regularization is all that is shown in the Table of Contents, so the reader gets little help finding specifics by reading the Table. Similarly, the organization of the chapter contents did not seem compelling for use in problem solving. The additional “How Can I” matrix (on page 10) and its associated discussion text is intended as the primary means for locating relevant discussions.

Comparisons

User software comparisons are given throughout the book, for Microsoft Outlook 4.01, Netscape Messenger 4.04, Qualcomm Eudora Pro 4.0, Lotus cc:Mail 8.1, CompuServe WinCIM 3.02, and America Online 3.0. Specific mailing lists, security, remote access, and directory software and services are also reviewed. Oddly, the discussion of remote access mentions only global, single-provider services—and their favorite is currently having financial problems—but did not mention the “association” style of service that integrates many independent providers, notably GRIC and iPass. (Full disclosure: iPass is a client.)

Most products are undergoing aggressive enhancement so that no printed text can be entirely up-to-date. Hence the Web site. For the software and services I know well, the book looked reasonable. Of course it is not entirely error free, but the errors are small and perfect detail is not required. I believe there are two major benefits to these comparisons. One is that the reader is given a very solid sense of the general capabilities and limitations of modern e-mail software. The second is to make a reasonable, first-pass filtering of candidate packages to be used in an organization. It would *not* be appropriate to attempt selecting among these packages according to subtle differences reported in the book.

Benefits

As one would expect of these authors, a very large, long-term benefit of their efforts is in their many excellent criticisms and suggestions. Unfortunately, many of them are in notes located at the end of each chapter. It’s hard to imagine a less-convenient place to put them, since I found myself constantly shifting back and forth between the main text and the notes. It would not have been so irritating if the comments were less interesting; they should have been true footnotes, with easy access on each page. The stellar example of direct utility from these comments is Figure 2.1 on page 38. It shows a systems structure for user software processing of incoming mail. Every vendor should study this discussion carefully and implement it immediately. Please!

—Dave Crocker
Brandenburg Consulting
dcrocker@brandenburg.com

Web Security *Web Security: A Step-by-Step Reference Guide*, by Lincoln D. Stein, ISBN 0-201-63489-9, Addison-Wesley, December 1997,
<http://www.awl.com/cseng/titles/0-201-63489-9>

Whenever the topic of the World Wide Web comes up, you can be sure that some mention of “security” will soon follow. Web users, Web creators, and even Web technology developers are all keenly aware of the security concerns. But what do we mean by “security?” The safety to use a credit card? Keeping a Web site safe from break-ins? Keeping the kids away from online erotica? And whose security are we concerned with, the user’s or the Web site operator’s?

This book covers most of what we might expect to find under the umbrella of security. In addition to dealing with the broad scope of Web security, the author also tries to cover the topic with sufficient simplicity for the novice and enough detail for the engineer. The good news is that this book succeeds in delivering a single volume that covers all we could possibly expect on the topic, and at levels suited for a broad audience range.

Organization

The author begins by making the distinction between security for the browser, the Web site, and the network between them. This division of the topic forms the basis for the organization of the book. Moving through each of the three parts, the author proceeds from the simple to the complex in a logical, additive order. He discusses topics introduced early in the book from a functional standpoint—how they affect the user. He may cover the same technology in later chapters, but in greater depth, detailing server and network configuration and discussing the underlying technology.

In the first part of the book, the author covers document confidentiality, including standard “text” documents as well as electronic commerce. A major theme in this section is cryptography. The author presents symmetric and public key encryption technologies from a functional standpoint. He presents various encryption standards, with a discussion of their strengths and weaknesses. In another chapter he provides a good primer on the *Secure Electronic Transaction* (SET) protocol handling, as well as other options (*Common Gateway Interface* [CGI] scripts and *Secure Sockets Layer* [SSL]) for credit card order processing.

In Part 2 we are introduced to issues of client-side security. The author devotes a full chapter to an in-depth explanation of SSL services. He also looks at issues associated with active content, and presents technologies such as Java, ActiveX, and other options, along with notes on their respective security implications. Finally, he covers issues of privacy—in this case, the personal privacy of the user. Throughout these chapters, the author emphasizes user-controllable settings such as browser configuration options.

Whereas the author focuses on user involvement in the first two parts, with an appropriate level of technical content, in part 3, targeted to Web masters and system administrators, he introduces the engineering side with an in-depth coverage of server-side security. He covers the two prominent Web-serving operating systems: UNIX and Windows NT, with good attention to various versions of each. Topics include basic system security, access control, and activity monitoring. Other chapters include an excellent discussion of encryption and certificate technology, safe CGI scripting, remote authoring of Web data, and firewalls.

Presentation and Style

The author illustrates his points with good examples. He also presents appropriate sidebar discussions and illustrations, which not only clarify the information, but also provide interest and variety in what could be a very dry volume. Each chapter ends with a listing of resources, both print and “online.” Where appropriate, the author includes checklists to help the reader apply the material just covered.

As a result of the practical, well-grounded presentation of material, we are continually able to see practical applicability to our own situation. For example, the author presents us with information about dangers to our privacy, and why that might be important to us. This is immediately followed by clear instruction on changing privacy-affecting settings in various versions of both Netscape and Internet Explorer. The author uses this technique throughout the book, and it is as useful with password management, CGI scripting, or firewall configuration as it is with privacy.

Recommended

Although experts in encryption and other specific security-related technologies will find this book too simple for their personal area of expertise, the strength of the book is not in its coverage of any one area, but in its well-integrated and cohesive coverage of a broad range of interrelated topics. The ability for any reader, first-time surfer or Web guru, to find practical, easily applied information makes this book a required item on any webmaster’s bookshelf, and a must-read for anyone who spends any serious time on the Web.

—Richard Perlman
Berkeley Internet Group
perl@berkinet.com

Internet Cryptography *Internet Cryptography*, by Richard E. Smith, ISBN 0-201-92480-3, Addison-Wesley, 1998, www.awl.com/cseng/titles/0-201-92480-3

The 1990s might easily be known as the decade of the Internet. The Internet came into the mainstream during this decade, a global frontier with frontier problems and rules. Seemingly overnight, everyone from government agencies to Chinese restaurants had a Web presence. Young children exchanged e-mail with their grandparents and friends, a big change from just a few years ago when it was the domain of technologies and a place where everybody knew your name.

The 1990s could also be known as the decade when cryptography became mainstream. Perhaps because of the change in the Internet community, people became more aware of the need to protect the privacy of internetwork communications. Certainly, the U.S. government's attempt to push government control of cryptographic keys in the Clipper controversy helped to move cryptography and its related issues from science journals to the front pages of our newspapers. Today, while not mainstream, terms such as *Virtual Private Networks* (VPNs), *Secure Sockets Layer* (SSL), *IP Security* (IPSec), *Pretty Good Privacy* (PGP), *Secure Multipurpose Internet Mail Extensions* (S/MIME), and related technologies are known among IT professionals, and cryptography is no longer a tool used only by spies and military communication officers.

The Author

Richard E. Smith is well-known to members of various security-related forums on the Internet, as well as to security conference attendees. A security consultant with Secure Computing Corporation, Smith's background is in military-grade security. His experience on the lecture circuit, explaining issues of firewalls, cryptography, and other computer and network security topics, has directly contributed to production of a book on a lofty subject that is reachable by the nonscientist.

Organization

The chapters of this book fall into three groupings: an introduction to the basics of cryptography, its terms, methods, and mechanisms; network encryption and a discussion of VPNs, focusing on IPSec; and finally public key cryptography as it is used with message and file encryption and "Web" transactions.

The discussion in the opening chapter on basics may scare some off; Smith tends to oscillate between various levels of complexity. Consequently, some members of the intended audience of (quoting from the Preface) "people who know very little about cryptography but need to make technical decisions about cryptographic security," may, for example, zone out during the discussion of IP protocols. My suggestion would be to press on, and not worry about the random item that might go over your head. Everything there has a purpose, and the important information will fall into place by the end of each chapter.

If this book ended with Chapter 4, it would still be a useful book. The complex basics of cryptography and the issues that should be of concern to an information security officer are clearly presented and explained. The only area that is given less than adequate coverage is that of key recovery. Smith makes no mention of legitimate business reasons for the recovery of encrypted data if the originator is unavailable (the proverbial question, “What if you got hit by a truck?”), nor does he mention any mechanism other than the escrow of secret keys, although there are other, safer, methods. Of particular use are Smith’s explanations of the various cryptographic algorithms and his discussions of safe key lengths and risks.

In the sections on VPNs and IPSec, Smith covers everything from mobile users and remote access, to point-to-point encryption, and the issues of key distribution, exchange, and the mechanisms used to automate encrypted communication. Everyone seems to know that IPSec will save the world and is the answer to all our security problems (and I have my tongue firmly planted in my cheek), but few know what IPSec really does, from a “features and benefits” point of view. Of particular use and interest are the sections labeled “Deployment Example.” These are small case studies that show the technology in action and discuss some of the decisions and processes that came before deployment.

The section covering public key cryptography along with file and message encryption is perhaps shorter than it should be, although much of the groundwork is done earlier in the book. Missing is a “how to” on setting up a public key infrastructure (PKI) for a corporation to use. There are “Product Examples” in this section, but not “Deployment Examples.” Perhaps those will have to wait for a second edition, for although this is a lack in the book, there are not many real-life examples from which to choose. Although discussed in theory for years, this is still “leading edge” in the real world. The chapter on Web servers should prove informative and useful to any organization thinking of deploying (or having already deployed) a Web server.

In the chapter entitled “Secure Electronic Mail,” the fact that Smith covers *Privacy Enhanced Mail* (PEM) as a technology more than he covers S/MIME is puzzling, but the basics of PEM are useful for discussion, even if PEM as a technology seems to be dead.

Cryptography Is Necessary

The advertisement on the back of the book (not written by the author, of course) states “Here, in one comprehensive, soup-to-nuts book, is the solution for Internet security: modern-day cryptography.” Obviously the claim that cryptography is *the* solution for Internet security is way overinflated; modern-day cryptography is not *the* solution, but, cryptography is an important part of a “balanced” security solution. Smith does an admirable job of making this heretofore...well, cryptic... subject, understandable, interesting, and even enjoyable.

—Frederick M. Avolio, Avolio Consulting, fred@avolio.com

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and quality of service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ contains standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

ICANN

The *Internet Corporation for Assigned Names and Numbers* (ICANN) was incorporated in late October. ICANN is a private, non-profit corporation, managed by an international board, formed to coordinate and administer policies and technical protocols relating to the domain name and address system that permits Internet communications to be routed to the correct person or entity. Its proposed duties include those now performed under U.S. Government contract by the *Internet Assigned Numbers Authority* (IANA), whose Director, Internet pioneer Jon Postel, died on October 16th. ICANN has elected its Initial Board and chosen Michael M. Roberts as its Interim President and Chief Executive Officer. In addition, the Board chose Esther Dyson as its Interim Chairman, and appointed an Executive Committee consisting of Dyson, Gregory L. Crew, Hans Kraaijenbrink and Roberts. The other Initial Board members include Geraldine Capdeboscq (France), George H. Conrades (United States), Gregory L. Crew (Australia), Frank Fitzsimmons (United States), Hans Kraaijenbrink (The Netherlands), Jun Murai (Japan), Eugenio Triana (Spain), and Linda S. Wilson (United States). ICANN was originally proposed by Postel on behalf of a broad coalition of Internet stakeholders in response to the request by the U. S. Government last June that the Internet community create a global consensus non-profit corporation to which the U.S. could transition the responsibility for overseeing and funding those coordination activities. For more information, see:

<http://www.iana.org/index2.html>

APRICOT '99

The *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) will be held in Singapore, March 1–5, 1999. APRICOT provides a forum for key Internet builders in the region to learn from their peers and other leaders in the Internet community from around the world. The week-long summit consists of seminars, workshops, tutorials, conference sessions, birds-of-a-feather sessions, and other forums—all with the goal of spreading and sharing the knowledge required to operate the Internet within the Asia Pacific region. For more information, see: <http://www.apricot.net>

Send us your comments!

We look forward to hearing your comments and suggestions regarding anything you read in this publication. Send us e-mail at: ipj@cisco.com

This publication is distributed on an "as-is" basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or noninfringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Engineering
MCI WorldCom, USA

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Edward R. Kozel, Sr. VP, Corporate Development
Cisco Systems, Inc., USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the Cisco News
Publications Group, Cisco Systems, Inc.
www.cisco.com*

*Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 1998 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-J4
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

Bulk Rate Mail
U.S. Postage
PAID
Cisco Systems, Inc.

The Internet Protocol Journal

March 1999

Volume 2, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Peering and Settlements	2
IPv6.....	17
Secure E-Mail	30
Book Review.....	44
Letter to the Editor	46
Fragments	47

FROM THE EDITOR

Today's Internet is comprised of numerous interconnected *Internet Service Providers* (ISPs), each serving many constituent networks and end users. Just as individual regional and national telephone companies interconnect and exchange traffic and form a global telephone network, the ISPs must arrange for points of interconnection to provide global Internet service. This interconnection mechanism is generally called "peering," and it is the subject of a two-part article by Geoff Huston. In Part I, which is included in this issue, he discusses the technical aspects of peering. In Part II, which will follow in our next issue, Mr. Huston continues the examination with a look at the business arrangements (called "settlements") that exist between ISPs, and discusses the future of this rapidly evolving marketplace.

In the early 1990s, concern grew regarding the possible depletion of the IP version 4 address space because of the rapid growth of the Internet. Predictions for when we would literally run out of IP addresses were published. Several proposals for a new version of IP were put forward in the IETF, eventually resulting in IP version 6 or IPv6. At the same time, new technologies were developed that effectively slowed address depletion, most notably *Classless Inter-Domain Routing* (CIDR) and *Network Address Translators* (NATs). Today there is still debate as to if and when IPv6 will be deployed in the global Internet, but experimentation and development continues on this protocol. We asked Robert Fink to give us a status report on IPv6.

We've already discussed the historical lack of security in Internet technologies and how security enhancements are being developed for every layer of the protocol stack. This time, Marshall Rose and David Strom examine the state of electronic mail security. We clearly have a way to go before we see "seamless integration" of security systems with today's e-mail clients.

Our first Letter to the Editor is included on page 46. As always, we would love to hear your comments and questions regarding anything you read in this journal. Please contact us at ipj@cisco.com

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download
previous issues of IPJ in
PDF format from:
www.cisco.com/ipj

Interconnection, Peering and Settlements—Part I

by Geoff Huston, Telstra

Technology and business models share a common evolution within the Internet. To enable deployment of the technology within a service environment, a robust and stable business model also needs to be created. This tied destiny of technology and business factors is perhaps most apparent within the area of the interconnection of *Internet Service Providers* (ISPs). Here there is an interaction at a level of technology, in terms of routing signaling and traffic flows, and also an interaction of business models, in terms of a negotiation of benefit and cost in undertaking the interconnection. This article examines this environment in some detail, looking closely at the interaction between the capabilities of the technical protocols, their translation into engineering deployment, and the consequent business imperatives that such environments create.

It is necessary to commence this examination of the public Internet with the observation that the Internet is not, and never has been, a single network. The Internet is a collection of interconnected component networks that share a common addressing structure, a common view of routing and traffic flow, and a common view of a naming system. This interconnection environment spans a highly diverse set of more than 50,000 component networks, and this number continues, inexorably, to grow and grow. One of the significant aspects of this environment is the competitive Internet service industry, where many thousands of enterprises, both small and large, compete for market share at a regional, national, and international level.

Underneath the veneer of a highly competitive Internet service market is a somewhat different environment, in which every ISP network must interoperate with neighboring Internet networks in order to produce a delivered service outcome of comprehensive connectivity and end-to-end service. No ISP can operate in complete isolation from others while still offering public Internet services, and therefore, every ISP not only must coexist with other ISPs but also must operate in cooperation with other ISPs.

This article examines both the technical and business aspects that surround this ISP interaction, commonly referred to as “interconnection, peering, and settlements.” It examines the business motivation for interconnection structures, and then the technical architectures of such environments. The second part looks at the business relationships that arise between ISPs in the public Internet space, and then examines numerous broader issues that will shape the near-term future of this environment.

[This article is based in part on material in *The ISP Survival Guide*, by Geoff Huston, ISBN 201-3-45567-9, published by Wiley. Used with permission.]

Interconnection: Retailing, Reselling, and Wholesaling

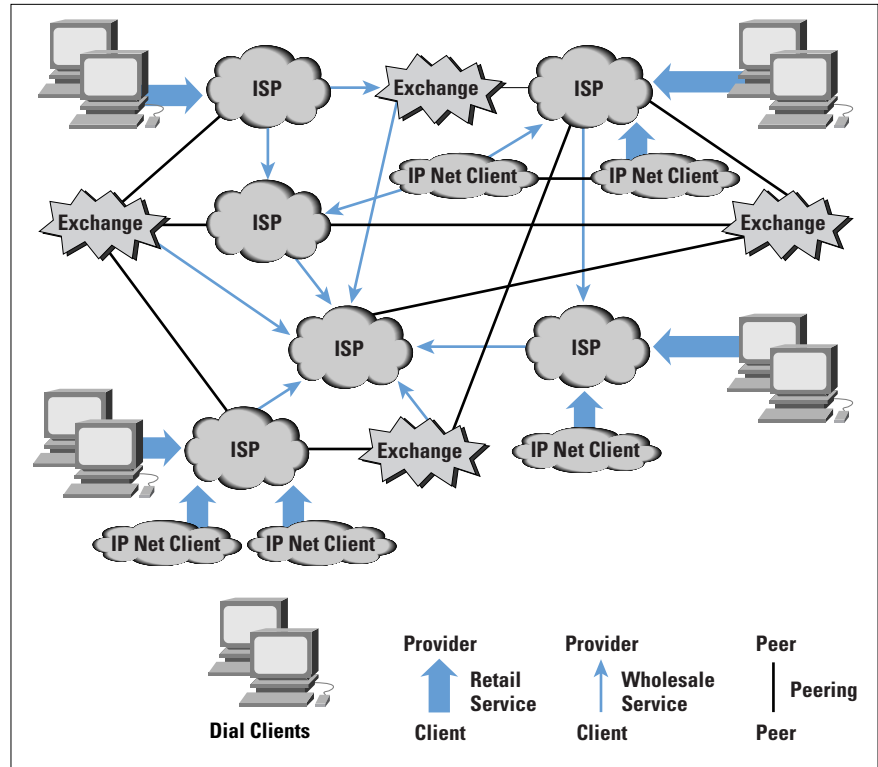
To provide some motivation for this issue of ISP interconnection, it is first appropriate to look at the nature of the environment. The regulatory framework that defined the traditional structure of other communications enterprises such as telephony or postal services was largely absent in the evolution of the Internet service industry. The resultant service industry for the Internet is most accurately characterized as an outcome of business and technology interaction, rather than a planned outcome of some regulatory process. This section examines this interaction between business and technology within the ISP environment.

A natural outcome of the Internet model is that the effective control of the retail service environment rests with a network client of an access service rather than with the access service provider. As such, a client of an ISP access service has the discretionary ability to resell the access service to third-party clients. In this environment, reselling and wholesaling are very natural developments within the ISP activity sector, with or without the explicit concurrence of the provider ISP. The provider ISP may see this reselling as an additional channel to market for its own Internet carriage services, and may adopt a positive stance by actively encouraging resellers into the market as a means of overall market stimulus, while tapping into the marketing, sales, and support resources of these reselling entities to continue to drive the volumes of the underlying Internet carriage service portfolio. The low barriers to entry to the wholesale market provide a means of increasing the scope of the operation, because to lift business cash-flow levels, the business enters into wholesale agreements that effectively resell the carriage components of the operation without the bundling of other services normally associated with the retail operation. This process allows the ISP to gain higher volumes of carriage capacity that in turn allow the ISP to gain access to lower unit costs of carriage.

Given that a retail operation can readily become a wholesale provider to third-party resellers at the effective discretion of the original retail client, is a wholesale transit ISP restricted from undertaking retail operations? Again, there is no such natural restriction from a technical or business perspective. An Internet carriage service is a commodity service that does not allow for a significant level of intrinsic product discrimination. The relatively low level of value added by a wholesale service operation implies a low unit rate of financial return for that operation. This low unit rate of financial return, together with an inability to competitively discriminate the wholesale product effectively, induces a wholesale provider into the retail sector as a means of improving the financial performance of the service operation. The overall result is that many ISPs operate both as clients and as providers. Few, if any, reasonable technical-based characterizations draw a clear and unambiguous distinction between a client and service provider when access services to networks are considered. A campus network may be a client of one or more ser-

vice providers, while the network is also a service provider to campus users. Indeed most networks in a similar situation take on the dual role of client and provider, and the ability to resell an access service can extend to almost arbitrary depths of the reselling hierarchy. From this technical perspective, very few natural divisions of the market support a stable segmentation into exclusively wholesale and exclusively retail market sectors. The overall structure of roles is indicated in Figure 1.

Figure 1:
ISP Roles and
Relationships



The resultant business environment is one characterized by a reasonable degree of fluidity, in which no clear delineation of relative roles or markets exists. The ISP market environment is, therefore, one of competitive market forces in which each ISP tends to create a retail market presence. However, no ISP can operate in isolation. Each client has the expectation of universal and comprehensive reachability, such that any client of any other ISP can reach the client, and the client can reach a client of any other ISP. The client of an ISP is not undertaking a service contract that limits connectivity only to other clients of the same ISP. Because no provider can claim ubiquity of access, every provider relies on every other provider to complete the user-provided picture of comprehensive connectivity. Because of this dependent relationship, an individual provider's effort to provide substantially superior service quality may have little overall impact on the totality of client-delivered service quality. In a best-effort public Internet, the service quality becomes something that can be impacted negatively by poor local engineering but cannot be uniformly improved beyond the quality provided by the network's peers, and their peers in turn. Internet wholesale carriage services in such an environment are constrained to be a com-

modity service, in which scant opportunity exists for service-based differentiation. In the absence of service quality as an effective service discriminator, the wholesale activity becomes a price-based service with low levels of added value, or in other words a commodity market.

The implication in terms of ISP positioning is that the retail operation, rather than the wholesale activity, is the major area in which the ISP can provide discriminating service quality. Within the retail operation, the ISP can offer a wide variety of services with a set of associated service levels, and base a market positioning on factors other than commodity carriage pricing.

Accordingly, the environment of interconnection between ISPs does not break down into a well-ordered model of a set of wholesale carriage providers and associated retail service providers. The environment currently is one with a wide diversity of retail-oriented providers, where each provider may operate both as a retail service operator, and a wholesale carriage provider to other retailers.

Peer or Client?

One of the significant issues that arises here is: Can an objective determination be made of whether an ISP is a peer to, or a client of, another ISP? This is a critical question, because if a completely objective determination cannot be readily made, the question then becomes one of who is responsible for making a subjective determination, and on what basis.

This question is an inevitable outcome of the reselling environment, where the reseller starts to make multiple upstream service contracts, with a growing number of downstream clients of the reselling service. At this point, the business profile of the original reseller is little distinguished from that of the original provider. The original reseller sees no unique value being offered by the original upstream provider and may conclude that it is, in fact, adding value to the original upstream provider by offering the upstream provider high-volume carriage and close access to the reseller's client base. From the perspective of the original reseller, the roles have changed, and the reseller now perceives itself as a peer ISP to the original upstream ISP provider.

This assertion of role reversal is perhaps most significant when the generic interconnection environment is one of "zero-sum" financial settlement, in which the successful assertion by a client of a change from client to peer status results in the dropping of client service revenue without any net change in the cost base of the provider's operation. The party making the successful assertion of peer interconnection sees the opposite, with an immediate drop in the cost of the ISP operation with no net revenue change.

The traditional public regulatory resolution of such matters has been through an administrative process of "licensed" communications service providers, who become peer entities through a process of

administrative fiat. In this model, an ISP becomes a licensed service provider through the payment of license fees to a communications regulatory body. The license then allows the service enterprise access to interconnection arrangements with other licensed providers. The determination of peer or client is now quite simple: A *client* is an entity that operates without such a carrier license, and a *peer* is one that has been granted such an instrument. However, such regulated environments are quite artificial in their delineation of the entities that operate within a market, and this regulatory process often acts as a strong disincentive to large-scale private investment, thereby placing the burden of underwriting the funding of service industries into the public sector. The regulatory environment is changing worldwide to shift the burden of communications infrastructure investment from the public sector, or from a uniquely positioned small segment of the private sector, to an environment that encourages widespread private investment. The Internet industry is at the leading edge of this trend, and the ISP domain typically operates within a deregulated valued-added communications service provider regulatory environment. Individual licenses are replaced with generic class licenses or similar deregulated structures in which formal applications or payments of license fees to operate in this domain are unnecessary. In such deregulated environments, no authoritative external entity makes the decision as to whether the relationship between two ISPs is that of a provider and client or that of peers.

If no public regulatory body wants to make such a determination, is there a comparable industry body that can undertake such a role? The early attempts of the *Commercial Internet eXchange* (CIX) arrangements in the United States in the early 1990s were based on a description of the infrastructure of each party, in which acknowledgments of peer capability were based on the operation of a national transit infrastructure of a minimum specified capability. This specification of peering within the CIX was subsequently modified so that CIX peer status for an ISP was simply based on payment of the CIX Association membership fee.

This CIX model was not one that intrinsically admitted bilateral peer relationships. The relationship was a multilateral one, in which each ISP executed a single agreement with the CIX Association and then effectively had the ability to peer with all other association member networks. The consequence of this multilateral arrangement is that the peering settlements can be regarded as an instance of “zero-sum” financial settlement peering, using a single-threshold pricing structure.

Other industry models use a functional peer specification. For example, if the ISP attaches to a nominated physical exchange structure, then the ISP is in a position to open bilateral negotiations with any other ISP also directly attached to the exchange structure. This model is inherently more flexible, as the bilateral exchange structure enables each represented ISP to make its own determination of whether to agree to a peer

relationship or not with any other colocated ISP. This model also enables each bilateral peer arrangement to be executed individually, admitting the possibility of a wider diversity of financial settlement arrangements.

The bottom line is that a true peer relationship is based on the supposition that either party can terminate the interconnection relationship and that the other party does not consider such an action a competitively hostile act. If one party has a high reliance on the interconnection arrangement and the other does not, then the most stable business outcome is that this reliance is expressed in terms of a service contract with the other party, and a provider/client relationship is established. If a balance of mutual requirement exists between both parties, then a stable basis for a peer interconnection relationship also exists. Such a statement has no intrinsic metrics that allow the requirements to be quantified. Peering in such an environment is best expressed as the balance of perceptions, in which each party perceives an acceptable approximation of equal benefit in the interconnection relationship in its own terms.

This conclusion leads to the various tiers of accepted peering that are evident in the Internet today. Local ISPs see a rationale to viewing local competing ISPs as peers, and they still admit the need to purchase trunk transit services from one or more upstream ISPs under terms of a client contract with the trunk provider ISP. Trunk ISPs see an acceptable rationale in peering with ISPs with a similar role profile in trunk transit but perceive an inequality of relationship with local ISPs. The conclusion drawn here is that the structure of the Internet is one in which there is a strong business pressure to create a rich mesh of interconnection at various levels, and the architecture of interconnection structures is an important feature of the overall architecture of the public Internet.

Physical Interconnection Architectures: Exchanges and NAPs

One of the physical properties of electromagnetic propagation is that the power required to transmit an electromagnetic pulse over a distance varies in accordance with this distance. The shorter the distance between the transmitter and the receiver, the lower the transmission power budget required; *closer is cheaper*.

This statement holds true not only for electrical power budgets but also for data protocol efficiency. Minimizing the delay between the sender and receiver allows the protocol to operate faster and operate more efficiently as well; *closer is faster*, and *closer is more efficient*.

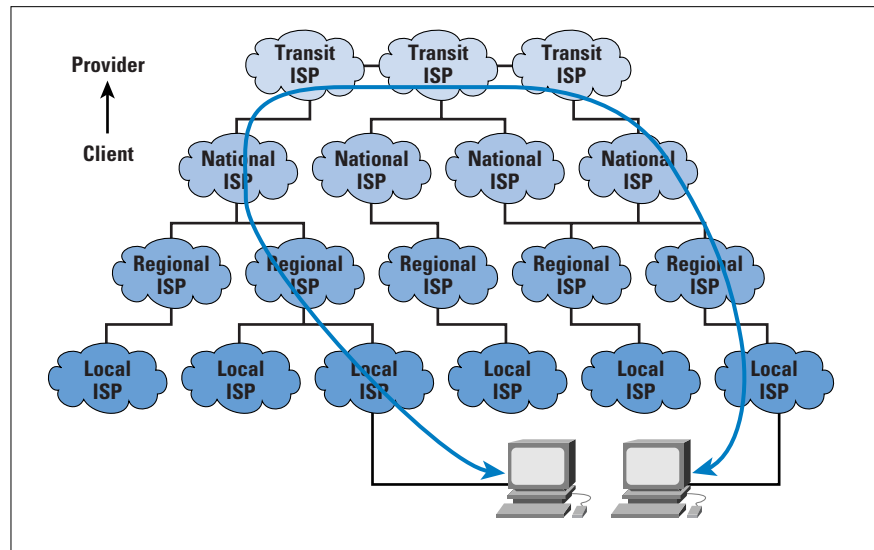
These observations imply that distinct and measurable advantages are gained by localizing data traffic; that is, by ensuring that the physical path traversed by the packets passed between the sender and the receiver is kept as physically short as possible. These advantages are realizable in terms of service performance, efficiency, and service cost.

How then are such considerations of locality factored into the structure of the Internet?

The Exchange Model

A strictly hierarchical model of Internet structure is one in which a small number of global ISP transit operators is at the “top;” a second tier is of national ISP operators; and a third tier consists of local ISPs. At each tier, the ISPs are clients of the tier above, as shown in Figure 2. If this hierarchical model is strictly adhered to, traffic between two local ISPs is forced to transit a national ISP, and traffic between two national ISPs transits a global ISP—even if both national ISPs operate within the same country. In the worst case, traffic between two local ISPs needs to transit a national ISP, then a global ISP from one hierarchy, then a second global ISP, and a second national ISP from an adjacent hierarchy in order to reach the other local ISP. If the two global providers interconnect at a remote location, the transit path of the traffic between these two local ISPs could be very long indeed.

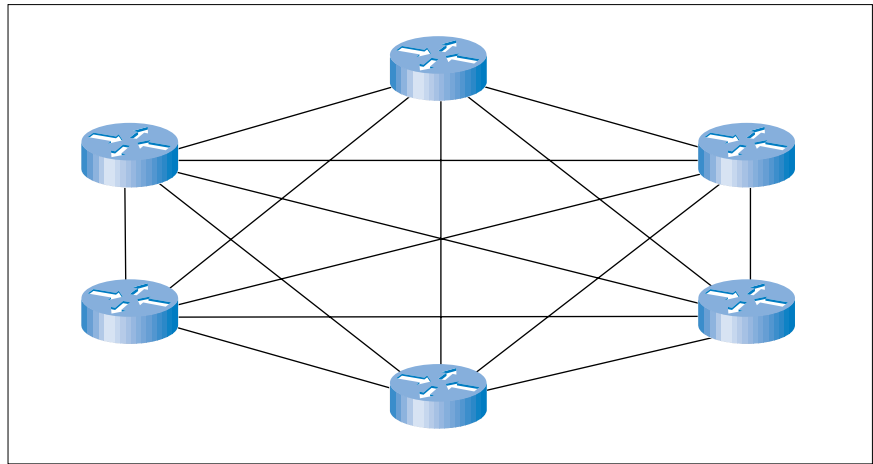
Figure 2:
A Purely Hierarchical
Structure for the
Internet



As noted above, such extended paths are inefficient and costly, and such costs are ultimately part of the cost component of the price of Internet access. In an open, competitive market, strong pressure always is applied to reduce costs. Within a hierarchical ISP environment, strong pressure is applied for the two national providers, who operate within the same market domain, to modify this strict hierarchy and directly interconnect their networks. Such a local interconnection allows the two networks to service their mutual connectivity requirements without payment of transit costs to their respective global transit ISP providers. At the local level is a similar incentive for the local ISPs to reduce their cost base, and a local interconnection with other local ISPs would allow local traffic to be exchanged without the payment of transit costs to the respective transit providers.

Although constructing a general interconnection regime based on point-to-point bilateral connections is possible, this approach does not exhibit good scaling properties. Between N providers who want to interconnect, the outcome of such a model of single interconnecting circuits is $(N^2 - N) / 2$ circuits and $(N^2 - N) / 2$ routing interconnections, as indicated in Figure 3. Given that interconnections exhibit the greatest leverage within geographical local situations, simplifying this picture within the structure of a local exchange is possible. In this scenario, each provider draws a single circuit to the local exchange and then executes interconnections at this exchange location. Between N providers who want to interconnect, the same functionality of complete interconnection can be constructed using only N point-to-point circuits.

Figure 3:
Fully Meshed Peering



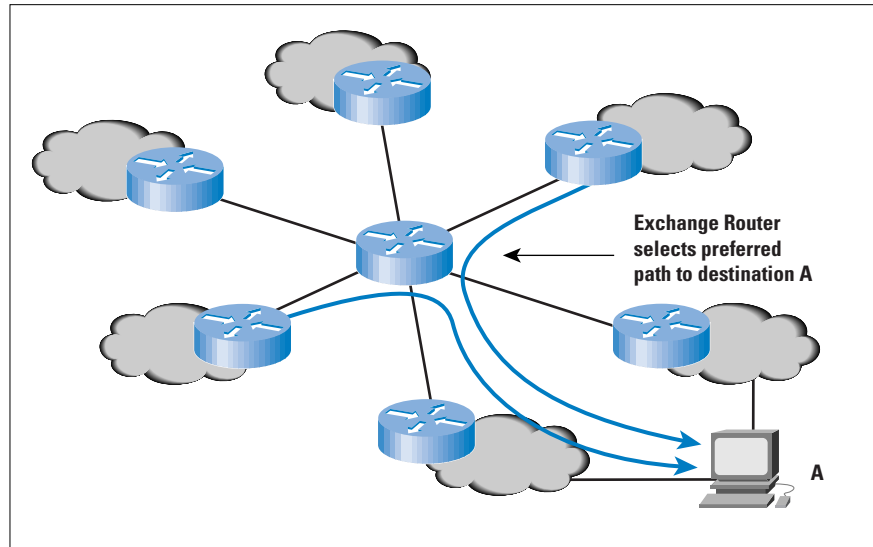
The Exchange Router

One model of an exchange is to build the exchange itself as a router, as indicated in Figure 4. Each provider's circuit terminates on the exchange router, and each provider's routing system peers with the routing process on the exchange router. This structure also simplifies the routing configuration, so that full interconnection of N providers is effected with N routing peer sessions. This simplification does allow greater levels of scaling in the interconnection architecture.

However, the exchange router model becomes an active component of the interconnect peering policy environment. In effect, each provider must execute a multilateral interconnection peering with all of the other connected providers. Selectively interconnecting with a subset of the providers present at such a router-based exchange is not easily achieved. In addition, this type of exchange must execute its own routing policy. When two or more providers are advertising a route to the same destination, the exchange router must execute a policy decision as to which provider's route is loaded in the router's forwarding table, making a policy choice of transit provider on behalf of all other exchange-connected providers.

Because the exchange is now an active policy element in the interconnection environment, the exchange is no longer completely neutral to all participants. This imposition on the providers may be seen as unacceptable, in that some of their ability to devise and execute an external transit policy is usurped by the exchange operator's policies.

Figure 4:
An Exchange Router



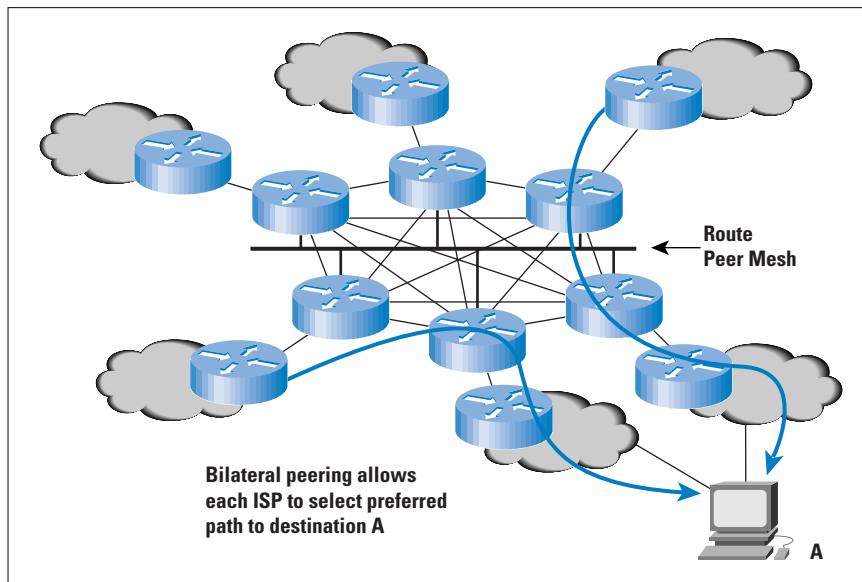
Typically, providers have a higher expectation of flexibility of policy determination from exchange structures than the base level of functionality that is provided by an exchange router. Providers want the flexibility to execute interconnections on a bilateral basis at the exchange, and to make policy decisions as to which provider to prefer when the same destination is advertised by multiple providers. They require the exchange to be neutral with respect to such individual routing policy decisions.

The Exchange Switch

The modification to the interprovider exchange structure is to use a local Layer 2 switch (or LAN) as the exchange element. In this model, a participating provider draws a circuit to the exchange and locates a dedicated router on the exchange LAN, as shown in Figure 5. Each provider executes a bilateral peering agreement with another provider by initiating a router peering session with the other party's router. When the same network destination is advertised by multiple peers, the provider can execute a policy-based preference as to which peer's route will be loaded in the local forwarding table. Such a structure preserves the cost efficiency of using N circuits to effect interconnection at the N provider exchange, while admitting the important policy flexibility provided by up to $(N^2 - N) / 2$ potential routing peer sessions.

Early interprovider exchanges were based on an Ethernet LAN as the common interconnection element. This physical structure was simple, and not all that robust under the pressures of growth as the LAN became congested.

Figure 5:
An Exchange LAN



Subsequent refinements to the model have included the use of Ethernet switches as a higher capacity LAN, and the use of *Fiber Distributed Data Interface* (FDDI) rings, switched FDDI hubs, Fast Ethernet hubs, and switched Fast Ethernet hubs. Exchanges are very-high-traffic concentration points, and the desire to manage ever-higher traffic volumes has led to the adoption of Gigabit Ethernet switches as the current evolutionary technology step within such exchanges.

The model of the exchange colocation accommodates a model of diversity of access media, in which the provider's colocated router undertakes the media translation between the access link protocol and the common exchange protocol.

The local traffic exchange hub does represent a critical point of failure within the local Internet topology. Accordingly, the exchange should be engineered in the most resilient fashion possible, using standards associated with a premium quality data center. This structure may include multiple power utility connections, uninterruptible power supplies, multiple trunk fiber connections, and excellent site security measures.

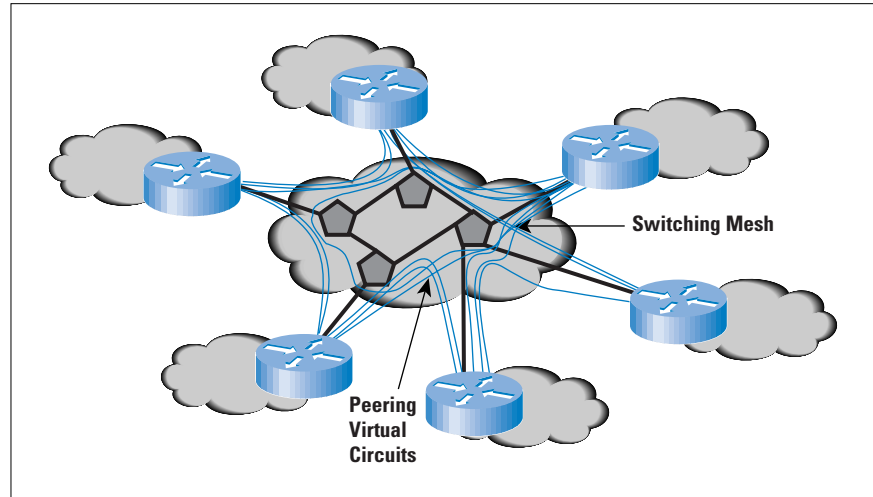
The exchange should operate neutrally with respect to every participating ISP, with the interests of all the exchange clients in mind. Thus, exchange facilities, which are operated by an entity that is not also a local or trunk ISP, enjoy higher levels of trust from the clients of the exchange.

There are also some drawbacks to an exchange, and a commonly cited example is that of imposed transit. If an exchange participant directs a default route to another exchange router, then in the absence of defensive mechanisms, the target router carries the imposed transit traffic even when there is no routing peering or business agreement between the two ISPs. Exchange-located routers do require careful configuration management to ensure that route peering and associated transit traffic matches the currently executed interconnection agreements.

Distributed Exchanges

Distributed exchange models also have been deployed in various locations. This deployment can be as simple as a metropolitan FDDI extension, in which the exchange comes to the provider's location rather than the reverse, as indicated in Figure 6. Other models that use an ATM-based switching fabric also have been deployed using *LAN Emulation* (LANE) to mimic the Layer 2 exchange switch functionality. Distributed exchange models attempt to address the significant cost of operating a single colocation environment with a high degree of resilience and security, but do so at a cost of enforcing the use of a uniform access technology between every distributed exchange participant.

Figure 6:
A Distributed Exchange



However, the major challenge of such distributed models is that of switching speed. Switching requires some element of contention resolution, in which two ingress data elements that are addressed to a common egress path require the switch to detect the resource contention and then resolve it by serializing the egress. Switching, therefore, requires signaling, in which the switching element must inform the ingress element of switch contention. To increase the throughput of the switch, the latency of this signaling must be reduced. The dictates of increased switching speed have the corollary of requiring the switch to exist within the confines of a single location, if exchange performance is a paramount concern.

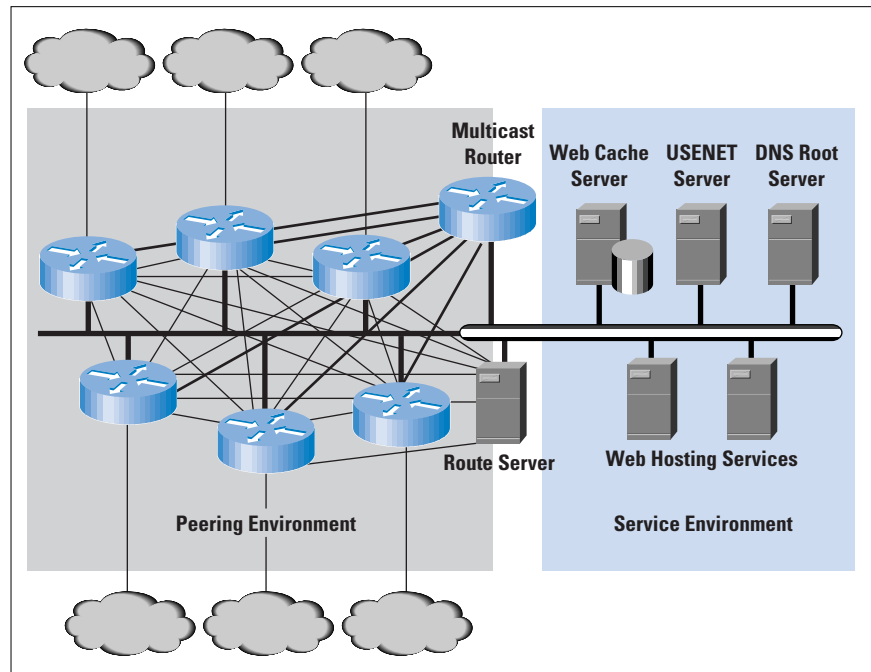
In addition to speed, the cost shift must be considered. In a distributed exchange model, the exchange operator operates the set of access circuits that form the distributed exchange. This process increases costs to providers, while it prevents the providers from using a specific access technology that matches their business requirements of cost and supportable traffic volume. Not surprisingly, to date the most prevalent form of exchange remains the third-party hosted colocation model. This model admits a high degree of diversity in access technologies, while still providing the substrate of an interconnection environment that can operate at high speed and therefore manage high traffic volumes.

Other Exchange-Located Services

The colocation environment is often broadened to include other functions, in addition to a pure routing and traffic exchange role. For a high-volume content provider, the exchange location offers minimal transit distance to a large user population distributed across multiple local service providers, as well as allowing the content provider to exercise a choice in selecting a nonlocal transit provider.

The exchange operator can also add value to the exchange environment by providing additional functions and services, as well as terminating providers' routers and large-volume content services. The exchange location within the overall network topology is an ideal location for hosting multicast services, because the location is optimal in terms of multicast carriage efficiency. Similarly, USENET trunk feed systems can exploit the local hub created by the exchange. The overall architecture of a colocation environment that permits value-added services, which can productively use the unique environment created at an exchange, is indicated in Figure 7.

Figure 7:
Exchange-Located
Service Platforms



Network Access Points

The role of the exchange was broadened with the introduction of the *Network Access Point* (NAP) in the architecture proposed by the National Science Foundation (NSF) in 1995 when the NSFNET backbone was being phased out.

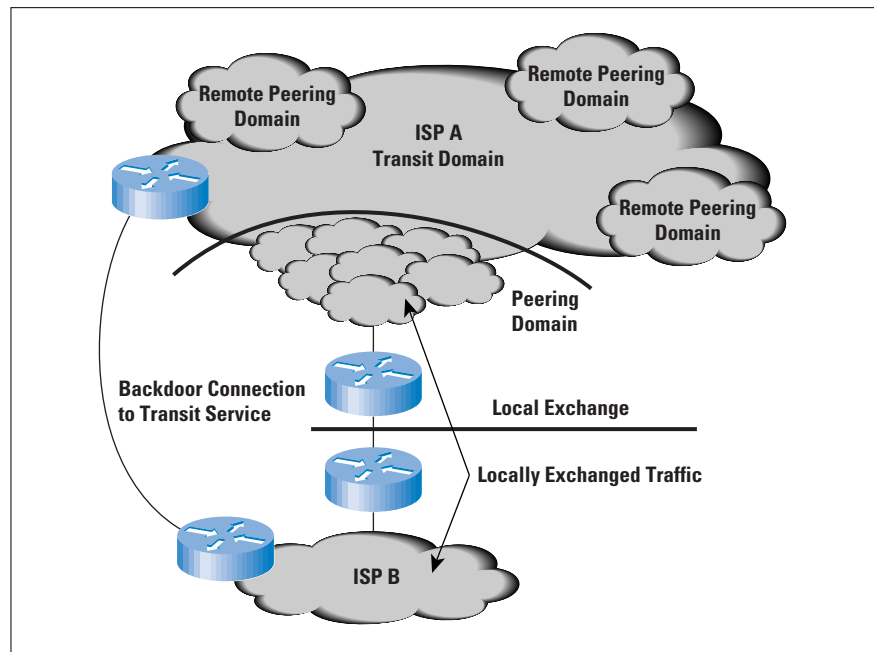
The NAP was seen to undertake two roles: the role of an exchange provider between regional ISPs who want to execute bilateral peering arrangements and the role of a transit purchase venue, in which regional ISPs could execute purchase agreements with one or more of a set of trunk carriage ISPs also connected at the NAP. The access point concept was intended to describe access to the trunk transit service.

This mixed role of both local exchange and transit operations leads to considerable operational complexity, in terms of the transit providers being able to execute a clear business agreement. What is the bandwidth of the purchased service in terms of requirements for trunk transit, versus the access requirements for exchange traffic? If a local ISP purchases a transit service at one of the NAPs, does that imply that the trunk provider is then obligated to present all the ISP's routes at remote NAPs as a peer? How can a trunk provider distinguish between traffic presented to it on behalf of a remote client versus traffic presented to it by a local service client?

The issue that the quality of the purchased transit service is colored by the quality of the service provided by the NAP operator should also be considered. Although the quality of the transit provider's network may remain constant, and the quality of the local ISP's network and ISP's NAP access circuit may be acceptable, the quality of the transit service may be negatively impacted by the quality of the NAP transit itself.

One common solution is to use the NAP colocation facility to execute transit purchase agreements and then use so-called *backdoor* connections for the transit service provision role. This usage restricts the NAP exchange network to a theoretically simpler local exchange role. Such a configuration is illustrated in Figure 8.

Figure 8:
Peering and Transit
Purchase



Exchange Business Models

For the ISP industry, many attributes are considered highly desirable for an exchange facility. The common model of an Internet exchange includes many, if not all, of the following elements:

- Operated by a neutral party who is not an ISP (to ensure fairness and neutrality in the operation of the exchange)
- Constructed in a robust and secure fashion
- Located in areas of high density of Internet market space
- Able to scale in size
- Operates in a fiscally sound and stable business fashion

A continuing concern exists about the performance of exchanges and the consequent issue of quality of services that traverse the exchange. Many of these concerns stem from an exchange business model that may not be adequately robust under pressures of growth from participating ISPs.

The exchange business models typically are based on a flat-fee structure. The most basic model uses a fee structure based on the number of rack units used by the ISP to colocate equipment at the exchange. When an exchange participant increases the amount of traffic presented over an access interface, under a flat-fee structure, this increased level of traffic is not accompanied by any increase in exchange fees. However, the greater traffic volumes do imply that the exchange itself is faced with a greater traffic load. This greater load places pressure on the exchange operator to deploy further equipment to augment the switching capacity, without any corresponding increase in revenue levels to the operator.

For an exchange operator to base tariffs on the access bandwidths is not altogether feasible, given that such access facilities are leased by the participating ISPs and the access bandwidth may not be known to the exchange operator. Nor is using a traffic-based funding model possible, because an exchange operator should refrain from monitoring individual ISP traffic across the exchange, given the unique position of the exchange operator. Accordingly, the exchange operator has to devise a fiscally prudent tariff structure at the outset that enables the exchange operator to accommodate large-scale traffic growth, while maintaining the highest possible traffic throughput levels.

Alternatively, there are business models in which the exchange is structured as a cooperative entity among numerous ISPs. In these models, the exchange is a nonprofit common asset of the cooperative body. Although widely used, these models are prone to the economic condition of the *Tragedy of the Commons*. It is in everyone's interest to maximize their exploitation of the exchange, while no single member wants to underwrite the financial responsibility for ensuring that the quality of the exchange itself is maintained.

The conclusion that can be drawn is that the exchange is an important component of Internet infrastructure, and the quality of the exchange is of paramount importance if it is to be of any relevance to ISPs. Using an independent exchange operator whose income is derived from the utility of the exchange is one way of ensuring that the exchange is managed proficiently and that the service quality is maintained for the ISP clients of the exchange.

A Structure for Connectivity

Enhancing the Internet infrastructure is quantified by the following objectives:

- Extension of reachability
- Enhancement of policy matching by ISPs
- Localization of connectivity
- Backup arrangements for reliability of operation
- Increasing capacity of connectivity
- Enhanced operational stability
- Creation of a rational structure of the connection environment to allow scalable structuring of the address and routing space in order to accommodate orderly growth

We have reached a critical point within the evolution of the Internet. The natural reaction of the various network service entities in response to the increasing number of ISPs will be to increase the complexity of the interconnection structure to preserve various direct connectivity requirements. Today, we are in the uncomfortable position of increasingly complex interprovider connectivity environments, a situation that is stressing the capability of available technologies and equipment. The inability to reach stable cost-distribution models in a transit arrangement creates an environment in which each ISP attempts to optimize its position by undertaking as many direct 1:1 connections with peer ISPs as it possibly can. Some of these connections are managed via the exchange structure. Many more are implemented as direct links between the two entities. Given the relative crudity of the inter-*Autonomous System* (AS) routing policy tools that we use today, this structure must be a source of considerable concern. The result of a combination of an increasingly complex mesh of inter-AS connections, together with very poor tools to manage the resultant routing space, is an increase in the overall instability of the Internet environment. In terms of meeting critical immediate objectives, however, such dire general predictions do not act as an effective deterrent to these actions.

The result is a situation in which the inter-AS space is the critical component of the Internet. This space can be viewed correctly as the *demilitarized zone* within the politics of today's ISP-based Internet. In the absence of any coherent policy, or even a commonly accepted set of practices, the lack of administration of this space is a source of paramount concern.

GEOFF HUSTON holds a B.Sc and a M.Sc from the Australian National University. He has been closely involved with the development of the Internet for the past decade. He was responsible for the initial build of the Internet within the Australian academic and research sector. Huston is currently the Chief Technologist in the Internet area for Telstra. He is also an active member of the IETF, and is a member of the Internet Society Board of Trustees. He is author of *The ISP Survival Guide*, and coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, a collaboration with Paul Ferguson. Both books are published by John Wiley & Sons. E-mail: gih@telstra.net

IPv6—What and Where It Is

by Robert L. Fink, *Energy Sciences Network*

The current Internet Protocol, known as IPv4 (for version 4), has served the Internet well for over 20 years, but is reaching the limits of its design. It is difficult to configure, it is running out of addressing space, and it provides no features for site renumbering to allow for an easy change of *Internet Service Provider* (ISP), among other limitations. Various mechanisms have been developed to alleviate these problems (for example, *Dynamic Host Configuration Protocol* [DHCP] and *Network Address Translation* [NAT]), but each has its own set of limitations.

The *Internet Engineering Task Force* (IETF) took on this problem in the early 1990s by starting an IPng (*Internet Protocol next generation*) project. After an over two-year-long process of defining goals and features, getting the best possible advice from industry and user experts, and sponsoring a protocol design competition, a new Internet Protocol was selected. Many proposed protocols were reviewed, analyzed, and evaluated. An evolved combination of several of them (*Simple Internet Protocol* [SIP], the “P” *Internet Protocol* [PIP], and *Simple Internet Protocol Plus* [SIPP]), each using fixed-length addressing, resulted in a final variation, called IPv6, which was selected over a version of the ISO OSI *Connectionless Network Protocol* (CLNP) (known as the *TCP and UDP with Bigger Addresses* (TUBA) IPng proposal).

Much work has been done since the selection of IPv6 in 1994. Over 50 implementations of IPv6 are believed to be under way or completed. A constantly growing international IPv6 testbed, called the *6bone*, now spans 260 sites in 39 countries, with over 25 different IPv6 implementations in use. Most router companies, including 3Com, Bay, Cisco Systems, Digital, Nokia, and Telebit support IPv6. IPv6 is also available for Digital, HP, IBM, Sun, WinTel, and many other end-user host systems.

IPv6 Addresses—Larger and Different

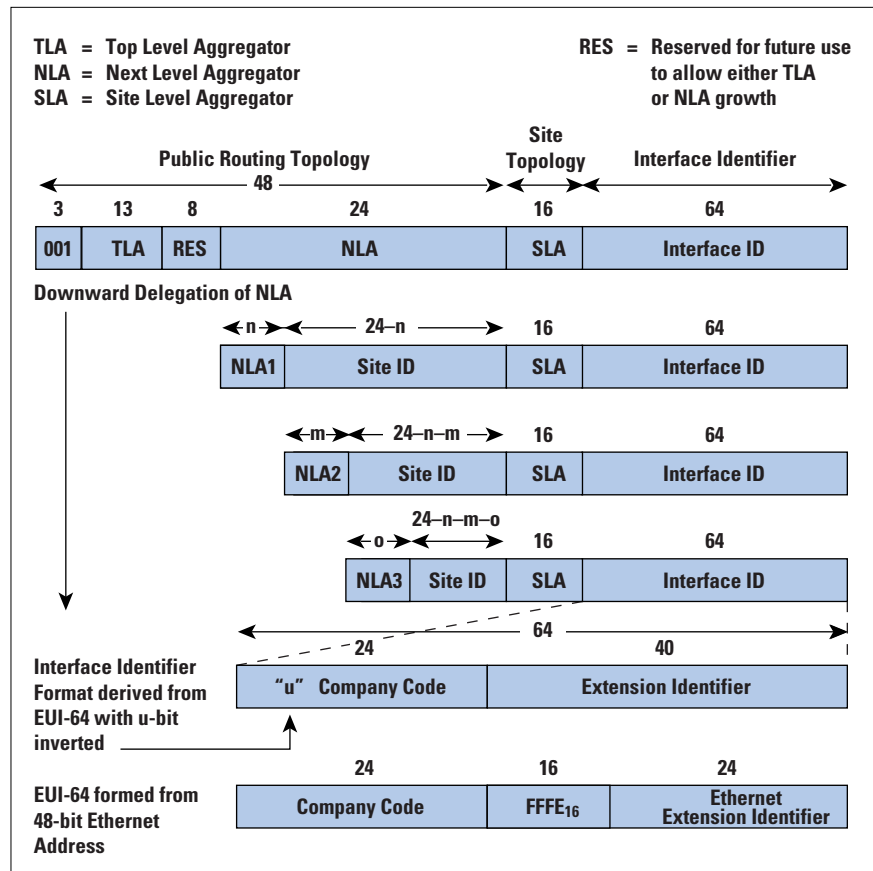
The larger 128-bit IPv6 address (versus the 32-bit IPv4 address) allows more flexibility in designing newer addressing architectures, as well as providing large enough address spaces for predicted future growth of the Internet and Internet-related technologies. A new addressing format, called the *Aggregatable Global Unicast Address Format*, has been developed to help solve route complexity scaling problems with the current IPv4 Internet. The current IPv4 provider-based addressing used in the Internet relies on separate IPv4 addresses being assigned to ISPs in contiguously numbered blocks for routing efficiency; that is, the routers need to carry fewer routes.

However, there is currently much fragmentation in the IPv4 address space. This situation, aggravated by sites not being able to easily renumber, causes many more separate routes than necessary, in turn leading to route computation complexity (too many routes, too many dynamic changes, too much computation in routers).

Public Routing Topology Prefixes

With the new aggregatable style addressing (see Figure 1), the left-most 48 bits of the address are defined as a *Public Routing Topology* (PRT) prefix. The first 3-bit field of this prefix specifies that the addressing format is aggregatable. The next 13-bit portion specifies the *Top Level Aggregator* (TLA) ID that constrains the top level of Internet routing to 8,192 major transit providers and a new concept of routing exchanges. Each TLA (top level transit ISP) is then responsible for all the remaining public routing topology assignment below it; that is, the *Next Level Aggregator* (NLA) ID. As shown in Figure 1, the NLA may have a tiered hierarchy to allow multiple levels (NLA1, NLA2, and so on) of other ISPs, each of which would then have control of the assignment of the space below it. The right-most portion of the NLA field, at whatever level it may be, would identify the end-user “leaf” site. An 8-bit reserved field has been defined to allow the growth of either the TLA or the NLA fields.

Figure 1:
Aggregatable Global
Unicast Address
Format



The advantage of this style of addressing is that it allows automatic address clustering, or aggregation, into a constrained set of routes, which are represented through the TLA field. If the initial assignment of 13 bits (8,192 TLAs) is insufficient in the future, either the reserved field or another piece of the IPv6 128-bit address space could be utilized. Note that only one-eighth of the current IPv6 address space has been assigned to aggregatable addressing.

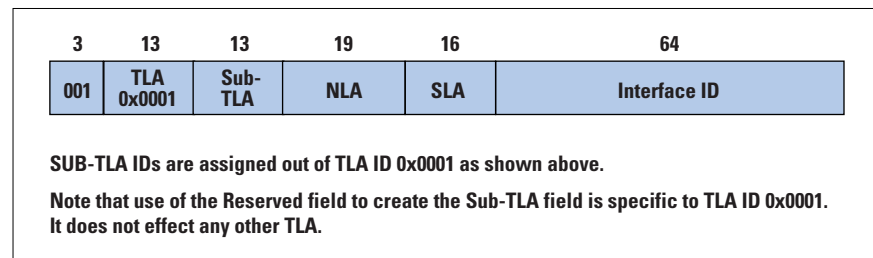
Even with this new concept of addressing, sites will still occasionally want to change their ISP (as in the current IPv4-based Internet) and thus will need to readdress to keep the addressing structure constrained. This is where *Site Renumbering*, which will be discussed later, comes in.

IPv6 TLA Assignment

To begin the production use of IPv6, ISPs providing IPv6 service need to be assigned TLAs so they may assign NLAs to transits and sites they are serving. Until recently, this was not possible. Recent discussions between the IETF, the IANA (*Internet Assigned Numbers Authority*), and the major address registries (APNIC, ARIN, and RIPE-NCC), have resulted in agreements that will provide a way to request and assign TLAs by early 2nd quarter 1999.

The process agreed upon is based on the above discussions that have been published as a recommendation in an Informational RFC on TLA assignments. The basic idea is to provide a slow start mechanism for TLAs by assigning one TLA ID to be used for defining a Sub-TLA field of 13 bits out of the reserved and NLA fields (see Figure 2). This will allow transits to demonstrate their need for a full TLA based on usage of the assigned Sub-TLA. These rules, based on much current practice with IPv4, are necessary to keep aggregatable addressing functional and effective for hierarchical routing as IPv6 comes into use.

Figure 2:
Sub-TLA Format for
IPv6 Address
Assignment



Rules for assigning these Sub-TLAs include:

- Must have a plan to offer native IPv6 service within three months from assignment; must have a verifiable track record providing Internet transit to other organizations
- Must make payment of a registration fee to the IANA and reasonable fees for services rendered by the address registry
- Must maintain registries of sites and next-level providers and make them available publicly and to the registries; must provide utilization statistics of NLA space below the assigned TLA (or Sub-TLA) and also show evidence of carrying TLA routing and transit traffic

These rules are intended to minimize route explosion and address assignment misuse to aid in the stability of the IPv6-based Internet.

Site Topology Prefixes

In addition to identifying the address of the site with the PRT prefix, aggregatable addressing provides for a site to have aggregation as well using a 16-bit *Site Level Aggregator* (SLA). The SLA might be as simple as a subnet number (more than 64,000 of them!), or a tiered hierarchy such as the NLA provides. However it is structured, the SLA is under the control of the site, and identifies the subnet that a host interface is attached to (IPv6's addressing, as IPv4's, specifies interfaces on systems, not the entire system).

It is very unlikely that an organization will ever need more than one PRT prefix, given the size and flexibility of the SLA and the *System Interface Identifier* field (described below).

System Interface Identifiers

Now that we have identified how to reach the site and the subnet a system is attached to, an interface identifier (ID) specifies the local logical address of the interface on the local subnet (or *link* as it is often called). The interface ID is formed and derived from the new IEEE EUI-64 media-level address that is an expansion of the well-known Ethernet 48-bit address format that allows for more device identifiers to be assigned by each manufacturer. The global/local bit is also inverted to make manually assigned (that is, local) addresses easy to form with only leading zeros.

If the IPv6 node is attached to an Ethernet "link," then the 48-bit address is turned into 64 bits by a filler field inserted in the middle (see Figure 1).

This enlarged Interface ID will allow newer technologies, such as *FireWire*, and newer applications, such as traffic lights and PCS/PDA telephones, to have unique interface identifiers assigned to them from a global address space.

The use of a media-level address for a network-level Interface ID allows the very important IPv6 Stateless Address Autoconfiguration Protocol to work.

Stateless Address Autoconfiguration

Automatic configuration of IPv6 end systems (hosts) is one of the most important features of IPv6. In the current IPv4 Internet, you must either manually configure IP address, network mask, and default gateway, or rely on having a DHCP server. With IPv6, this process can take place automatically, with no reliance on outside systems, using the IPv6 *Stateless Address Autoconfiguration Protocol*.

This can be done because the *Media Access Control* (MAC) address is used to form the host's interface ID. For example, if a host has an Ethernet interface that it is trying to configure for use with IPv6, the 48-bit Ethernet MAC address is formed into a 64-bit interface ID, which is the right-most 64 bits of the IPv6 address (see Figure 1). Then, using the *Neighbor Discovery* (ND) protocol, which is unique to IPv6, this formed interface ID is checked to see that it does not have a duplicate on this link (that is, subnet). If it does, a randomly generated token can be used (though a rare occurrence, it is a necessary protection against illegal Ethernet address usage and situations where the same address may be used on multiple interfaces for legitimate reasons).

At this point, an *ND Router Solicitation* multicast message is sent out to discover if there is a local IPv6 capable router, what the local site's topology ID for the host's subnet is, and what the site's public topology routing prefix is. Neighbor Discovery can also be used to control whether the site then wishes to continue with further configuration using Stateful Autoconfiguration with DHCPv6.

IPv6 Autoconfiguration thus provides for standalone operation of two or more hosts on a local LAN link with no router present, provides for operation within a site with no outside Internet connectivity present, and allows for easy changing of the site's public topology routing prefix, either when external connectivity comes on line, or when the external connectivity is changed, such as when a different ISP is chosen.

Domain Name System—Forward and Reverse

The *Domain Name System* (DNS) is an essential component of the Internet. To provide a mapping from a domain name to an IPv6 address, as well as an IPv4 address, a new DNS record type of "AAAA," or "quad A," is defined. This is a clever word play on the "A" record type that the original DNS specification defines for 32-bit IPv4 addresses, because IPv6 addresses are four times larger (128-bits), hence "AAAA"!

Most existing implementations of DNS already support AAAA records and existing IPv4 queries of DNS can access these records; that is, you don't need a DNS operating over IPv6 to retrieve these new AAAA records. This support also includes reverse lookups, similar to IPv4s, although a new reverse lookup proposal that will allow automatic partitioning of the delegation information on arbitrary bit boundaries is under consideration. This new capability should make for more reliable reverse registry than exists with IPv4, and easier maintenance when sites change their PRT prefix.

When a host with both IPv4 and IPv6 operating on it ("dual stack") queries the DNS for the address of a remote host, the A and AAAA records returned are used to indicate what protocol to use in communicating with that remote host. If no AAAA record is returned, IPv4 must be used. If only a AAAA record is returned, IPv6 must be used. If both A and AAAA are returned, either IPv4 or IPv6 may be used.

A new modification of the IPv6 DNS extensions is nearing completion that allows the automatic joining of the routing prefixes and Interface IDs when a host's IPv6 address is returned, thus making it easier to renumber a site. This new IPv6 DNS feature makes changing a site's PRT prefix (renumbering) very easy as only one entry, the PRT prefix, needs to be changed. This setup also facilitates easy support of multiple addresses for each host. These enhancements are very useful; IPv4 does not have this feature.

Renumbering Sites When ISPs Change

Because IPv6 addressing is based on the PRT prefix assigned by its ISP, it is essential that it be easy for a site to renumber itself when its choice of ISP changes. To aid in this, a new *Router Renumbering* (RR) protocol, in conjunction with Autoconfiguration, Neighbor Discovery and the new Aggregatable Unicast addressing PRT prefix are used.

RR allows a site's network administrator to set new PRT prefixes into the site's routers, as well as lower the lifetime of existing ISP PRT prefixes to specify an overlap interval, after which the old ISP's service is discontinued.

Hosts learn their new routing prefixes either when they restart, and thus are automatically configured with Autoconfiguration, or when they are informed by their local router that a new prefix is to be used during periodic router notification updates using ND.

For example, a new ISP service is readied for service while the old ISP is notified that it will provide service for just 60 more days. After the new PRT prefix is announced to the site's routers by RR, hosts will use the new prefix (that is, new ISP) for all new connections, while existing connections continue to work until the old prefix is withdrawn (that is, after 60 days in this example).

The easy renumbering of an IPv6 site will make easy a task that is currently very painful for an IPv4 site because hosts are often manually configured in many networks.

The 6bone—An IPv6 Testbed

The 6bone is an international IPv6 testbed network that is overseen and directed through the IETF *IPng Transition Working Group* (ngtrans) that provides:

- Testing of IPv6 implementations and standards
- Testing of IPv6 transition strategies
- A place to gain early applications and operations experience
- Motivation and a place for implementers, users, and ISPs to try IPv6
- An experimental first step toward transition

In the early phases of IPv6 deployment, most native IPv6 transport is restricted to site LANs with the ability to experiment with it locally. Some sites in Great Britain, The Netherlands, and Japan are using native IPv6 over WAN links.

ISPs and various other private IPv4 transit providers may not place IPv6 in their production routers in this early phase of IPv6 deployment, leaving early IPv6 testers with the need to use the existing IPv4 Internet infrastructure to deliver IPv6 packets among themselves when remotely located. Thus an IPv6 transition feature, IPv6 encapsulation (that is, *tunneling*) over IPv4, is used for parts of the 6bone where native IPv6 may not be available. In this way, the 6bone is also thoroughly testing out its own transition technology as well as providing IPv6 service.

The 6bone is a diverse community of users, ISPs, and developer organizations, many of whom provide transit on the public spirited basis of promoting and gaining early experience with IPv6. It is expected that production variations of the 6bone will also be created to more formally carry production IPv6 traffic.

Components of the 6bone

The 6bone provides this needed IPv6 transport over the public Internet infrastructure, relying on:

- Dual IPv4/IPv6 stacks in the client host
- IPv6 packets encapsulated (tunneled) in IPv4 packets
- Dual IPv4/IPv6 stack backbone routers that know IPv6 routes of 6bone participants
- DNS that supports IPv6 AAAA records
- A 6bone Routing Registry to keep track of sites and their tunnels
- A mailing list, various IPv6 tools, and a 6bone Web site at:
www.6bone.net

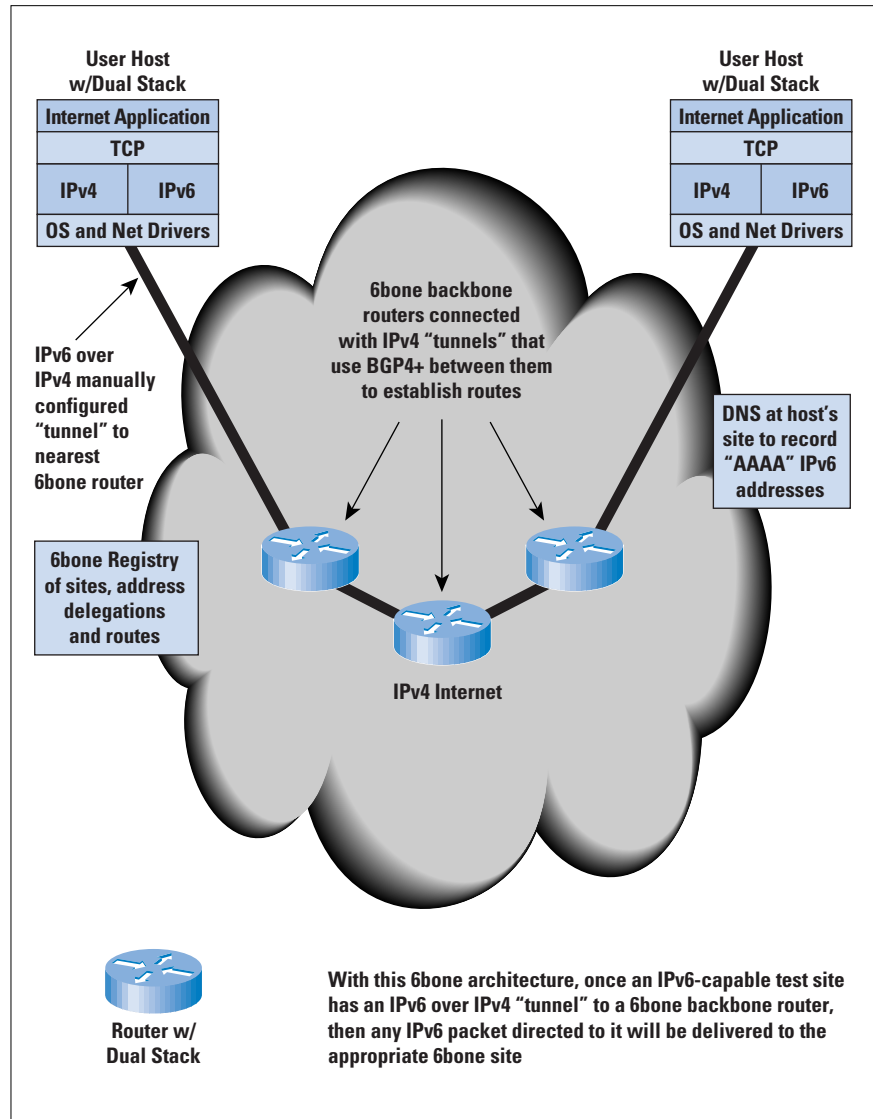
Figure 3 shows a conceptual overview of how a basic 6bone is structured and a picture of the current 6bone backbone structure can be seen at:

<http://www.cs-ipv6.lancs.ac.uk/ftp-archive/6Bone/Maps/full-backbone.gif>

...with the pseudo TLA site-to-site peering indicated by various colored links.

To date, the 6bone has spread to 260 organizations in 39 countries (see Table 1 on page 25).

Figure 3:
6bone Conceptual
Architecture



6bone History

Serious work to evolve and refine the IPv6 protocols sufficient to allow the start of various implementations of IPv6 began in 1994. By early 1996, it was obvious that a testing environment was needed, so in March 1996, several implementers and users met and agreed to start an international testbed called the 6bone.

By June 1996, two groups raced to provide the first IPv6 connectivity: the University of Lisbon (Portugal), the Naval Research Laboratory (U.S.), and Cisco Systems (U.S.); a Danish universities consortium (UNI-C), a French universities consortium (G6), and a Japanese universities consortium (WIDE).

Table 1: Countries with Sites Participating in the 6bone

AT-Austria	FI-Finland	NL-The Netherlands
AU-Australia	FR-France	NO-Norway
BE-Belgium	GB-United Kingdom	PL-Poland
BG-Bulgaria	GR-Greece	PT-Portugal
BR-Brazil	HK-Hong Kong	RO-Romania
CA-Canada	HU-Hungary	RU-Russian Federation
CH-Switzerland	IE-Ireland	SE-Sweden
CM-Cameroon	IT-Italy	SG-Singapore
CN-China	JP-Japan	SI-Slovenia
CZ-Czech Republic	KR-Korea	SK-Slovakia
DE-Germany	KZ-Kazakhstan	TW-Taiwan
DK-Denmark	LT-Lithuania	US-United States
ES-Spain	MX-Mexico	ZA-Zaire

6bone Backbone and Addressing

By the end of 1997, the 6bone converted to the new aggregatable addressing format, a change necessitated by having originally adopted an early prototype provider-based addressing format discussed during early IPv6 design efforts.

Along with the change to a new addressing format was the need to clean up the routing used among the 6bone backbone transit sites. It was originally thought that IDRPv6 (a new Internet Domain Routing Protocol based on earlier IPv4 work) would be the prevailing *Exterior Gateway Protocol* (EGP) used for IPv6 Internet peering.

By mid 1996, various ISPs made it known that a new EGP for IPv6 was not a practical alternative, given the explosive growth of the Internet and the current evolution and widespread use of the *Border Gateway Protocol 4* (BGP4) by ISPs. There was a need to allow for multiprotocol extensions to BGP4, allowing ISPs to more easily adapt their operations to IPv6. This situation led to the rapid evolution of BGP4+, an extension of BGP4 to include IPv6 and IPv4 multiprotocol routing.

By mid 1997, the decision was made to convert the 6bone backbone to BGP4+ for its EGP. See <http://www.cs-ipv6.lancs.ac.uk/ftp-archive/6Bone/Maps/full-backbone.gif> for a recent picture of the 6bone backbone sites using the new aggregatable addressing format and the current status of the conversion to BGP4+.

6bone Future Plans

To date, most 6bone efforts have been to prove out basic IPv6 interoperability among the many implementations, and to create a reliable international testbed infrastructure. This has included making its backbone operationally ready with the new aggregatable addressing format and use of BGP4+ for high-reliability routing and transit.

Now that the 6bone has completed these conversions, serious work can begin on testing site renumbering, security, applications, and transition mechanisms.

Other IPv6 Trials and Testing

Other testing venues have also been very important to the evolution of IPv6: the University of New Hampshire *Inter Operability Laboratory* (IOL), various trade show demonstration networks, for example, Net-World+Interop, and various vendor-sponsored interoperability testing.

By early 1998, the UNH IOL had hosted five IPv6 test sessions, though specific details about participating vendors are not released.

In a positive sign of industry response to evolving IPv6 specifications, the late July 1997 UNH testing resulted in the successful interoperability of all participants using the new aggregatable addressing format, no more than two months from its first Internet Draft.

Implementations

To date, over 50 different IPv6 host and router implementations are either completed or under way. More than 30 implementations have been tested and used on the 6bone.

Router implementations to date include: 3Com, Bay, Cisco Systems, Digital, Fujitsu LR550, Hitachi NR60, Inria BSD, Linux, Merit MRT, Nokia, NRL for BSD, Telebit, WIDE KAME and ZETA for BSD, and WIDE v6d.

Host implementations to date include: Apple MacOS OpenTransport demo version, Digital OpenVMS, Digital UNIX, FTP Software Windows95, Fujitsu LR450, 460, and 550, Hitachi NR60, IBM AIX, Inria BSD, Linux, HP-UX (SICS), Microsoft Research WindowsNT versions 4 and 5, Sony CSL Apertos IPv4/v6 stack, Sun Solaris, Trumpet Winsock for IPv6, UNH for BSD, NRL for BSD, WIDE KAME and ZETA for BSD, and WIDE v6d.

Several new Windows implementations that will operate under Windows95/98/NT are under way.

Transition from IPv4 to IPv6—A Seamless Approach

IPv6 is unlikely to become the Internet network-layer protocol of choice unless there is literally no choice to be made by the end user, little effort by network and system administrators, and it can operate alongside IPv4 for the indefinite future. Therefore, it must be very easy for the private network (your corporate net) and public network (your ISP) operators to equip, enable, and operate IPv6, while operating IPv4, in such a way that the user doesn't notice that IPv6 is there at all.

A system administrator, but not the user, must be conscious of IPv6 in a minimal sense. It is just another protocol stack that any Internet-based applications will operate over if the system is configured and distributed to do so by the system administrator.

At the network operator level, IPv6 is just another routing stack that can easily be turned on in the site's and ISP's routers (many sites certainly support IPX, AppleTalk, DECnet,...). IPv6 interdomain routing can be operated just like IPv4s because it uses BGP4+.

With the aid of the new *Dynamic DNS Registration Protocol* and IPv6's Stateless Autoconfiguration, users can boot up their system after it has been enabled with an IPv6 stack, in addition to its IPv4 stack, and become IPv6-ready without being aware of it at all. The system would automatically be configured with an IPv6 address, have itself registered automatically in the DNS with the host's existing name alongside its new IPv6 address (in addition to its DNS IPv4 address registration), and when finding a remote host with IPv6, start talking IPv6—all this without the user being required to consciously take action.

Early Production IPv6 Networks

In October of 1998, the 6REN initiative, was established by the U.S. Energy Sciences Network (ESnet). The 6REN is a voluntary coordination initiative of *Research and Education Networks* (RENs) that provide production IPv6 transit service to facilitate high quality, high performance, and operationally robust IPv6 networks.

The first participants were ESnet (the U.S. Dept. of Energy's Energy Sciences Network), Internet2 (the advanced Internetworking development collaboration comprised of many large U.S. research universities), CANARIE (the Canadian joint government and industry initiative for advanced networking), vBNS (the MCI network for NSF advanced networking) and WIDE (the Japanese research effort to establish a "Widely Integrated Distributed Environment").

Other profit and not-for-profit networks worldwide have been invited to join the 6REN. It is expected that during 1999 a sizable production environment capable of advanced demonstrations and deployment of Internet applications over IPv6 networks will be in place.

The Future for IPv6

It is too early to predict with total certainty that the Internet will adapt to the use of the IPv6 protocol. However, it should be obvious that IPv6 offers many important features for a next-generation Internet: automatic configuration, greatly expanded addressing, easy site renumbering, built-in security, and more.

One possible scenario for IPv6 is where it becomes the protocol of choice for newer applications not currently using Internet technology; for example, controlling traffic lights, reading electric meters, and so on. In these uses, IPv6 does not require coexistence with IPv4 because some form of gateway function would provide interconnection to the current Internet.

Another scenario (which doesn't exclude the previous one) is that Microsoft provides IPv6 support for a future version of Windows Networking on Windows OS, and promotes it within corporate America for its better features in supporting advanced corporate application/networking needs. In this scenario, the Internet will learn to carry IPv6 somehow, even if it is via automatically created tunnels that operate over IPv4 (somewhat similar to the 6bone's tunneling, but with dynamic creation of the tunnels as needed). It is expected that after Microsoft ships IPv6 and large corporations begin using it, ISPs will deploy IPv6 to get their business.

Yet another possibility is that the Internet telephony revolution will come to the conclusion that only IPv6 can provide cost-effective, scalable, end-to-end worldwide telephony implementations. This may be even more important as new classes of wireless networked devices, for example, PDAs and PCS phones, are integrated and built in very large volume.

Also, in parts of Asia and China, where there is little Internet connectivity at present, and very few IPv4 addresses assigned, IPv6 may become very popular because it will allow rapid growth without concerns about address space.

The probability is high that not just one of the above scenarios will happen, but that all will occur, in addition to others not yet imagined.

Whatever the implementation scenario, the probability that IPv6 will augment IPv4 as a part of the Internet of the future is very high!

References

- [1] IPng and IPv6 information, including formal specifications can be found at: <http://playground.sun.com/pub/ipng/html/>
- [2] 6BONE information, including diagrams, hookup info, and registry access is at: <http://www.6bone.net>
- [3] An IEEE EUI-64 overview can be found at:
<http://standards.ieee.org/db/oui/tutorials/EUI64.html>
- [4] “Internet Protocol, Version 6 (IPv6) Specification,” RFC 2460, December 1998.
- [5] “Neighbor Discovery for IP Version 6 (IPv6),” RFC 2461, December 1998.
- [6] “IPv6 Stateless Address Autoconfiguration,” RFC 2462, December 1998.
- [7] “Internet Control Message Protocol (ICMPv6) for the Internet Protocol Version 6 (IPv6),” RFC 2463, December 1998.
- [8] “IP Version 6 Addressing Architecture,” RFC 2373, July 1998.
- [9] “An IPv6 Aggregatable Global Unicast Address Format,” RFC 2374, July 1998.
- [10] “DNS Extensions to support IP version 6,” RFC 1886, December 1995.
- [11] “Proposed TLA and NLA Assignment Rules,” RFC 2374, December 1998.
- [12] “Transition Mechanisms for IPv6 Hosts and Routers,” RFC 1933, April 1996.
- [13] “Router Renumbering for IPv6,” Internet Draft, Work in Progress, **draft-ietf-ipngwg-router-renum-06.txt**, November 1998.
- [14] *IPv6: The New Internet Protocol*, Christian Huitema, ISBN 0-13-850505-5, Prentice Hall, 1998.
- [15] *IPng: Internet Protocol Next Generation*, Edited by Scott O. Bradner and Allison Mankin, ISBN 0-201-63395-7, Addison-Wesley, 1996.

ROBERT FINK is a network researcher with ESnet (the U.S. Dept. of Energy’s Energy Sciences Network) at the Berkeley Lab (the Ernest Orlando Lawrence Berkeley National Laboratory). He is cochair of the IETF ngtrans (IPng Transition) Working Group, and leads the 6bone project. You can reach him at: fink@es.net

Secure E-Mail: Problems, Standards, and Prospects

by Marshall T. Rose and David Strom

As we spend more and more time using e-mail, most of us eventually find that we need to be able to prove our identity to our correspondents and secure the contents of our messages so that others can't view them readily. Proving your identity is called *authentication*. In the physical world, this is accomplished by photo identification, such as a driver's license, passport, or corporate identity card. When the time comes to prove who you are (for example, before a major purchase), you show your card. Your appearance and signature match the photo and signature on your card, and the purchase is made.

On the Internet, however, the process isn't as easy. Does e-mail from **sidney@example.com** really originate from our friend Sidney at the Example Corporation? Maybe it's from someone else, who just happens to be using Sidney's machine when he is out to lunch. Or, worse, someone trying to impersonate Sidney illicitly. And even if the message actually is from the "real" Sidney, how can we be sure: Is there an electronic analog to a signature?

Most of us are trusting individuals; we tend to believe that people are who they say they are unless we have particular reasons to doubt their identity. But on the Internet, we have to look beyond face value. And proving that someone indeed did send a particular message is a very difficult problem.

This may be one of the main reasons why corporations employ Lotus Notes and other Internet-based messaging systems that are not 100-percent pure. They want to ensure that all messages carry the appropriate authentication with them at all times. In order for new users of Notes to start using the software, they must first obtain an electronic certificate that authenticates them to the system. The certificate is created by the Notes system administrator, who works in conjunction with that particular Notes server owned by that particular corporation.

Securing the message contents is also a challenge: all e-mail sent over the Internet, unless otherwise protected, is sent in clear ASCII text. If you have the tools, the time, and the technical expertise, you can capture this traffic and read anyone's correspondence. It isn't simple, but it is quite possible.

Besides being sent as clear text, e-mail can also be intercepted and its contents changed between the time the sender composes the message and the recipient reads it. Again, this task is neither likely nor simple, but it can be accomplished if someone is determined enough to do it. Therefore, senders can neither prove nor deny that they sent a particular message to you; it could be real or a forgery, and you have no way of knowing which.

Cryptography Standards

It would be great if we could say that the future for secure e-mail is bright, and that there will be standards in place that will help. However, the state of secure e-mail standards for the Internet is best described as a “terrible mess”! (Ed.: a less charitable phrase is used in the book from which this material is adopted.) Think that characterization is unprofessional? It is actually quite detached, considering the amount of culpability enjoyed by the principals of the Internet’s secure e-mail debacle. We would love to write an article describing the high crimes and misdemeanors of these scoundrels, but that would only publicize the guilty, not punish them. So, instead we’ll survey the horizon and try to make sense of what little terrain there is.^[1]

In brief, no technologies for secure e-mail in the Internet meet all of the following criteria:

- Multivendor
- Interoperable
- Approved or endorsed by the Internet’s standardization body

There are two competing technologies, each of which satisfies at most one of these criteria. However, for any 100-percent-pure Internet solution to succeed, we feel it must be based on technologies that satisfy all three.

Basic Concepts

In order to understand secure e-mail, you need to know only three concepts:

- Data encryption (privacy)
- Message integrity (authentication)
- Key management

Everything else is a matter of data formats.

Data Encryption

When the contents of a message are to be protected from third-party disclosure, it is necessary to agree upon an encryption algorithm. Because cryptographic algorithms are constantly being scrutinized, a secure e-mail standard must be extensible with respect to the algorithms that it allows.

Historically, *symmetric encryption algorithms* are used for this purpose. A symmetric algorithm is one in which the same key is used to both encrypt and decrypt the data. Symmetric algorithms are chosen because they are computationally less burdensome (in other words, faster to execute) than asymmetric algorithms.

As such, each time a message is to be encrypted, a new session key is generated for that purpose. Although one could send the session key via some secure path, it is easier to include the session key along with the message, but encrypted so that only the intended recipient can decipher it. Upon deciphering the session key, the recipient can apply the encryption algorithm and retrieve the original contents.

For example, Network Associates' Pretty Good Privacy (PGP), one of the two technologies we'll examine, uses an asymmetric algorithm to encrypt the session key and a symmetric algorithm to encrypt the user's data.

Message Integrity

When the contents of a message are to be verified as authored by a particular user and unaltered by any other user, it is necessary to agree upon a *signature* and *hash algorithm*. The former is used to verify the authenticity of the message, and the latter is used to verify the integrity of the message. Again, any secure e-mail standard must be extensible with respect to the algorithms that it uses for these purposes.

For signature algorithms, asymmetric algorithms are typically used. These algorithms utilize a public key and a secret key. A signature algorithm combined with a secret key allows someone to generate a digital signature for the contents of a message. A signature algorithm combined with a public key allows someone to verify the digital signature for a message. As you might expect, signature algorithms are one-way functions: You can't reconstruct the input to a signature function by looking at its output.

Hash algorithms are often called *message digest algorithms*. They simply compute a checksum on their input; no keys are involved. Hash algorithms are also one-way functions, and a good hash algorithm is one in which very similar inputs produce dramatically different outputs. Hence, if even a single bit is altered or corrupted in transit, the hash value will be different.

Key Management

All discussion now hinges on how keys are used for asymmetric algorithms. Specifically, how do you trust the identity of the secret key used to make a digital signature? To start, we have to introduce the notion of a *public key certificate*. Although the actual formats vary, at its heart a certificate contains three things:

- The identity of the "owner" of the certificate
- A public key
- Zero or more guarantees to the validity of the binding between the identity contained in the key and the owner in the "real world"

So, the next step is to ask what these identities and guarantees look like. Unfortunately, we now enter the realm of sociology rather than technology. The only theoretical limitation on an identity is that you have to be able to represent it digitally. It could be a name (for example, “Jim Bidzos”) or an e-mail address (for example, `prz@pgp.com`) or a key in some database (for example, the name of an object in a directory). More interesting examples could include a series of assertions (for example, your driver’s license number is this, your passport number is that, and so on).

Fortunately, the guarantees are a bit simpler to describe—they are digital signatures from other public keys that vouch for the veracity of the binding. For example, if you encountered a public key certificate in which the identity was someone’s passport number, it would be natural to expect that the certificate contains a digital signature from the government entity (or its agent) that issued the passport. However, this begs another question: Why should you trust the entities that have signed someone’s public key? It turns out that our two contending technologies have different answers to that question.

As you might expect, certificates have some additional properties, such as a date the certificate becomes valid, the date the certificate expires, and a “fingerprint.” The fingerprint is simply a hash of the identity and public key so you can tell if it has been altered in transit.

Finally, *certificate revocation lists* identify certificates that are no longer valid. For example, if the secret key associated with a certificate is accidentally disclosed, then the corresponding certificate is revoked.

Pretty Good Privacy: The Web of Trust

Pretty Good Privacy (PGP) is encryption for the masses. Despite the fact that it required a couple of complete rewrites in order to achieve stability, it gets the job done.

An effort is under way to provide a “standards-based” version of the PGP technology, termed *OpenPGP*. The “pre-standards” version of PGP uses the RSA algorithm for signatures and the IDEA algorithm for encryption. The version being developed is more flexible with respect to the algorithms it supports.

The most remarkable thing about PGP is its trust model. Remember the earlier question: How do you know whether you should believe the identity in a public key certificate? To answer this in the context of PGP, each user assigns two attributes to the PGP certificates that they encounter: *trust* and *validity*. Trust is a measure as to how accurate the certificate’s owner is with respect to signing other certificates. Validity indicates whether or not you think the identity in the certificate refers to the certificate’s owner.

So, initially your local collection of certificates starts out with one—your own PGP certificate. You then sign your friend’s certificate and he or she signs yours. Because you trust yourself when signing those certificates, your friend’s certificates are automatically considered valid. Then, based on your judgment of your friend’s abilities to sign other certificates accurately, you assign a level of trust to his or her PGP certificates. As you receive messages containing other people’s certificates, if they are signed by you, or any of your trustworthy friends, they are automatically deemed valid. This organic, highly decentralized approach toward validating public key certificates is termed the *web of trust*.

Key servers are also available that are repositories of PGP certificates. If you need to send e-mail to someone, but don’t have his or her certificate, you can query a server to see if a copy is there. Of course, the usual rules apply with respect to assigning trust and validity—it’s up to you! Key servers also help when you receive e-mail from someone new. Although the message will contain a copy of someone’s PGP certificate, you may not know about any of the signatories. So, you can go to a key server and fetch the certificates for the signatories; you might decide to trust them after seeing who signed their certificates.

We’ve simplified the web of trust in that validity isn’t “all or nothing,” as we implied previously. Rather, PGP offers a flexibility spectrum of possibilities; for example, requiring two trustworthy signatories before considering a certificate to be valid. But the one thing that should be clear is that trust and validity are *different*. You will probably have many keys in your local collection of certificates that are considered valid, but probably only a few of those will be considered authorized to vouch for others.

Secure MIME: The Hierarchy of Trust

There is an interesting concept in advertising called “ambush marketing.” The basic idea is that your advertising campaign leverages off the brand and promotion of a competitor. *Secure Multipurpose Internet Mail Extensions*, or S/MIME, is an example of ambush marketing in the Internet. Although MIME is an Internet standard, which has been implemented by hundreds of vendors and provisioned in tens of thousands of networks, S/MIME is the product of a closed vendor consortium.

S/MIME has two versions: version 2 and version 3. As of this writing, products that claim to implement S/MIME implement version 2. They use the RSA algorithm for signatures and a weak algorithm for encryption (RC2 with 40-bit keys). An effort is under way to provide a “standards-based” version of the S/MIME technology—version 3. The version being developed is more flexible with respect to the algorithms it supports. S/MIME uses a hierarchical model for establishing trust. For example, if your employer assigns you an S/MIME certificate, he will act as a certification authority and sign that certificate. As a consequence, trust is established on the basis of a hierarchical relationship between the *subject* of a certificate (the identity) and the *issuer* (the signatory).

This model has some strengths: users rely on the certification authorities implicitly. However, a bootstrapping problem still exists: How do you know to trust the issuer? The answer is that your local collection of certificates also has some “top-level” certificate authorities, and it is these authorities that sign the public key certificates of the issuers. If the hierarchy of trust can be kept to one or two levels, this is manageable in practice.

The web and hierarchical models of trust share many attributes in common. For example, when you receive a message, it contains a copy of the certificate that was used to make the digital signature. If you aren’t familiar with the signatories, you can look in a remote repository of keys. The only difference between the two models here is that the hierarchical model needs key servers to make its key infrastructure work. Because of this, keys are usually stored in a directory service accessed via the *Lightweight Directory Access Protocol* (LDAP).

Data Formats

The **multipart/encrypted** and **multipart/signed** contents are used to convey secure e-mail. Fortunately, they are both very simple content types.

A **multipart/signed** content has two subordinate body parts. The first contains the data that is being authenticated and can be any MIME content type (**text/HTML**, **multipart/mixed**, and so on). The second contains the digital signature used to authenticate the content. The **multipart/signed** content has two mandatory parameters. The *protocol* parameter defines the technology used to generate the digital signature, and the *micalg* (for “MIC algorithm”) parameter defines the hashing algorithm used (for “MIC” read: *message integrity check*). The value of the protocol parameter is also the content type used for the second body part. The only tricky part is that the digital signature is calculated on the data before a transfer encoding, if any, is applied.

Let’s make this a little more concrete. If we assume that the OpenPGP effort produces an Internet standard based on the current draft (a reasonable assumption at 50,000 feet), then the structure of a **multipart/signed** message created using PGP technology would look like the following:

- The protocol parameter would be **application/pgp-signature**
- The micalg parameter would be **pgp-md5**
- The first body part would be labeled as whatever you wanted to sign
- The second body part would be labeled as **application/pgp-signature**

The second body part, a data structure defined by the OpenPGP document, contains the digital signature along with any supporting material (for example, a copy of the sender’s PGP certificate).

Note that you don't encrypt the first body part in a **multipart/signed** content. In this way, if only some of your recipients have secure e-mail, but you still want to sign it for those who do, everyone can still read the first body part.

A **multipart/encrypted** content has two subordinate body parts. The first contains the information needed to decipher the encrypted data (for example, the encrypted session key along with an indication as to the certificate needed to decipher the session key). The second contains the encrypted data, labeled as **application/octet-stream**. The **multipart/encrypted** content has one mandatory parameter, **protocol**, which defines the technology used to encrypt the data. The value of the **protocol** parameter is also the content type used for the first body part.

To further define this concept, if we use OpenPGP as the basis for a hypothetical example, then the structure of a **multipart/encrypted** would look like the following:

- The **protocol** parameter would be **application/pgp-encrypted**
- The first body part would be labeled as **application/pgp-encrypted**
- The second body part would be labeled as **application/octet-stream**. In practice, the input to the encryption algorithm would be **multipart/signed**.

Finally, one or more MIME content types might be defined for sending certificates, certificate revocation lists, and so on. These are all specific to the particular secure e-mail technology being used.

Encrypting Your Messages

If we look at popular commercial e-mail products, many of them include support for some kind of encryption. Both Microsoft's Outlook Express and Netscape Messenger include support for S/MIME, although we'll see in a moment that the two have radically different capabilities. And Qualcomm's Eudora Pro package comes with an add-on module for supporting PGP, which you may or may not have installed when you installed the software. In order to encrypt a message, you need to go through the following process:

1. Choose which of the two competing technologies (and specific e-mail software) you wish to use for your encrypted correspondence. Both methods have advantages and disadvantages.
2. Choose whether you want to just digitally sign your messages or encrypt their entire contents, or both.
3. Either choose an enterprise certificate authority and set up the appropriate server software, or obtain a certificate from a public authority. Again, both methods have advantages and disadvantages.
4. Enroll with this certificate authority and obtain an encryption certificate or key for a particular machine and a single e-mail address.

5. Exchange keys with your correspondents, and manage where these keys are stored on your machine.
6. Encrypt and decrypt messages.

If this process seems rather involved and complex, it is. The process is not nearly where it should be to enable encryption to be useful by most e-mail users, and won't be for some time. If all of this seems overwhelming to you, we certainly understand.^[2] It is to us, too! But let's go through these six steps in more detail.

PGP vs. S/MIME

Our discussion in the standards section might have convinced you that encryption technology is still very much a work in progress, and after you begin to use the encryption features of your own e-mail software, you'll be further convinced. Nevertheless, unless you plan to test lots of different software products, you should first decide on which product and which encryption technology you intend to use. You definitely want to limit yourself to as small a universe as possible, because running more than one e-mail software product will only make your encryption life miserable. So which to choose?

PGP is everyman's product. It was designed for single individuals to use and still remains the easiest method to set up and get going, although it is far from simple. The version of PGP that comes with the Eudora Pro box is the individual version; a separate and more capable version is available for workgroups or businesses, called *PGP for Business Security*. This business version is the one we recommend, even if you are the only person in your corporation that will use encryption. You'll find that after you start, others will follow, and you might as well start off with the more capable version.

If you want to use PGP, you will need to run a separate piece of software to encrypt and decrypt your messages. If you already use software such as Messenger or Outlook Express, that is certainly more cumbersome than using the built-in S/MIME features of those two products.

In 1999, PGP is more capable than S/MIME when it comes to setting up an enterprise encryption policy and putting it into practice on a daily basis. For example, with PGP you could establish that all outgoing and incoming encrypted messages are first copied to a special archive, and that all outgoing messages are encrypted with a special administrator's key that can be used in an emergency to read the message if the sender forgets his key or leaves the company. S/MIME doesn't have this ability yet, although this feature is being developed for the future.

PGP is a single-vendor solution: All your software must eventually come from Network Associates to run the various certificate servers and encryption modules. With S/MIME, you'll have some degree of choice, although we found that in practice you probably want to make use of

the same e-mail product when exchanging encrypted messages if you want them to be read with a minimum of difficulty. Not all S/MIME packages can exchange encrypted messages with each other because of differences in their implementations. When Dan Backman of *Network Computing* magazine tested five different products, he found several that couldn't read messages sent by others.^[3]

Part of the problem with S/MIME is the various choices of "strength" of cryptographic algorithms that are in use in today's browsers and e-mail software. This debate is more about politics than technology, because the U.S. government places restrictions on various algorithms, as mentioned earlier. Two different parameters are of interest: the length of the key itself used in any certificate and the type of encryption technology used. Netscape software supports key lengths ranging from 512 to 1024 bits, for example. In addition, several choices are available for encryption technology; they are labeled RC2 (which can either be 40-bit encryption, the only one allowed for export by the U.S. government, or more complex encryption of 64, 128, or even 255 bits), and *Data Encryption Standard* (DES). RSA, Inc., developed RC2. On the other hand, the U.S. government developed DES. Debate abounds as to which is the better or more or less proprietary technology.

These details are outside the scope of this article, but you should know that the larger the key size and encryption algorithm, the more difficult it is for someone to decode an intercepted message.

Digital Signature Required?

Your next choice is to consider whether to just make use of a digital signature, to encrypt the entire message, or to make use of both technologies. All encryption products can do both, but in somewhat different ways.

Digital signatures guarantee that your recipients have received your message without any tampering and that they can trust that the message came from you. The actual message body, and any attachments, arrive without any encryption, meaning that someone could still capture this traffic and read your correspondence. You might want to use a digital signature without encrypting the message, if you care that your message was received intact and that your correspondents can know that you sent it.

There are two different types of signed messages: *clear* and *opaque*. With clear-signed messages, you can still read the message text, even if you don't have any encryption functions in your e-mail software. The signature is carried along with the message in a separate MIME portion of the message from the message body, which remains untouched and still readable. This feature can be handy, especially if you correspond with many people and they probably haven't adopted any particular encryption product, or if they are using older versions of e-mail software that don't support encryption. Clear signing is also useful in circum-

stances where your encryption technology isn't compatible with your correspondents' technology. PGP supports only clear signing in its products.

One problem with clear signing is e-mail gateways. They often will break the encryption of the signature, because they will either add or remove characters from the message, and that sloppiness could invalidate the signature block. After all, part of the role of the signature is to ensure that the message was delivered intact and unaltered!

Opaque signing means that your recipients will get a blank message if they aren't running any encryption software, or if their encryption software doesn't work with yours. Opaque signing wraps the entire message in a Base64 encoding, which is usually left alone by most e-mail gateways. This encoded message then gets transmitted and then decoded by the S/MIME recipient.

PGP places its signature inside the encrypted envelope when it sends messages, making it difficult to determine the signature of such a message until you first decrypt it. The PGP producers claim that this feature offers extra protection in case the message is compromised or copied en route. Newer versions of PGP offer a MIME option that places the signature outside the encrypted envelope. This is how S/MIME products work, making it easier to determine who sent it.

Choose Your Certificate Authority

Now you have another decision to face, and that is how to set up what is called the *certificate authority* (CA) for your enterprise. This software runs on a UNIX or NT server and manages the keys or certificates of everyone in your corporation. It serves as a central place of trust and signs all of your users' certificates. If you trust your CA, in theory you should be able to trust the certificates that are signed by the CA, called *inherited trust*.

The problem is that there isn't any "central" CA for the entire universe of e-mail users. Although there are several public CAs that anyone can use, either for free or for a fee, they don't necessarily trust each other, nor should they. What happens if an employee of VeriSign becomes disgruntled and starts issuing bad certificates? There should be checks and audits to ensure that these types of problems can't undermine the entire CA system, just as there are checks and audits to prevent rogue banking employees from crediting their own accounts.

Setting up a CA is the beginning of setting up a very complex security infrastructure for your enterprise. Your CA needs to establish a link of trust from all your users to the administrator or operator of the CA itself, and from your CA to other CAs with which you communicate.

There are two different kinds of CAs: One uses software that you install on your own server inside your enterprise and you maintain; the other is

public servers. Having your own server places the burden on creating and revoking certificates on your security administrator, or whoever is going to operate the CA server. In many cases, these products can be administered from a Web browser after they are installed, and the servers can handle certificates from a wide variety of S/MIME products, one of the few shining spots on the interoperability scene at the moment.

PGP for Business comes with its own version of a certificate server. It runs on a Windows desktop machine and typically is used by the administrator of the entire security apparatus to handle certificates. It can handle only PGP certificates.

Some popular software products that function as certificate servers are listed below.

Vendor	URL	Product
<i>Enterprise CAs:</i> Netscape Xcert	www.netscape.com www.xcert.com	Certificate Server Sentry CA
<i>Public CAs:</i> VeriSign Thawte	www.verisign.com www.thawte.com	Secure Server ID Public CA

Enroll and Acquire Your Certificate

When you have your certificate authority either in mind or installed, you next have to set up how you want to acquire your own certificate.

You have two broad methods: by Web or by e-mail. Actually, you don't have any choice: If you have picked your e-mail product and CA at this point in the process, you have to use whatever method comes with that choice. Netscape Messenger and Microsoft's Outlook Express, among others, make use of their related Web browsers to enroll certificates, as you might suspect. And other products make use of e-mail to send and enroll certificates. For example, Xcert's Sentry CA sends you a message telling you that your certificate has been granted, but in the e-mail it has URLs for both Communicator and Internet Explorer where you can download the certificate and place it inside the appropriate software. Why two different links? Because each product supports a different way of acquiring certificates, of course. So much for standards.

Exchange and Manage Certificates

Now comes the hard part—dealing with the certificates of your correspondents, and managing both theirs as well as other certificates around your corporation.

As we mentioned in our standards section, you need to exchange certificates with your correspondents before you can begin to exchange encrypted e-mail. And that means sending your public key to them, and getting their public keys from them, before you can exchange actual encrypted messages. If you are corresponding with someone who doesn't have the same CA in common, you'll first need to establish a trust relationship and exchange root CA certificates before you can exchange the individual certificates. This is somewhat painful, but when you get the hang of it, it isn't that difficult.

After you begin to exchange more than a few of these certificates, you might think that this is a job for a directory server, and, thankfully, the vendors are already there. The CA server can set up entries in an LDAP directory to keep track of who is issued a certificate, and you can query this LDAP server to find who has them. That is the good news, and indeed the PGP product makes use of its own LDAP server to keep track of its certificates. However, the LDAP server is only used by PGP; if you want a general-purpose LDAP server to keep track of your users, you'll have to install something else.

As a challenge for open systems and interoperability, we installed the Xcert Sentry CA and Netscape's Directory Server on a test network. The Xcert was used to create and manage our certificates for our test corporation, and the entries were placed in the Netscape LDAP directory. We created the certificates using the Netscape browser and stored the information in our Messenger e-mail software. After going through the process described previously, we had a valid certificate and could see it in the SecurityMessenger settings. Although the Sentry CA couldn't automatically deposit a certificate in the Netscape LDAP server, we (operating as the security administrator) could do so with a few simple Web forms and keystrokes. So far, so good.

The challenge was trying to pry these certificates loose using other products, such as Outlook Express. There we ran into trouble, mainly because the Netscape software creates the certificate in a nonstandard place in the LDAP directory. According to the standards documents, the certificate should be placed in a particular spot in the LDAP directory schema, called *usercertificate*. Netscape, for whatever reason, places them at a location called *usersmimecertificate*. This meant that non-Netscape products couldn't view the certificates in our directory, because they were looking in the wrong place.

This brings up a very good point: The connection between a user and his or her certificate is tenuous at best. Just because you know that **david@strom.com** is the e-mail address of David Strom and you have his certificate, it doesn't mean that any of your expensive software tools can make this connection. This situation will create all sorts of headaches for your security administrators, and it means that you need to maintain at least two directories on your own machine—one for users and one for certificates.

It would be nice if the address books of our e-mail software could handle this automatically, but they don't.

That's not the only issue with managing certificates. What if someone leaves the company? Or changes his or her e-mail address? Or if you want to use the same certificate, but on several different machines? Most certificates are tied to a particular machine and a particular e-mail address, meaning that any new address will require a new certificate. Again, we find this situation unacceptable.

Encrypt and Decrypt Messages

Now you can finally go and encrypt your messages. Various options are available in your e-mail software to do this, and you can choose to sign a message as well as to encrypt it.

That is the encryption portion. What about the decryption side? If you have done your homework and exchanged certificates as we discussed earlier, then when you receive your encrypted message, it should automatically decrypt and display in plain text. You shouldn't have to do anything else—unless the encryption system is broken by a gateway or product incompatibility.

Futures

The obvious question is whether the Internet needs two standards for secure e-mail.

Proponents for both sides can make superficially compelling arguments. PGP proponents point to a grassroots constituency and a huge installed base of legacy systems. PGP emphasizes privacy for individuals. S/MIME proponents, on the other hand, point to some major vendors and an emphasis on nonrepudiation.

If history is any judge, the PGP side will win because less infrastructure is required to make it work. S/MIME has to solve all the problems that PGP has to solve, plus a few more. However, these things aren't decided overnight. So, our prediction is rather straightforward: The two sides will compete in the Internet marketplace for a couple of years, but ultimately the game is PGP's to lose. It requires less infrastructure and fewer broad agreements to achieve ubiquity.

Endnotes

- [0] Our thanks to Dan Backman of *Network Computing* magazine for his help in sharing his lab and providing many valuable insights in the preparation of this article. This article is based, in part, on *Internet Messaging: From the Desktop to the Enterprise*, ISBN 0-13-9786100-4 Prentice-Hall, 1998.
- [1] See <http://strom.com/places/smime.html> for details regarding product interoperability testing for encrypted e-mail packages.
- [2] There is an alternative to this process. The United Parcel Service has produced a file transfer utility called NetDox, available at www.netdox.com. It requires special software to be installed on each computer, and it simplifies the certificate and encryption process somewhat. But this is yet another proprietary solution to the encrypted e-mail problem—something we think goes in the wrong direction.
- [3] The article has more in-depth examination of testing MIME interoperability and features of Messenger, Outlook Express, Baltimore's MailSecure, OpenSoft's ExpressMail, and two Worldtalk plug-ins for Eudora and Outlook Express. See "Secure E-Mail Clients: Not Quite Ready for S/MIME Prime Time. Stay Tuned." *Network Computing*, February 1, 1998, techweb.cmp.com/nc/902/902r2.html.

MARSHALL T. ROSE is Chief Technology Officer of MessageMedia Inc. (formerly First Virtual Holdings, Inc.). He is responsible for the design, specification, and implementation of several Internet-standard technologies and is an author of over 60 of the Internet's RFCs, and several books on Internet technologies. He can be reached at mrose@dbc.mtview.ca.us

DAVID STROM is an independent consultant and frequent speaker at NetWorld+Interop shows around the world, where he teaches a class on e-commerce and Web storefronts. He was founding editor-in-chief of *Network Computing* magazine and has written over a thousand articles for various computer trade publications. He is also publisher of the e-mail newsletter *Web Informant*, an almost-weekly series of essays on Web marketing, technology, and culture. He can be reached at david@strom.com

Book Review

IP Multicasting *IP Multicasting: The Complete Guide to Interactive Corporate Networks*, by Dave Kosiur, ISBN 0-471-24359-0 Wiley Computer Publishing, 1998, <http://www.wiley.com/compbooks/kosiur>

There is nothing remarkable about the statement: As technology becomes more affordable, applications once limited to power users find their way to the mainstream desktop. Video streaming, audio streaming, collaborative applications, and videoconferencing are all examples of applications once found exclusively on high-end workstations but now making their way to the mainstream desktop. If widespread deployment of these applications is to occur, we must be prepared to supply a supporting infrastructure.

The use of IP multicasting is gaining popularity, but many of the fundamentals that drive this and other network technologies, such as routing protocols and transport protocols, are still being debated. This book supplies a comprehensive view of the state-of-the-art as well as practical procedures one can follow in order to incorporate multicasting into existing network topologies.

Organization

Chapter 2 presents an introduction to TCP/IP basics and routing. Chapter 3, *The Basics of Multicasting*, addresses three sender-based multicasting protocols (ST-II, XTP, and MTP) and concentrates on IP multicast (a receiver-based multicasting protocol). The book would be much easier to follow if this chapter had been combined with Chapter 6.

Chapter 4, *Multicast Routing Concepts*, Chapter 5, *Multicast Routing Protocols*, and Chapter 6, *Transport Protocols*, constitute the heart of this book.

Beginning with basic concepts of unicast routing and routing algorithms, the author extends the models to deal with the problems of routing multicast data. Tree maintenance techniques form the bulk of Chapter 4.

Chapter 5 covers four multicast routing protocols: *Distance Vector Multicast Routing Protocol* (DVMRP); *Multicast Open Shortest Path First* (MOSPF); *Protocol Independent Multicast* (PIM); and *Core-Based Trees* (CBT). Placing the emphasis on PIM, Kosiur covers both PIM-SM (*sparse mode*) and PIM-DM (*dense mode*). He does a nice job of describing each of the protocols and summarizes each by reviewing its advantages and disadvantages. Finally, the author concludes by examining ways of achieving interdomain routing and protocol interoperability.

In Chapter 6, Kosiur provides an overview of the *Real-Time Transport Protocol (RTP)/Real-Time Transport Control Protocol (RTCP)* and the *Real-Time Streaming Protocol (RTSP)*.

In addition, he discusses a dozen or more multicast protocols, all trying to answer the question: “How is retransmission of lost packets handled?” He classifies the protocol approaches into *receiver-based* or *sender-based*. In my opinion, this is the most interesting problem of multicasting. Answer this question wrong, and you find yourself with a nonscalable network cluttered with acknowledgments (ACKs).

Chapters 4 through 7 all consider delivering Quality of Service and so I was a little surprised to see Chapter 7 devoted to the subject.

Kosiur provides a good introduction to RSVP (*Resource ReserVation Protocol*), but until we see RSVP in wide deployment I would look at the previous three chapters for practical knowledge on the topic. In Chapter 7, and then in Chapter 11 he covers a lot of practical issues concerning Quality of Service, as well as ways to support multicasting over various networks, such as ATM, Frame Relay, and ISDN/dialup networks.

Chapter 9 is a compilation of some free and commercial software packages that use multicasting. Chapter 10 covers *Mbone* (the Multicast backbone), a popular experimental multicasting network. It is arguable that the state of multicasting wouldn’t be where it is today without the Mbone.

A C+

This book rates a C+. Kosiur certainly has an understanding of the material, but his descriptions are neither clear nor concise. Reading this book is difficult, and learning from it even more so, but better organization could turn it into a gem.

—Neophytos Iacovou
University of Minnesota
Academic & Distributed Computing Services
`iacovou@boombox.micro.umn.edu`

Letter to the Editor

I just read the September 1998 issue of *The Internet Protocol Journal* and thoroughly enjoyed it. It was well written with excellent technical detail but more importantly, the contributors wrote in an understandable and organized method. This is not always the norm for good technical resources; so many times it is simply the reprint of a vendor's documentation.

"What is a VPN—Part II," written by Paul Ferguson and Geoff Huston, was a great article which described the various components and methodologies of VPNs. The information and explanation of the Virtual Private Dial Networking implementations, voluntary versus compulsory tunneling, subscriber's perspectives and real world applications clarified my understanding and knowledge on this subject. I also appreciate an article that ends with a conclusion. I have already located Part I of this article and will be reading it soon. There is one comment; it would be interesting to know which vendor when an example is used, regarding specifically the Frame Relay service provider.

The "Reliable Multicast Protocols and Applications" article was useful and informative, including the scaling issues and the information regarding the new reliable multicast protocols. The details of the *Pretty Good Multicast* (PGM) protocol and how it may improve scaling for multicast was very interesting.

The *Gigabit Ethernet* book review written by Ed Tittel was one of the most informative and well structured book reviews that I have read, especially in a smaller publication. Thanks for providing three pages for book reviews in a forty-seven page publication. This review provided all the information that would assist with the determination of purchasing the book or not.

I hope you continue to publish IPJ in hard copy. I do read and gather information from the Web like everyone else, but I prefer a physical copy to carry with me if I am traveling or at my home. Thanks again for a great publication and I can hardly wait for the next issue.

—Joe Brannan

joe.brannan@pepsi.com

Ed.: We appreciate your comments about our publication. Regarding your question about the Frame Relay example, it is our policy to avoid as much as possible any discussion of products, but we encourage readers to contact the authors directly for that kind of information.

We certainly plan to continue the print edition of IPJ. We are also developing a companion Web site (at www.cisco.com/ipj) that will contain additional information such as glossaries, links to other documents, updates, corrections, and so on. Thanks for writing.

—Ole Jacobsen

ole@cisco.com

ICANN Update

Back in the summer of 1997, the Clinton Administration decided that it was time to privatize the remaining Internet functions that were being managed within the federal research establishment, mostly dealing with Internet names and addresses. These functions had been handled very successfully over many years by the *Internet Assigned Numbers Authority* (IANA) under the direction Dr. Jon Postel and his staff at the Information Sciences Institute of the University of Southern California under contract to DARPA. But it was clear from the rapid expansion of the Internet, the emergence of important players on the industry side, and rising controversy over issues such as Network Solutions' monopoly in issuing domain names for **.com**, that change was necessary.

After two major policy papers and months of argumentative debate, the government recognized the *Internet Corporation for Assigned Names and Numbers* (ICANN) as the new body to assume responsibility for these largely technical management functions. Working from plans drawn up by Dr. Postel, his advisors and the Jones Day law firm, ICANN is endeavoring to satisfy the many constituencies that seek a voice in future decisions on Internet naming and addressing. Sadly, Jon died last fall just as his plan was approaching endorsement by the federal government.

The young organization, incorporated at the end of September, 1998, began operation in early November, has an initial Board of nine appointed Directors headed by Chairman Esther Dyson, and an interim President/CEO Mike Roberts. They are responsible for completing organizational details, devising a representation structure for electing their successors, and beginning to deal with a backlog of undone policy work, such as a determination on if, how and when new top level domains (TLDs) will be created. The new Board has Directors from six countries and plans to hold meetings quarterly in locations throughout the world, beginning with Singapore in March, 1999 and Berlin in May, 1999.

Being neither a Congressionally chartered corporation nor an industry trade association, but something in between, ICANN is an international organization that faces a tough political future with many skeptics challenging the notion that the Internet community can successfully govern itself in the important naming and addressing area. But with a startup fund from corporate contributions, Chairman Dyson and President Roberts, both short timers by design, are determined to get ICANN off the ground and into operation in coming months. More information is available at: **www.icann.org**

This publication is distributed on an "as-is" basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or noninfringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Engineering
MCI WorldCom, USA

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Edward R. Kozel, Sr. VP, Corporate Development
Cisco Systems, Inc., USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the Cisco News
Publications Group, Cisco Systems, Inc.
www.cisco.com*

*Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 1999 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-J4
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

Bulk Rate Mail U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol Journal

June 1999

Volume 2, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

F R O M T H E E D I T O R

In This Issue

From the Editor 1

Peering and Settlements 2

Firewalls and Internet
Security 24

Was the Melissa Virus
So Different? 33

Book Review 36

Call for Papers 38

Fragments 39

In this issue, Geoff Huston concludes his two-part article on Interconnection, Peering, and Settlements. Last time Geoff discussed the technical aspects for Internet Service Provider (ISP) interconnection. This time he examines the associated business relationships that arise out of ISP peering arrangements. He also looks at some future directions for the ISP interconnection environment, particularly with respect to Quality-of-Service considerations.

A recurring theme in this journal has been the traditional lack of security in Internet technologies and systems. We have examined several ways in which security has been added at all levels of the protocol stack. This time we look at *firewalls*, a popular way to segregate internal corporate intranet traffic from Internet traffic while still maintaining Internet connectivity. Fred Avolio gives the history of firewalls, their current state, and future directions.

Computer viruses have probably existed for as long as we have had computers. However, the ease with which viruses can be distributed as Internet e-mail attachments has made the problem more prevalent. Recently, the *Melissa* virus achieved some notoriety because of its “self-replication” properties. Barbara Fraser, Lawrence Rogers, and Linda Pesante of the Software Engineering Institute at Carnegie Mellon University examines some of the issues raised by this kind of virus.

This issue is the first anniversary issue of *The Internet Protocol Journal* (IPJ). You can find all of our back issues in PDF format at the IPJ Web site: www.cisco.com/ipj. Please let us know if you have suggestions for articles, books you want to review, or general feedback for this journal. Our contact address is: ipj@cisco.com.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Interconnection, Peering and Settlements—Part II

by Geoff Huston, Telstra

In Part I we examined the business drivers behind the adoption of the exchange model as the common basis of interconnection, and also examined the advantages and pitfalls associated with the operation of such exchanges within the public Internet. (See *The Internet Protocol Journal*, Volume 2, No. 1, March 1999.) In continuing our examination of the technology and business considerations that are significant within the subject of Internet Service Provider (ISP) interconnection, in this part we focus on the topic from a predominately business perspective.

Interaction Financials: Peering and Settlements

Any large multiprovider distributed service sector has to address the issue of cost distribution at some stage in its evolution. Cost distribution is the means by which various providers can participate in the delivery of a service to a customer who purchases a service from a single provider, and providers can each be compensated for their costs in an equitable structure of interprovider financial settlement.

As an example, when an airline ticket is purchased from one air service provider, various other providers and service enterprises may play a role in the delivery of the service. The customer does not separately pay the service fee of each airport baggage handler, caterer, or other form of service. The customer's original fare, paid to the airline, is distributed to other providers who incurred cost in providing components of the total service. These costs are incurred through sets of service contracts, and are the subject of various forms of interprovider financial settlements, all of which are invisible to the customer.

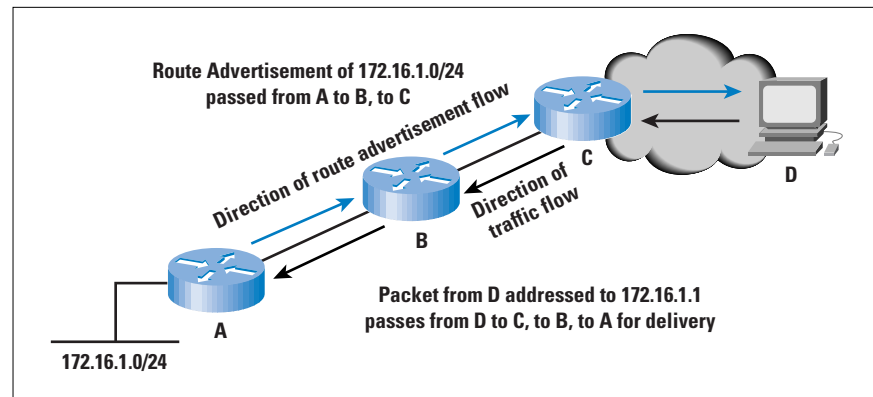
The Internet is in a very similar situation. Some 50,000 constituent networks must interconnect in one fashion or another to provide comprehensive end-to-end service to each client. In supporting a data transaction between two clients, the two parties often are not clients of the same network. Indeed, the two-client service networks often do not directly interconnect, and one or more additional networks must act in a transit provider role to service the transaction. Within the Internet environment, how do all the service parties to a transaction who incur cost in supporting the transaction receive compensation for their cost? What is the cost distribution model of the Internet?

Here, we examine the basis for Internet interprovider cost distribution models and then look at the business models currently used in the interprovider Internet environment. This area commonly is termed *financial settlement*, a term the Internet has borrowed from the telephony industry.

The Currency of Interconnection

What exactly is being exchanged between two ISPs who want to interconnect? In the sense of the meaning of currency as the circulating medium, the question is: What precisely is being circulated at the exchange and within the realm of interconnection? The technical answer to the question is: *routing entries*. When two parties exchange routing entries, the outcome is that traffic flows in response to the flow of routing entries. The route advertisement and traffic flows move in opposite directions, as indicated in Figure 1, and a bilateral routing-mediated flow occurs only when routes are passed in both directions.

Figure 1: Routing and Traffic Flows



Within the routing environment of an ISP there are many different classes of routes, with the classification based predominately on the way in which the route has been acquired by the ISP:

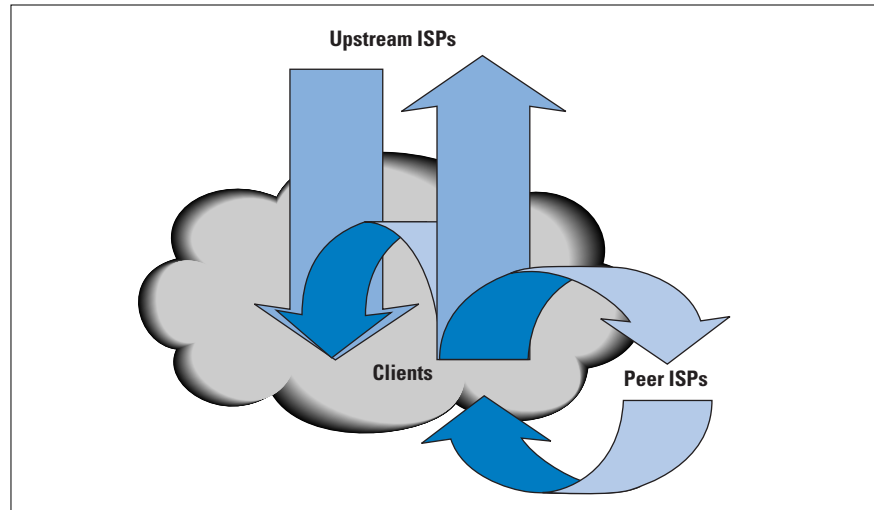
- *Client routes* are passed into the ISP's routing domain by virtue of a service contract with the client. The routes may be statically configured at the edge of the ISP's network, learned by a *Border Gateway Protocol* (BGP) session with the client, or they may constitute part of an ISP pool of addresses that are dynamically assigned to the client as part of the dialup session.
- *Internal ISP routes* fall into numerous additional categories. Some routes correspond to client services operated by the ISP, solely for access to the clients of the ISP, such as Web caches, *Post Office Protocol* (POP) mail servers, and game servers. Some routes correspond to ISP-operated client services that require Internet-wide access, such as *Domain Name System* (DNS) forwarders and *Simple Mail Transfer Protocol* (SMTP) relay hosts. Lastly are internal services with no visibility outside the ISP network, such as *Simple Network Management Protocol* (SNMP) network management platforms.
- *Upstream routes* are learned from upstream ISPs as part of a transit service contract the ISP has executed with the upstream provider.
- *Peer routes* are learned from exchanges or private interconnections, corresponding to routers exported from the interconnected ISP.

How then should the ISP export routes so that the inbound traffic flow matches the outbound flows implied by this route structure? The route export policy is generally structured along the following lines:

- *Clients*: All available routes in the preceding four categories, with the exception of internal ISP service functions, should be passed to clients, either in the form of a *default route* or as *explicit route entries* passed via a BGP session.
- *Upstream providers*: All client routes and all internal ISP routes corresponding to Internet-wide services should be passed to upstream providers. Some clients may want further restrictions placed on their routes being advertised in such a fashion. The ability for a client to specify such caveats on the routing structure, and the mechanism used by the ISP to allow this to happen, should be clearly indicated in the service contract.
- *Peer ISPs*: All client routes and all ISP routes corresponding to Internet-wide service should be passed to peer ISPs. Again the clients may want to place a restriction on such an advertisement of their routes as a qualification to the ISP's own route export policy.

This structure is shown in Figure 2.

Figure 2: External Routing Interaction



The implicit outcome of this routing policy structure is that the ISP does not act in a transit role to peer ISPs and permits neither peer-to-peer transit nor peer-to-upstream transit. Peer ISPs have visibility only to clients of the ISP. From the service visibility perspective, client-only services are not visible to peer ISPs or upstream ISPs, and, therefore, value-added client services are implicitly visible only to clients and only when they access the service through a client channel.

Settlement Options

Financial settlements have been a continual topic of discussion within the domain of Internet interconnection. To look at the Internet settlement environment, let's first look at the use of interprovider financial settlements within the international telephony service industry. Then, we will look at the application of these generic principles to the Internet environment.

Within the traditional telephony model, interprovider peering takes place within one of three general models:

Bilateral Settlements

The first, and highly prevalent, international peering model is that of bilateral settlements. A *call-minute* is the unit of settlement accounting. A call is originated by a local client, and the local client's service provider charges the client for the duration of the entire end-to-end call. The call may pass through, or transit, many providers, and then terminate within the network of the remote client's local provider. The cost distribution mechanism of settlements is handled bilaterally. In the most general case of this settlement model, the originating provider pays the next hop provider to cover the costs of termination of the call. The next hop provider then either terminates the call within the local network, or undertakes a settlement with the next hop provider to terminate the call. The general telephony trunk model does not admit many multiparty transit arrangements. Most telephony settlements are associated with trunk calls that involve only two providers: the originating and terminating providers.

Within this technology model, the bilateral settlement becomes easier, because the model simplifies to the case where the terminating provider charges the originating provider a per-call-minute cost within an accounting rate that has been bilaterally agreed upon between the two parties. Because both parties can charge each other using the same accounting currency, the ultimate financial settlement is based on the net outcome of the two sets of call-minute transactions with the two call-minute termination accounting rates applied to these calls. (There is no requirement for the termination rates for the two parties to be set at the same level.) Each provider invoices the originating end user for the entire call duration, and the financial settlements provide the accounting balance intended to ensure equity of cost distribution in supporting the costs of the calls made between the two providers. Where there is equity of call accounting rates between the two providers, the bilateral interprovider financial settlements are used in accordance with originating call-minute imbalance, in which the provider hosting the greater number of originating call-minutes pays the other party according to a bilaterally negotiated rate as the mechanism of cost distribution between the two providers.

As a side note, the *Federal Communications Commission* of the United States (FCC) asserts that U.S. telephone operators paid out some \$5.6 billion in settlement rates in 1996, and the FCC is voicing the view that accounting rates have now shifted into areas of non-cost-based settings, rather than working as a simple cost distribution mechanism.

This accounting settlement issue is one of the drivers behind the increasing interest in voice-over-IP solutions, because typically no accounting rate settlement component exists in such solutions, and the call termination charges are cost-based, without bilateral price setting. In those cases

where accounting rates have come to dominate the provider's call costs, voice-over-IP is perceived as an effective lever to bypass the accounting rate structure and introduce a new price point for call termination in the market concerned.

Sender Keeps All

The second model, rarely used in telephony interconnection, is that of *Sender Keeps All* (SKA), in which each service provider invoices its originating client's user for the end-to-end services, but no financial settlement is made across the bilateral interconnection structure. Within the bilateral settlement model, SKA can be regarded as a boundary case of bilateral settlements, where both parties simply deem the outcome of the call accounting process to be absolutely equal, and consequently no financial settlement is payable by either party as an outcome of the interconnection.

Transit Fees

The third model is that of transit fees, in which one party invoices the other party for services provided. For example, this arrangement is commonly used as the basis of the long-distance/local access provider interconnection arrangements. Again, this case can be viewed as a boundary case of a general bilateral settlement model, where in this case the parties agree to apply call accounting in only one direction, rather than bilaterally.

Telephony Settlement Trends

The international telephony settlement model is by no means stable, and currently, significant pressure is being placed on the international accounting arrangements to move away from bilaterally negotiated uniform call accounting rates to rates separately negotiated for calls in each direction of a bilateral interconnection. Simultaneously, communications deregulation within many national environments is changing the transit fee model, as local providers extend their network into the long-distance area and commence interconnection arrangements with similar entities. Criticism also has been directed at the bilaterally negotiated settlement rates, because of the observation that in many cases the accounting rates are not cost-based rates but are based on a desire to create a revenue stream from accounting settlements.

Internet Considerations

Numerous critical differences exist between the telephony models of interconnection and the Internet environment; these differences have confounded all attempts to cleanly map telephony interconnection models into the Internet environment.

Internet Settlement Accounting by the Packet

Internet interconnection accounting is a packet-based accounting issue, because there is no "call-minute" in the Internet architecture. Therefore, the most visible difference between the two environments is the replacement of the *call* with the *packet* as the currency unit of interconnection.

Although we can argue that a TCP session has much in common with a call, this concept of an originating TCP call-minute is not always readily identified within the packet forwarding fabric, and accordingly it is not readily apparent that this is a workable settlement unit. Unlike a telephony call, no concept of state initiation exists to pass a call request through a network and lock down a network transit path in response to a call response. The network undergoes no state change in response to a TCP session, and therefore, no means is readily available to the operator to identify that a call has been initiated, and by which party. Of course the use of *User Datagram Protocol* (UDP), and various forms of tunnelling traffic, also confound any such TCP call-minute accounting mechanism.

Packets may be dropped

When a packet is passed across an interconnection from one provider to another, no firm guarantee is given by the second provider that the packet will definitely be delivered to the destination. The second provider, or subsequent providers in the transit path, may drop the packet for quite legitimate reasons, and will remain within the protocol specification in so doing. Indeed, the TCP protocol uses packet drop as a rate-control signal. For the efficient operation of the TCP protocol, some level of packet drop is a useful and anticipated event. However, if a packet is used as the accounting unit in a general cost distribution environment, should the provider who receives and subsequently drops the packet be able to claim an accounting credit within the interconnection? The logical response is that such accounting credits should apply only to successfully delivered packets, but such an accounting structure is highly challenging to implement accurately within the Internet environment.

Packet paths are not predetermined

Packet transit paths can be within the explicit control of the end user, not the provider. Users can exercise some significant level of control of the path a packet takes to transit the Internet if source routing is honored, so that the relative packet flows between two providers can be arbitrarily manipulated by any client, if so desired.

Routing and traffic flow are not paired

Packet forwarding is not a verified operation. A provider may choose to forward a packet to a second provider without reference to the particular routes the second provider is advertising to the first party. A packet may also be forwarded to the second provider with a source address that is not being advertised to the second provider. Given that the generic Internet architecture strives for robustness under extreme conditions, attempts to forward a packet to its addressed destination are undertaken irrespective of how the packet may have arrived at this location in the first place, and irrespective of how a packet with reverse header IP addresses will transit the network.

Comprehensive routing information is not uniformly available

Complete information is not available to the Internet regarding the status and reachability of every possible Internet address. Only as a packet is forwarded closer to the addressed destination does more complete information regarding the status of the destination address become apparent to the provider. Accordingly, a packet may have incurred some cost of delivery before its ultimate undeliverability becomes evident. An intermediate transit provider can never be completely assured that a packet is deliverable.

Settlement Models for the Internet

Where a wholesale or retail service agreement is in place, one ISP is, in effect, a customer of the other ISP. In this relationship, the customer ISP (downstream ISP) is purchasing transit and connectivity services from the supplier ISP (upstream ISP). The downstream ISP resells this service to its clients. The upstream ISP must announce the downstream ISP's routes to all other customers and other egress points of the ISP's networks to honor the service contract to the downstream ISP customer.

However, given two ISPs who interconnect, the decision as to which party should assume the upstream provider role and which party should assume the downstream customer role is not always immediately obvious to either party, or even to an outside observer. Greater geographic coverage may be the discriminator here that allows the customer/provider determination. However, this factor is not the only possible one within the scope of the discussion. One ISP may host significant content and may observe that access to this content adds value to the other party's network, which may be used as an offset against a more uniform customer relationship. In a similar vein, an ISP with a very large client population within a limited geographic locality may see this large client base as an offset against a more uniform customer relationship with the other provider. In many ways, the outcome of these discussions can be likened to two animals meeting in the jungle at night. Each animal sees only the eyes of the other, and from this limited input, they must determine which animal should attempt to eat the other!

An objective and stable determination of which ISP should be the provider and which should be the client is not always possible. In many contexts, the question is inappropriate, given that for some traffic classes the respective roles of provider and client may swap over. The question often is rephrased along the lines of, "Can two providers interconnect without the implicit requirement to cast one as the provider and the other as the client?" Exploration of some concepts of how the question could possibly be answered is illustrative of the problem space here.

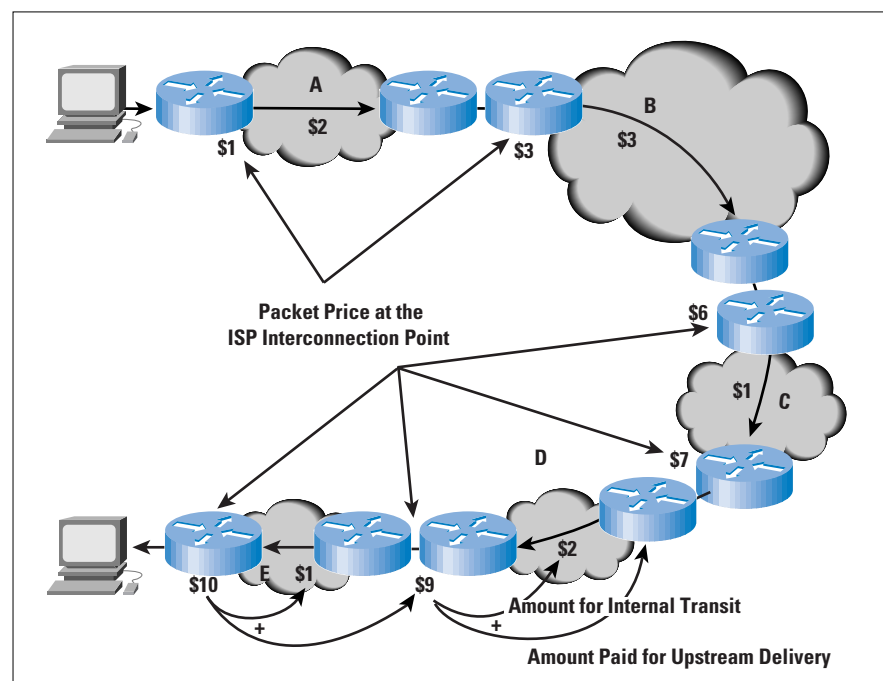
Packet Cost Accounting

One potential accounting model is based on the observation that a packet incurs cost as it passes through the network. For a small interval of time, the packet occupies the entire transmission capacity of each circuit over which it passes.

Similarly, for a brief interval of time, the packet is exclusively occupying the switching fabric of the router. The more routers the packet passes through, and the greater the number and distance of transmission hops the packet traverses, the greater the incurred cost in carrying the packet.

A potential settlement model could be constructed from this observation. The strawman model is that whenever a packet is passed across a network boundary, the packet is effectively sold to the next provider. The sale price increases as the packet transits through the network, accumulating value in direct proportion to the distance the packet traverses within the network. Each boundary packet sale price reflects the previous sale price, plus the value added in transiting the ISP's infrastructure. Ultimately, the packet is sold to the destination client. This model is indicated in Figure 3.

Figure 3: Financial Interprovider Settlement via Packet Cost Accounting



As with all strawman models, this one has numerous critical weaknesses, but let's look at the strengths first. An ISP gains revenue from a packet only when delivered on egress from the network, rather than in network ingress. Accordingly, a strong economic incentive exists to accept packets that will not be dropped in transit within the ISP, given that the transmission of the packet generates revenue to the ISP only on successful delivery of the packet to the next hop ISP or to the destination client. This factor places strong pressure on the ISP to maintain quality in the network, because dropped packets imply foregone revenue on local transmission. Because the packet was already purchased from the previous provider in the path, packet loss also implies financial loss. Strong pressure also is exerted to price the local transit function at a commodity price level, rather than attempt to undertake opportunistic pricing. If the chosen transit price is too great, the downstream provider has the opportunity to extend its network to reach the next upstream

provider in the path, resulting in bypassing the original upstream ISP and purchasing the packets directly from the next hop upstream source. Accordingly, this model of per-packet pricing, using a settlement model of egress packet accounting, and locally applied value increments to a cumulative per-packet price, based on incremental per-hop transmission costs, does allow for some level of reasonable stability and cost distribution in the interprovider settlement environment.

However, weaknesses of this potential model cannot be ignored. First, some level of packet drop is inevitable, irrespective of traffic load. Generally, the more remote the sender from the destination, the less able the sender is to ascertain that the destination address is a valid IP address, and the destination host is available. To minimize the liability from such potential packet loss, the ISP should maintain a relatively complete routing table and accept only packets in which a specific route is maintained for the network. More critical is the issue that the mechanism is open to abuse. Packets that are generated by the upstream ISP can be transmitted across the interface, which in turn results in revenue being generated for the ISP. Of course, per-packet accounting within the core of the network is a significant refinement of existing technology. Within a strict implementation of this model, packets require the concept of an attached value that ISPs augment on an ingress-to-egress basis, which could be simplified to a hop-by-hop value increment. Implementations feasibly can use a level of averaging to simplify this process by using a tariff for domestic transit and a second for international transit.

TCP Session Accounting

These traffic-based metrics do exhibit some weaknesses because of their inability to resist abuse and the likelihood of exacting an interprovider payment even when the traffic is not delivered to an ultimate destination. Of more concern is that this settlement regime has a strong implication in the retail pricing domain, where the method of payment on delivered volume and distance is then one of the more robust ways that a retail provider can ensure that there is an effective match between the interprovider payments and the retail revenue. Given that there is no intrinsic match of distance, and therefore cost, to any particular end-to-end network transaction, such a retail tariff mechanism would meet with strong consumer resistance.

Does an alternative settlement structure that can address these weaknesses exist? One approach is to perform significantly greater levels of analysis of the traffic as it transits a boundary between a client and the provider, or between two providers, and to adopt financial settlement measures that match the type of traffic being observed. As an example, the network boundary could detect the initial TCP SYN handshake, and all subsequent packets within the TCP session could be accounted against the session initiator, while UDP traffic could be accounted against the UDP source. Such detailed accounting of traffic passed across a provider boundary could allow for a potential settlement structure based on duration (*call-minutes*), or volume (*call-volumes*).

Although such settlement schemes are perhaps limited more by imagination in the abstract, very real technical considerations must be borne to bear on this speculation. For a client-facing access router to detect a TCP flow and correctly identify the TCP session initiator requires the router to correctly identify the initial SYN handshake, the opening packet, and then record all in-sequence subsequent packets within this TCP flow against this accounting element. This identification process may be completely impossible within the network at an interprovider boundary. The outcome of the routing configuration may be an asymmetric traffic path, so that a single interprovider boundary may see only traffic passing in a single direction.

However, the greatest problem with this, or any other traffic accounting settlement model, is the diversity of retail pricing structures that exist within the Internet today. Some ISPs use pricing based on received volume, some on sent volume, some on a mix of sent and received volume, and some use pricing based on the access capacity, irrespective of volume. This discussion leads to the critical question when considering financial settlements: Given that the end client is paying the local ISP for comprehensive Internet connectivity, when a client's packet is passed from one ISP to another at an interconnection point, where is the revenue for the packet? Is the revenue model one in which the packet sender pays or one in which the packet receiver pays? The packet egress model described here assumes a uniform retail model in which the receiver pays for Internet packets. The TCP session model assumes the session initiator pays for the entire traffic flow. This uniformity of retail pricing is simply not mirrored within the retail environment of the Internet today.

Although this session-based settlement model does attempt to promote a quality environment with fair carriage pricing, it cannot address the fundamental issue of financial settlements.

Internet Settlement Structures

For a financial settlement structure to be viable and stable, the settlement structure must be a uniform abstraction of a relatively uniform retail tariff structure. This conclusion is critically important to the entire Internet financial settlement debate.

The financial structure of interconnection must be an abstraction of the retail models used by the two ISPs. If the uniform retail model is used, the party originating the packet pays the first ISP a tariff to deliver the packet to its destination within the second ISP; then the first ISP is in a position to fund the second ISP to complete the delivery through an interconnection mechanism. If, on the other hand, the uniform retail model is used in which the receiver of the packet funds its carriage from the sender, then the second ISP funds the upstream ISP. If no uniform retail model is used, when a packet is passed from one provider to the other, no understanding exists about which party receives the revenue for the carriage of the packet and accordingly, which party settles with

the other party for the cost incurred in transmission of the packet. The answer to these issues within the Internet environment has been to commonly adopt just two models of interaction. These models sit at the extreme ends of the business spectrum, where one is a customer/provider relationship, and the other is a peering relationship without any form of financial settlement, or SKA. These models approximately correspond to the second and third models described previously from traditional models of interconnection within the communications industry. However, an increasing trend has moved toward models of financial settlement in a bilaterally negotiated basis within the Internet, using non-cost-based financial accounting rates within the settlement structure. Observing the ISP industry repeat the same well-trodden path, complete with its byways into various unproductive areas and sometimes mistakes of the international telephony world, is somewhat interesting to say the least. Experiential learning is often observed to be a rare commodity in this area of Internet activity.

No Settlement and No Interconnection

Examining the option of complete autonomy of operation, without any form of interaction with other local or regional ISPs, is instructive within this examination of settlement options.

One scenario for a group of ISPs is that a mutually acceptable peering relationship cannot be negotiated, and all ISPs operate disconnected network domains with dedicated upstream connections and no interconnection. The outcome of such a situation is that third-party connectivity would take place, with transit traffic flowing between the local ISPs being exchanged within the domain of a mutually connected third-party ISP (or via transit across a set of third-party ISPs). For example, for an Asian country, this situation would result in traffic between two local entities, both located within the same country, being passed across the Pacific, routed across numerous network domains within the United States, and then passed back across the Pacific. Not only is this scenario inefficient in terms of resource utilization, but this structure also adds a significant cost to the operation of the ISPs, a cost that ultimately is passed to the consumer in higher prices for Internet traffic.

Note that this situation is not entirely novel; the Internet has seen such arrangements appear in the past; and these situations are still apparent in today's Internet. Such arrangements have arisen, in general, as the outcome of an inability to negotiate a stable local peering structure.

However, such positions of no interconnection have proved to be relatively short-lived because of the high cost of operating international transit environments, the instability of the significantly lengthened interconnection paths, and the unwillingness of foreign third-party ISPs to act (often unwittingly) as agents for domestic interconnection in the longer term. As a result of these factors, such off-shore connectivity structures generally have been augmented with domestic peering structures.

The resultant general operating environment of the Internet is that effective isolation is not in the best interests of the ISP, nor is isolation in the interests of other ISPs or the consumers of the ISPs' services. In the interests of a common desire to undertake rational and cost-effective use of communications resources, each national (or regional) collection of ISPs acts to ensure local interconnectivity between such ISPs. A consequent priority is to reach acceptable ISP peering arrangements.

Sender Keeps All

Sender Keeps All (SKA) peering arrangements are those in which traffic is exchanged between two or more ISPs without mutual charge (an interconnection arrangement with no financial settlement). Within a national structure, typically the marginal cost of international traffic transfer to and from the rest of the Internet is significantly higher than domestic traffic transfer. In these cases, any SKA peering is likely to relate to only domestic traffic, and international transit would be provided either by a separate agreement or independently by each party.

This SKA peering model is most stable where the parties involved perceive equal benefit from the interconnection. This interconnection model generally is used in the context of interconnection or with providers with approximate equal dimension, as in peering regional providers with other regional providers, national providers with other national providers, and so on. Oddly enough, the parties themselves do not have to agree on what that value or dimension may be in absolute terms. Each party makes an independent assessment of the value of the interconnection, in terms of the perceived size and value of the ISP and the value of the other ISP. If both parties reach the conclusion that in their terms a net balance of value is achieved, then the interconnection is on a stable basis. If one party believes that it is larger than the other and SKA interconnection would result in leverage of its investment by the smaller party, then an SKA interconnection is unstable.

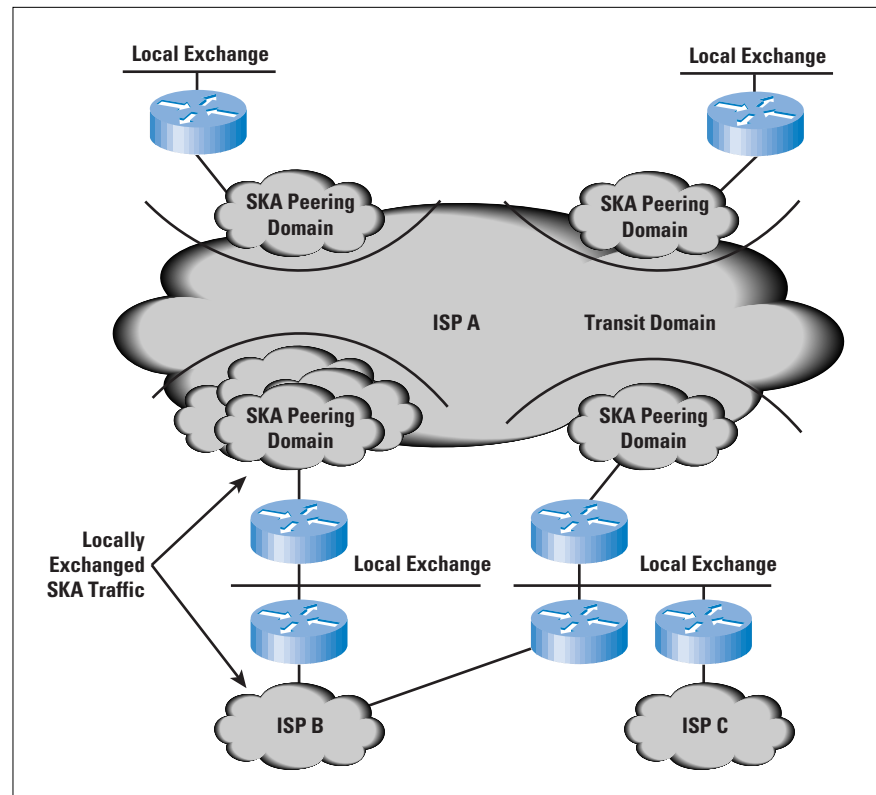
The essential criterion for a stable SKA peering structure is perceived equality in the peering relationship. This criterion can be achieved in many ways, including the use of entry threshold pricing into the peering environment or the use of peering criteria, such as the specification of ISP network infrastructure or network level of service and coverage areas as eligibility for peering.

A typical feature of the SKA peering environment is to define an SKA peering in terms of traffic peering at the client level only. This definition forces each peering ISP to be self-sufficient in the provision of transit services and ISP infrastructure services that would not be provided across a peering point. This process may not result in the most efficient or effective Internet infrastructure, but it does create a level of approximate parity and reduces the risks of leverage within the interconnection. In this model, each ISP presents at each interconnection or exchange only those routes associated with the ISP's customers and accepts only traffic

from peering ISPs at the interconnection or exchange directed to such customers. The ISP does not accept transit traffic destined to other remote exchange locations, nor to upstream ISPs, nor traffic directed to the ISP's infrastructure services. Equally, the ISP does not accept traffic that is destined to peering ISPs, from upstream transit providers. The business model here is that clients of an ISP are contracting the ISP to present their routes to all other customers of the ISP, to the upstream providers of the ISP, and to all exchange points where the ISP has a presence. The particular tariff model chosen by the ISP in servicing the customers is not material to this interconnection model. Traffic passed to a peer ISP at the exchange becomes the responsibility of the peer ISP to pass to its customers at its cost.

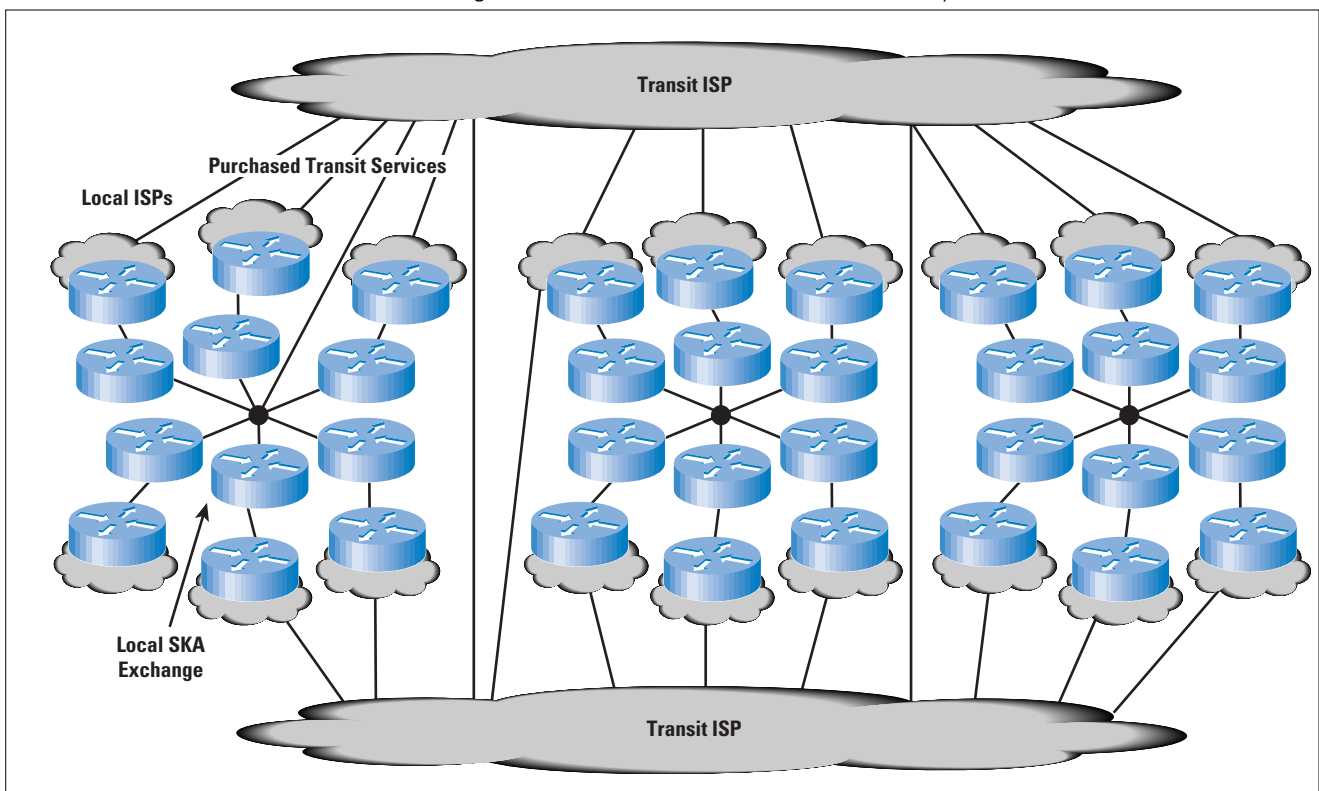
Another means of generating equity within an SKA peering is to peer only within the terms of a defined locality. In this model, an ISP would present routes to an SKA peer in which the routes correspond to customers located at a particular access POP, or a regional cluster of access POPs. The SKA peer's ability to leverage advantage from the greater level of investment (assuming that the other party is the smaller party) is now no longer a factor, because the smaller ISP sees only those parts of the larger ISP that sit within a well-defined local or regional zone. This form of peering is indicated in Figure 4.

Figure 4: SKA Peering
Using Local Cells



The probable outcome of widespread use of SKA interconnections is a generalized ISP domain along the lines of Figure 5. Here, the topology is segregated into two domains consisting of a set of transit ISPs, whose predominate investment direction is in terms of high-capacity carriage infrastructure and high-capacity switching systems, and a collection of local ISPs, whose predominate investment direction is in service infrastructure supporting a string retail focus. Local ISPs participate at exchanges and announce local routes at the exchange on an SKA basis of interconnection with peer ISPs. Such ISPs are strongly motivated to prefer to use all routes presented at the exchange within such peering sessions, because the ISP is not charged any transit cost for the traffic under an SKA settlement structure. The exchange does not provide comprehensive connectivity to the ISP, and this connectivity needs to be complemented with a separate purchase of transit services. In this role, the local ISP becomes a client of one or more transit ISPs explicitly for the purpose of access to transit connectivity services.

Figure 5: ISP Structure of Local and Transit Operations



In this model, the transit ISP must have established a position of broad-ranging connectivity, with a well-established and significant market share of the wholesale transit business. A transit ISP also must be able to present customer routes at a carefully selected set of major exchange locations and have some ability to exchange traffic with all other transit ISPs. This latter requirement has typically been implemented using private interconnection structures, and the associated settlements often are negotiated bilaterally. These settlements possibly may include some element of financial settlement.

Negotiated Financial Settlement

The alternative to SKA and provider/client role selection is the adoption of a financial settlement structure. The settlement structure is based on both parties effectively selling services to each other across the interconnection point, with the financial settlement undertaking the task of balancing the relative sales amounts.

The simplest form of undertaking this settlement is to measure the volume of traffic being passed in each direction across the interconnection and to use a single accounting rate for all traffic. At the end of each accounting period, the two ISPs would financially settle based on the agreed accounting rate applied to the net traffic flow.

Which way the money should flow in relationship to traffic flow is not immediately obvious. One model assumes that the originating provider should be funding the terminating provider to deliver the traffic, and therefore, money should flow in the same direction as traffic. The reverse model assumes that the overall majority of traffic, is traffic generated in response to an action of the receiver, such as web page retrieval or the downloading of software. Therefore, the total network cost should be imposed on the discretionary user, so that the terminating provider should fund the originating provider. This latter model has some degree of supportive evidence, in that a larger provider often provides more traffic to a smaller attached provider than it receives from that provider. Observation of bilateral traffic flow statistics tends to support this, indicating that traffic-received volumes typically coincide with the relative interconnection benefit to the two providers.

The accounting rate can be negotiated to be any amount. There is a caveat on this ability to set an arbitrary accounting rate, because where an accounting rate is not cost-based, business instability issues arise. For greater stability, the agreed settlement traffic unit accounting rate would have to match the average marginal cost of transit traffic in both ISP networks for the settlement to be attractive to both parties. Refinements to this approach can be introduced, although they are accompanied by significant expenditure on traffic monitoring and accounting systems. The refinements are intended to address the somewhat arbitrary determination of financial settlement based on the receiver or the sender. One way is to undertake flow-based accounting, in which the cost accounting for the volume of all packets associated with a TCP flow is directed to the initiator of the TCP session. Here, the cost accounting for all packets of a UDP flow is directed to the UDP receiver. The session-based accounting is significantly more complex than simple volume accounting, and such operational complexity would be reflected in the cost of undertaking such a form of accounting. However, asymmetric paths are a common feature of the inter-AS environment, so that it may not always be possible to see both sides of a TCP conversation and perform an accurate determination of the session initiator.

Another refinement is to use a different rate for each provider, where the base rate is adjusted by some agreed size factor to ensure that the larger provider is not unduly financially exposed by the arrangement. The adjustment factor can be the number of Points of Presence, the range of the network, the volume carried on the network, the number of routes advertised to the peer, or any other metric related to the ISP's investment and market share profile. Alternatively, a relative adjustment factor can simply be a number, without any basis in a network metric, to which both parties agree.

Of course, such a relative traffic volume balance is not very robust either, and the metric is one that is vulnerable to abuse. The capability to adjust the relative traffic balance comes from the direct relationship between the routes advertised and the volume of traffic received. To reduce the amount of traffic received, the ISP reduces the number of routes advertised to the corresponding peer. Increasing the number of routes, and at the same time increasing the number of specific routes, increases the amount of received traffic. When there is a rich mesh of connectivity, the primary objective of routing policy is no longer that of supporting basic connectivity, but instead the primary objective is to maximize the financial return to the operator. If the ISP is paying for an "upstream" ISP service, the motivation is to minimize the cost of this contract, either by maximizing the amount of traffic covered under a fixed cost, or minimizing the cost by minimizing the traffic exchanged with the upstream ISP. Where there is a financially settled interconnection, the ISP will be motivated to configure its routing policies to maximize its revenue from such an arrangement. And of course an ISP will always prefer to use customer routes wherever possible, as a basic means of maximizing revenue into the operation.

Of greater concern is the ability to abuse the interconnection arrangements. One party can generate and then direct large volumes of traffic to the other party. Although overt abuse of the arrangements is often easy to detect, greed is a wonderful stimulant to ingenuity, and more subtle forms of abuse of this arrangement are always possible. To address this, both parties would typically indicate in an interconnection agreement their undertaking not to indulge in such forms of deliberate abuse.

Notwithstanding such undertakings by the two providers, third parties can still abuse the interconnection in various ways. Loose source routing can generate traffic flows that pass across the interconnection in either direction. The ability to remotely trigger traffic flows through source address spoofing is possible, even where loose source routing is disabled. This window of financial vulnerability is far wider than many ISPs are comfortable with, because it opens the provider to a significant liability over which it has a limited ability to detect and control. Consequently, financial settlement structures based on traffic flow metrics are not a commonly deployed mechanism, because they introduce significant financial risks to the ISP interconnection environment.

The Settlement Debate

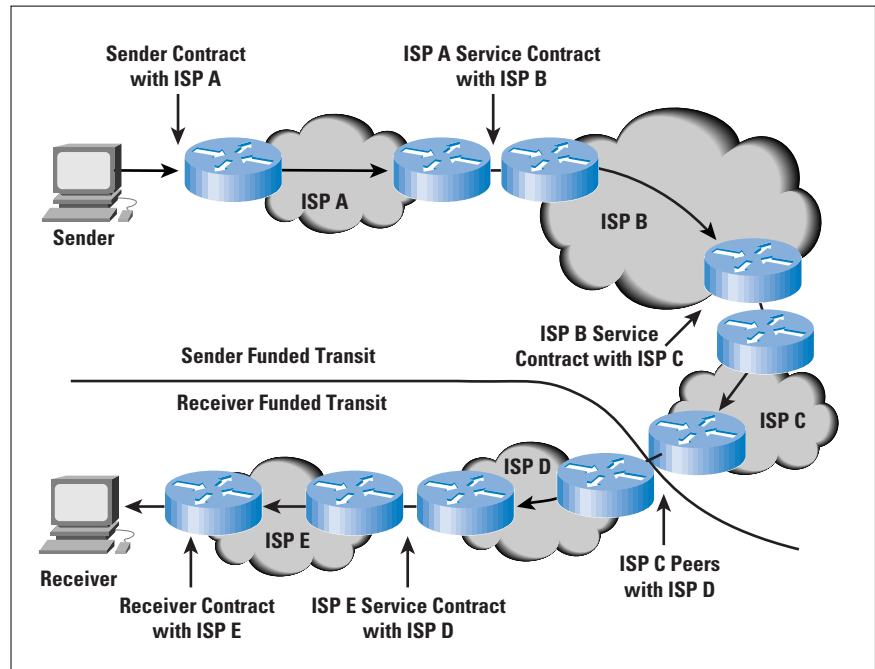
The issue of Internet settlements, and associated financial models of settlement, has occupied the attention of a large number of ISPs, traditional communications carriers, public regulators, and many other interested bodies for many years now. Despite these concentrated levels of attention and analysis, the Internet interconnection environment remains one where there are no soundly based models of financial settlement in widespread use today.

It is useful to look further into this matter, and pose the question: “Why has the Internet managed to pose such a seemingly intractable challenge to the ISP industry?” The prime reason is likely to be found within the commonly adopted retail model of ISP services. The tariff for an ISP retail service does not implicitly cover the provision of an Internet transmission service from the client to all other Internet-connected hosts. In other words, the Internet service, as retailed to the client, is not a comprehensive end-to-end service.

In a simple model of the operation of the Internet, each ISP owns and operates some local network infrastructure, and may choose to purchase services from one or more upstream service providers. The service domain offered to the clients of this network specifically encompasses an Internet subdomain limited to the periphery of the ISP network together with the periphery of the contracted upstream provider’s service domain. This is a recursive domain definition, in that the upstream provider in turn may have purchased services from an upstream provider at the next tier, and so on. After the client’s traffic leaves this service domain, the ISP ceases to directly, or indirectly, fund the carriage of the client’s traffic, and the funding burden passes over to a funding chain linked to the receiver’s retail service.

For example, when traffic is passed from an ISP client to a client of another provider, the ISP funds the traffic as it transits through the ISP and indirectly funds the cost of carriage through any upstream provider’s network. When the traffic leaves the provider’s network, to be passed to either a different client, another ISP, or to a peer provider, the sender’s ISP ceases to fund the further carriage of the traffic. This scenario is indicated in Figure 6. In other words, these scenarios illustrate the common theme that the retail base of the Internet is not an end-to-end tariff base. The sender of the traffic does not fund the first hop ISP for the total costs of carriage through the Internet to the traffic’s destination, nor does the ultimate receiver pay the last hop ISP for these costs. The ISP retail pricing structure reflects an implicit division of cost between the two parties, and there is no consequent structural requirement for inter-provider financial balancing between the originating ISP and the terminating ISP.

Figure 6: Partial-Path
Paired Services



An initial reaction to this partial service model would be to wonder why the Internet works at all, given that no single party funds the carriage of traffic on the complete path from sender to receiver. Surely this would imply that once the traffic had passed beyond the sending ISP's service funded domain the traffic should be discarded as unfunded traffic? The reason why this is not the case is that the receiver implicitly assumes funding responsibility for the traffic at this handover point, and the second part of the complete carriage path is funded by the receiver. In an abstract sense, the entire set of connectivity paths within the Internet can be viewed as a collection of bilaterally funded path pairs, where the sender funds the initial path component and the receiver funds the second terminating path component. This underscores the original observation that the generally adopted retail model of Internet services is not one of end-to-end service delivery, but instead one of partial path service, with no residual retail price component covering any form of complete path service.

Financial settlement models typically are derived from a different set of initial premises than those described here. The typical starting point is that the retail offering is a comprehensive end-to-end service, and that the originating service provider utilizes the services of other providers to complete the delivery of all components of the retailed service. The originating service provider then undertakes some form of financial settlement with those providers who have undertaken some form of an operational role in providing these service elements. This cost-distributed business structure allows both small and large providers to operate with some degree of financial stability, which in turn allows a competitive open service market to thrive. Through the operation of open competition, the consumer gains the ultimate price and service benefit of cost-efficient retail services.

The characteristics of the Internet environment tend to create a different business environment to that of a balanced cost distribution structure. Here there is a clear delineation between a customer/provider relationship and a peer relationship, with no stable middle ground of a financially settled inter-ISP bilateral relationship. An ISP customer is one that assumes the role of a customer of one or a number of upstream providers, with an associated flow of funding from the customer to the upstream provider, whereas an ISP upstream service provider views the downstream provider as a customer. An ISP peer relationship is where the two ISPs execute a peering arrangement, where traffic is exchanged between the two providers without any consequent financial settlement, and such peering interactions are only stable while both providers perceive some degree of parity in the arrangement; for example, when the two providers present to the peering point Internet domains of approximate equality in market coverage and market share. An ISP may have multiple simultaneous relationships, being a customer in some cases, an upstream provider in others, and a peer in others. In general, the relationships are unique within an ISP pairing, and efforts to support a paired relationship which encompasses elements of both peering and customer/provider pose significant technical and business challenges.

The most natural business outcome of any business environment is for each provider to attempt to optimize its business position. For an ISP, this optimization is not simply a case of a competitive impetus to achieve cost efficiency in the ISP's internal service operation, because the realization of cost efficiencies within the service provider's network does not result in any substantial change in the provider's financial position with respect to upstream costs or peering positioning. The ISP's path toward business optimization includes a strong component of increasing the size and scope of the service provider operation, so that the benefits of providing funded upstream services to customers can be maximized, and non-financially settled peering can be negotiated with other larger providers.

The conclusion drawn is that the most natural business outcome of today's Internet settlement environment is one of aggregation of providers, a factor quite evident in the Internet provider environment at present.

Quality of Service and Financial Settlements

Within today's ISP service model, strong pressure to change the technology base to accommodate more sophisticated settlement structures is not evident. The fundamental observation is that any financial settlement structure is robust only where a retail model exists that is relatively uniform in both its nature and deployment, and encompasses the provision of services on an end-to-end basis. Where a broad diversity of partial-service retail mechanisms exists within a multiprovider environment, the stability of any form of interprovider financial settlement structure will always be dubious at best.

If paired partial path service models and SKA peering interconnection comfortably match the requirements of the ISP industry today, is this entire financial settlement issue one of simple academic interest?

Perhaps the strongest factor driving change here is the shift towards an end-to-end service model associated with the current technology impetus toward support of distinguished *Quality of Service* (QoS) mechanisms. Where a client signals the requirement for some level of preemption or reservation of resources to support an Internet transaction or flow, the signal must be implemented on an end-to-end basis in order for the service request to have any meaning or value. The public Internet business model to support practical use of such QoS technologies will shift to that of the QoS signal initiator undertaking to bear the cost of the entire end-to-end traffic flow associated with the QoS signal. This is a retail model where the application initiator undertakes to fund the entire cost of data transit associated with the application. This model is analogous to the end-to-end retail models of the telephony, postal, and freight industries. In such a model, the participating agents are compensated for the use of their services through a financial distribution of the original end-to-end revenue, and a logical base for inter-agent financial settlements is the outcome. It is, therefore, the case that meaningful inter-provider financial settlements within the Internet industry are highly dependent on the introduction of end-to-end service retail models. These financial settlements are, in turn, dependent on a shift from universal deployment of a best effort service regime with partial path funding to the introduction of layered end-to-end service regimes that feature both end-to-end service-level undertakings and end-to-end tariffs applied to the initiating party.

The number of conditionals in this argument is not insignificant. If QoS technologies are developed that scale to the size of the public Internet, that provide sufficiently robust service models to allow the imposition of service level agreements with service clients, and are standardized such that the QoS service models are consistent across all vendor platforms, then this area of inter-provider settlements will need to change as a consequence. The pressure to change will be emerging market opportunities to introduce interprovider QoS interconnection mechanisms and the associated requirement to introduce end-to-end retail QoS services. The consequence is that there will be pressure to support this with inter-provider financial settlements where the originating provider will apportion the revenue gathered from the QoS signal initiator with all other providers that are along the associated end-to-end QoS flow path.

Such an end-to-end QoS settlement model assumes significant proportions that may in themselves impact on the QoS signaling technologies. It is conceivable that each provider along a potential QoS path may need to signal not only their capability of supporting the QoS profile of the potential flow, but also the unit settlement cost that will apply to the flow. The end user may then use this cost feedback to determine

whether to proceed with the flow given the indication of total transit costs, or request alternate viable paths in order to choose between alternative provider paths so as to optimize both the cost and the resultant QoS service profile. The technology and business challenges posed by such an end-to-end QoS deployment model are certainly an impressive quantum change from today's best effort Internet.

With this in mind, one potential future is that the public Internet environment will adopt a QoS mediated service model that is capable of supporting a diverse competitive industry through interprovider financial settlements. The alternative is the current uniform best effort environment with no logical role for interprovider settlements, with the associated strong pressures for provider aggregation. The reliance on Internet QoS technologies to achieve not only Internet service outcomes, but also to achieve desired public policy outcomes in terms of competitive pressures, is evident within this perspective. It is unclear whether the current state of emerging QoS technologies and QoS interconnection agreements will be able to mature and be deployed in time to forge a new chapter in the story of the Internet interconnection environment. The prognosis for this is, however, not good.

Futures

Without the adoption of a settlement regime that supports some form of cost distribution among Internet providers, there are serious structural problems in supporting a diverse and well populated provider industry sector. These problems are exacerbated by the additional observation that the Internet transmission and retail markets both admit significant economies of scale of operation. The combination of these two factors leads to the economic conclusion that the Internet market is not a sustainable open competitive market. Under such circumstances, there is no natural market outcome other than aggregation of providers, leading to the establishment of monopoly positions in the Internet provider space. This aggregation is already well underway, and direction of the Internet market will be forged through the tension between this aggregation pressure and various national and international public policy objectives that relate to the Internet industry.

The problem stated here is not in the installation of transmission infrastructure, nor is it in the retailing of Internet services. The problem faced by the Internet industry is in ensuring that each provider of infrastructure is fairly paid when the infrastructure is used. In essence, the problem is how to distribute the revenue gained from the retail sale of Internet access and services to the providers of carriage infrastructure. While explosive growth has effectively masked these problems for the past decade, after market saturation occurs and growth tapers off, these issues of financial settlement between the various Internet industry players will then shape the future of the entire global ISP industry.

[This article is based in part on material in *The ISP Survival Guide*, by Geoff Huston, ISBN 0-471-31499-4, published by John Wiley & Sons in 1998. Used with permission.]

Annotated Reading List

The following articles and publications address various aspects of Internet interconnection and peering, and the underlying issues of the economics of Internet carriage.

- [0] Huston, G., "Interconnection, Peering and Settlements—Part I," *The Internet Protocol Journal*, Volume 2, Number 1, March 1999.
The first part of this article.

- [1] Huston, G., *ISP Survival Guide*, ISBN 0-471-31499-4, John Wiley & Sons, November 1998.
A more comprehensive view of the technology, business and strategy behind the Internet service sector.

- [2] Halabi, B., *Internet Routing Architectures*, ISBN 1-56205-652-2, Cisco Press, April 1997.
An excellent information resource on how to configure BGP to express policies for interconnecting networks.

- [3] Frieden, R., "Without Public Peer: The potential Regulatory and Universal Service Consequences of Internet Balkanization," *Virginia Journal of Law and Technology*, ISSN 1522-1687, Vol. 3, Sept. 1998.
http://vjolt.student.virginia.edu/graphics/vol3/vol3_art8.html.
A good briefing paper from an economic perspective on interconnection issues, with particular attention to the domestic situation in the United States.

- [4] Cukier, K., "Peering and Fearing: ISP Interconnection and Regulatory Issues," Presented paper at the Harvard Information Infrastructure Project Conference on the Impact of the Internet on Communication Policy, December 3–5 1997.
Conference program is at:
<http://ksgwww.harvard.edu/iip/iicompol/agenda.html>
The Cukier paper is at:
<http://ksgwww.harvard.edu/iip/iicompol/Papers/Cukier.html>

- [5] Shapiro, C., Varian, H., *Information Rules: A Strategic Guide to the Information Economy*, ISBN 087584863X, Harvard Business School Press, November 1998.
A broader look at the Internet from an economic perspective, looking at both content and service provider economics.

- [6] Varian, H., "The Information Economy—The Economics of the Internet, Information Goods, Intellectual Property and Related Issues," <http://www.sims.berkeley.edu/resources/infoecon/>
This is a collection of references to other online resources, and is a useful starting point for further reading on this topic.

GEOFF HUSTON holds a B.Sc and a M.Sc from the Australian National University. He has been closely involved with the development of the Internet for the past decade. He was responsible for the initial build of the Internet within the Australian academic and research sector. Huston is currently the Chief Technologist in the Internet area for Telstra. He is also an active member of the IETF, and is a member of the Internet Society Board of Trustees. He is author of *The ISP Survival Guide*, and coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, a collaboration with Paul Ferguson. Both books are published by John Wiley & Sons.
E-mail: gih@telstra.net

Firewalls and Internet Security, the Second Hundred (Internet) Years

by Frederic Avolio,
Avolio Consulting

Interest and knowledge about computer and network security is growing along with the need for it. This interest is, no doubt, due to the continued expansion of the Internet and the increase in the number of businesses that are migrating their sales and information channels to the Internet. The growth in the use of networked computers in business, especially for e-mail, has also fueled this interest. Many people are also presented with the post-mortems of security breaches in high-profile companies in the nightly news and are given the impression that some bastion of defense had failed to prevent some intrusion. One result of these influences is that many people feel that Internet security and Internet firewalls are synonymous. Although we should know that no single mechanism or method will provide for the entire computer and network security needs of an enterprise, many still put all their network security eggs in one firewall basket.

Computer networks may be vulnerable to many threats along many avenues of attack, including:

- *Social engineering*, wherein someone tries to gain access through social means (pretending to be a legitimate system user or administrator, tricking people into revealing secrets, etc.)
- *War dialing*, wherein someone uses computer software and a modem to search for desktop computers equipped with modems that answer, providing a potential path into a corporate network
- *Denial-of-service attacks*, including all types of attacks intended to overwhelm a computer or a network in such a way that legitimate users of the computer or network cannot use it
- *Protocol-based attacks*, which take advantage of known (or unknown) weaknesses in network services
- *Host attacks*, which attack vulnerabilities in particular computer operating systems or in how the system is set up and administered
- *Password guessing*
- *Eavesdropping* of all sorts, including stealing e-mail messages, files, passwords, and other information over a network connection by listening in on the connection.

Internet firewalls have been around for a hundred years—in Internet time. Firewalls can help protect against some of these attacks, but certainly not all. Firewalls can be very effective at what they do. The people who set up and use them must have the knowledge of how they work, and also be aware of what they can and cannot protect. In this article, we examine the Internet firewall, touch on its history, see how firewalls are used today, and discuss changes that are in place for the next hundred years.

Internet History

In the beginning, there was no Internet. There were no networks. There was no e-mail, and people relied on postal mail or the telephone to communicate. The very busy sent telegrams. Few people used ugly names to refer to others whom they had never met. Of course, the Internet has changed all this. The Internet, which started as the *Advanced Research Projects Agency Network* (ARPANET), was a small, almost closed, community. It was a place, to borrow a line from the theme to *Cheers*, “where everybody knows your name, and they’re always glad you came.”

On November 2, 1988, something happened that changed the Internet forever. Reporting this incident, Peter Yee at the NASA Ames Research Center sent a note out to the TCP/IP Internet mailing list that reported, “We are currently under attack from an Internet VIRUS! It has hit Berkeley, UC San Diego, Lawrence Livermore, Stanford, and NASA Ames.” Of course, this report was the first documentation of what was to be later called *The Morris Worm*. The researchers and contributors that had built the Internet, as well as the organizations that were starting to use it, realized at that moment that the Internet was no longer a closed community of trusted colleagues. In fact, it hadn’t been for years. To their credit, the Internet community did not overreact to this situation. Rather, they started sharing information on their practices to prevent future disruptions.

(One of the results of this problem was a growth in the number of Internet mailing lists dedicated to security and bug tracking. The *firewalls* list—subscribe with e-mail to **Majordomo@lists.gnac.net**—and the *bugtraqs* list—**LISTSERV@netspace.org**—are two examples, as well as the *CERT Coordination Center*—<http://www.cert.org/>.)

Other famous, and general, attacks followed:

- Bill Cheswick’s “evening with Berferd”^[4]
- Clifford Stoll’s run-in with German spies^[7]
- The massive password capture of the winter of 1994
- The IP spoofing attack that Kevin Mitnick used against Tsutomu Shimomura^[6]
- The rash of denial-of-service attacks in January 1996, and the “Web site break-in of the week.”

All these viruses have made it into the popular press, and all have raised awareness of the need for good computer and network security. As these, and other, events were unfolding, the firewall was starting its rapid evolution. Although the development of firewall technology and products may be seen as very fast, it sometimes seems that firewalls are just barely keeping up with the new applications and services that spring up and immediately become a “requirement” for many Internet users.

Firewall History

We are used to firewalls in other disciplines, and, in fact, the term did not originate with the Internet. We have firewalls in housing, separating, for example, a garage from a house, or one apartment from another. Firewalls are barriers to fire, meant to slow down its spread until the fire department can put it out. The same is true for firewalls in automobiles, segregating the passenger and engine compartments.

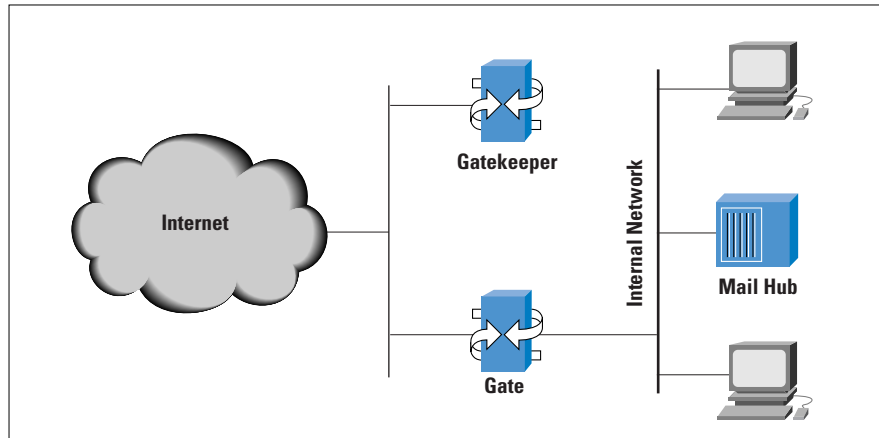
Cheswick and Bellovin, in the definitive text on Internet firewalls^[4], said an Internet firewall has the following properties: it is a single point between two or more networks where all traffic must pass (choke point); traffic can be controlled by and may be authenticated through the device, and all traffic is logged. In a talk, Bellovin later stated, “Firewalls are barriers between ‘us’ and ‘them’ for arbitrary values of ‘them.’”

The first network firewalls appeared in the late 1980s and were routers used to separate a network into smaller LANs. In these scenarios—and using Bellovin’s definition, above—“us” might be—well, “us.” And “them” might be the English Department. Firewalls like this were put in place to limit problems from one LAN spilling over and affecting the whole network. All this was done so that the English Department could add any applications to its own network, and manage its network in any way that the department wanted. The department was put behind a router so that problems due to errors in network management, or noisy applications, did not spill over to trouble the whole campus network. The first security firewalls were used in the early 1990s. They were IP routers with filtering rules. The first security policy was something like the following: allow anyone “in here” to access “out there.” Also, keep anyone (or anything I don’t like) “out there” from getting “in here.” These firewalls were effective, but limited. It was often very difficult to get the filtering rules right, for example. In some cases, it was difficult to identify all the parts of an application that needed to be restricted. In other cases, people would move around and the rules would have to be changed.

The next security firewalls were more elaborate and more tunable. There were firewalls built on so-called *bastion hosts*. Probably the first commercial firewall of this type, using filters and application gateways (proxies), was from Digital Equipment Corporation, and was based on the DEC corporate firewall. Brian Reid and the engineering team at DEC’s Network Systems Lab in Palo Alto originally invented the DEC firewall. The first commercial firewall was configured for and delivered to the first customer, a large East Coast-based chemical company, on June 13, 1991. During the next few months, Marcus Ranum at Digital invented security proxies and rewrote much of the rest of the firewall code. The firewall product was produced and dubbed DEC SEAL (for *Secure External Access Link*). The DEC SEAL was made up of an external system, called *Gatekeeper*, the only system the Internet could talk to, a filtering gateway, called *Gate*, and an internal *Mailhub* (see Figure 1).

In this same time frame, Cheswick and Bellovin at Bell Labs were experimenting with circuit relay-based firewalls. Raptor Eagle came out about six months after DEC SEAL was first delivered, followed by the ANS InterLock.

Figure 1: DEC SEAL—
First Commercial
Firewall



On October 1, 1993, the Trusted Information Systems (TIS) *Firewall Toolkit* (FWTK) was released in source code form to the Internet community. It provided the basis for TIS' commercial firewall product, later named *Gauntlet*. At this writing, the FWTK is still in use by experimenters, as well as government and industry, as a basis for their Internet security. In 1994, Check Point followed with the *Firewall-1* product, introducing “user friendliness” to the world of Internet security. The firewalls before Firewall-1 required editing of ASCII files with ASCII editors. Check Point introduced icons, colors, and a mouse-driven, X11-based configuration and management interface, greatly simplifying firewall installation and administration.

Early firewall requirements were easy to support because they were limited to the Internet services available at that time. The typical organization or business connecting to the Internet needed secure access to remote terminal services (Telnet), file transfer (*File Transfer Protocol* [FTP]), electronic mail (*Simple Mail Transfer Protocol* [SMTP]), and USENET News (the *Network News Transfer Protocol*—NNTP). Today, we add to this list of “requirements” access to the World Wide Web, live news broadcasts, weather information, stock quotes, music on demand, audio and videoconferencing, telephony, database access, file sharing, and the list goes on.

What new vulnerabilities are there in these new “required” services that are daily added to some sites? What are the risks? Too often, the answer is “we don’t know.”

Types of Firewalls

There are four types of Internet firewalls, or, to be more accurate, three types plus a hybrid. The details of these different types are not discussed here because they are very well covered in the literature.^[1, 3, 4, 5]

Packet Filtering

One kind of firewall is a packet filtering firewall. Filtering firewalls screen packets based on addresses and packet options. They operate at the IP packet level and make security decisions (really, “to forward, or not to forward this packet, that is the question”) based on the headers of the packets.

The filtering firewall has three subtypes:

- *Static Filtering*, the kind of filtering most routers implement—filter rules that must be manually changed
- *Dynamic Filtering*, in which an outside process changes the filtering rules dynamically, based on router-observed events (for example, one might allow FTP packets in from the outside, if someone on the inside requested an FTP session)
- *Stateful Inspection*, a technology that is similar to dynamic filtering, with the addition of more granular examination of data contained in the IP packet

Dynamic and stateful filtering firewalls keep a dynamic state table to make changes to the filtering rules based on events.

Circuit Gateways

Circuit gateways operate at the network transport layer. Again, connections are authorized based on addresses. Like filtering gateways, they (usually) cannot look at data traffic flowing between one network and another, but they do prevent direct connections between one network and another.

Application Gateways

Application gateways or proxy-based firewalls operate at the application level and can examine information at the application data level. (We can think of this as the *contents* of the packets, though strictly speaking proxies do not operate with packets.) They can make their decisions based on application data, such as commands passed to FTP, or a URL passed to HTTP. It has been said that application gateways “break the client/server model.”

Hybrid firewalls, as the name implies, use elements of more than one type of firewall. Hybrid firewalls are not new. The first commercial firewall, DEC SEAL, was a hybrid, using proxies on a bastion host (a fortified machine, labeled “Gatekeeper” in Figure 1), and packet filtering on the gateway machine (“Gate”). Hybrid systems are often created to quickly add new services to an existing firewall. One might add a circuit gateway or packet filtering to an application gateway firewall, because it requires new proxy code to be written for each new service provided. Or one might add strong user authentication to a stateful packet filter by adding proxies for the service or services.

No matter what the base technology, a firewall still basically acts as a controlled gateway between two or more networks through which all traffic must pass. A firewall enforces a security policy and it keeps an audit trail.

What a Firewall Can Do

A firewall intercepts and controls traffic between networks with differing levels of trust. It is part of the network perimeter defense of an organization and should enforce a network security policy. By Cheswick's and Bellovin's definition, it provides an audit trail. A firewall is a good place to support strong user authentication as well as private or confidential communications between firewalls. As pointed out by Chapman and Zwicky^[2], firewalls are an excellent place to focus security decisions and to enforce a network security policy. They are able to efficiently log internetwork activity, and limit the exposure of an organization.

The exposure to attack is called the "zone of risk." If an organization is connected to the Internet without a firewall (Figure 2), every host on the private network can directly access any resource on the Internet. Or to put it as a security officer might, every host on the Internet can attack every host on the private network. Reducing the zone of risk is better. An internetwork firewall allows us to limit the zone of risk. As we see in Figure 3, the zone of risk becomes the firewall system itself. Now every host on the Internet can attack the firewall. With this situation, we take Mark Twain's advice to "Put all your eggs in one basket—and watch that basket."

Figure 2: Zone of Risk
for an Unprotected
Private Network

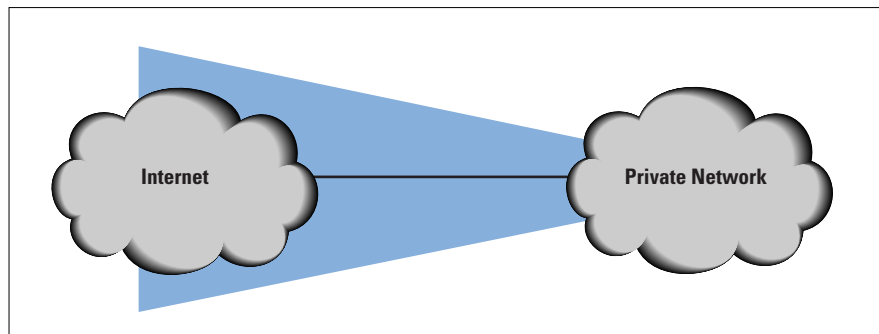
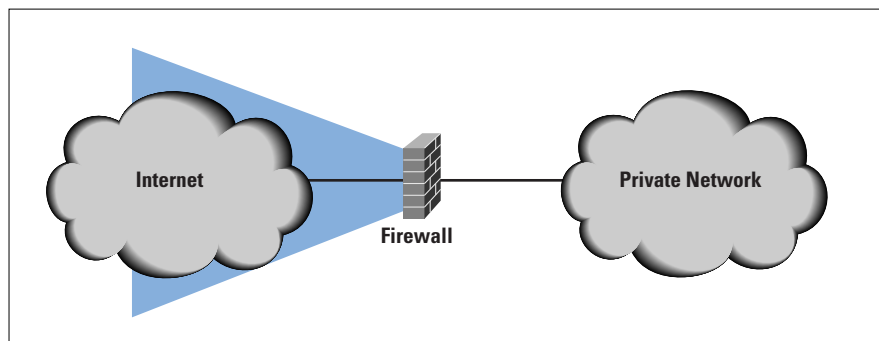


Figure 3: Zone of Risk
with a Firewall



What a Firewall Cannot Do

Firewalls are terrible at reading people's minds or detecting packets of data with "bad intent." They often cannot protect against an insider attack (though might log network activity, if an insider uses the Internet gateway in his crime). Firewalls also cannot protect connections that do not go through the firewall. In other words, if someone connects to the Internet through a desktop modem and telephone, all bets are off. Firewalls provide little protection from previously unknown attacks, and typically provide poor protection against computer viruses.

Firewalls Today: Additions

The first add-on to Internet firewalls was strong user authentication. If your security policy allows access to the private network from an outside network, such as the Internet, some kind of user authentication mechanism is required. User authentication simply means "to establish the validity of a claimed identity." A username and password provides user authentication, but not *strong* user authentication. On a nonprivate connection, such as an unencrypted connection over the Internet, a username and password can be copied and replayed. Strong user authentication uses cryptographic means, such as certificates, or uniquely keyed cryptographic calculators. These certificates prevent "replay attacks"—where, for example, a username and password are captured and "replayed" to gain access. Because of where it sits—on both the "trusted" and "untrusted" networks—and because of its function as a controlled gateway, a firewall is a logical place to put this service.

The next add-on to Internet firewalls was firewall-to-firewall encryption, first introduced on the ANS InterLock Firewall. Today, such an encrypted connection is known as a Virtual Private Network, or VPN. It is "private" through the use of cryptography. It is "virtually" private because the private communication flows over a public network—the Internet, for example. Although VPNs were available before firewalls via encrypting modems and routers, they came into common use running on firewalls. Today, most people expect a firewall vendor to offer a VPN option. Firewalls act as the endpoint for VPNs between the enterprise and mobile users or telecommuters, keeping communication confidential from notebook PC, home desktop, or remote office.

In the past two years, it has become popular for firewalls to also act as content screening devices. Some additions to firewalls in this area include virus scanning, URL screening, and key word scanners (also known in U.S. government circles as "guards"). If the security policy of your organization mandates screening for computer viruses—and it should—it makes sense to put such screening at a controlled entry point for computer files, such as the firewall. In fact, standards exist for plugging antivirus software into the data flow of the firewall, to intercept and analyze data files. Likewise, URL screening—firewall controlled access to the World Wide Web—and content screening of files and messages seem like logical additions to a firewall. After all, the data is

flowing through the fingers of the firewall system, so why not examine it and allow the firewall to enforce the security policies of the organization? The downside to this scenario is performance. Also virus scanning must ultimately be performed on each desktop because data may come in to the desktops from paths other than through the firewall—for instance, the floppy.

Recently, some firewall and router vendors have been making the case for a relatively new firewall add-on called “flow control” to deliver Quality of Service (QoS). QoS, for example, can limit the amount of network bandwidth any one user can take up, or limit how much of the network capacity can be used for specific services (such as FTP or the Web). Once again, because the firewall is the gateway, it is the logical place to put a QoS arbitrating mechanism.

Firewalls Tomorrow

In 1997, The Meta Group, and others, predicted that firewalls would be the center of network and internetwork security^[7]. After all, firewalls were the first big security item, the first successful Internet security product, and the most visible security device. They quickly became a “must have”—this is good—and a “good enough”—this is not good because firewalls alone are not sufficient. Firewalls became synonymous with security, as mentioned above. The firewall console becoming the network security console seemed natural at that time. But this scenario has not happened, nor will it happen. The reason? The firewall is just another mechanism used to enforce a security policy. This specific enforcement device will not be the policy management device.

As organizations broaden the base of measures and countermeasures used to implement a comprehensive network and computer security policy, firewalls will need to communicate with and interact with other devices. Intrusion detection devices—running on or separate from the firewall—must be able to reconfigure the firewall to meet a new perceived threat (just as dynamic filtering firewalls today “reconfigure” themselves to meet the needs of a user).

Firewalls will have to be able to communicate with network security control systems, reporting conditions and events, allowing the control system to reconfigure sensors and response systems. A firewall could signal an intrusion detection system to adjust its sensitivity, as the firewall is about to allow an authenticated connection from outside the security perimeter. A central monitoring station could watch all this, make changes, react to alarms and other notifications, and make sure that all antivirus software and other content screening devices were functioning and “up to rev.” Some products have started down this path already. The *Intrusion Detection System* (IDS) and firewall reconfiguration of network routers based on perceived threat is a reality today. Also, firewall-resident IDS and help-desk software enable another vendor’s system to expand from a prevention mechanism into detecting and re-

sponding. The evolution continues and firewalls are changing rapidly to address the next 100 (Internet) years.

In June 1994, the author wrote^[5], “Firewalls are a stopgap measure—needed because many services are developed that operate either with poor security or no security at all.” This statement is erroneous. Firewalls are *not* a stopgap measure. Firewalls play an important part in a multilevel, multilayer security strategy. Internet security firewalls will not go away, because the problem firewalls address—access control and arbitration of connections in light of a network security policy—will not go away.

As use of the Internet and internetworked computers continues to grow, the use of Internet firewalls will grow. They will no longer be the only security mechanism, but will cooperate with others on the network. Firewalls will morph—as they have—from what we recognize today, just as walls of brick and mortar were eventually replaced by barbed wire, motion sensors, and video cameras—and brick and mortar. But Internet firewalls will continue to be a required part of the methods and mechanisms used to enforce a corporate security policy.

References

- [1] Avolio, F. and Ranum, M., “A Network Perimeter with Secure External Access,” Proceedings of the ISOC NDSS Symposium, 1996.
(<http://www.avolio.com/netsec.html>)
- [2] Chapman, D. B. and Zwicky, E., *Building Internet Firewalls*, ISBN 1-56592-124-0, O'Reilly and Associates, 1995.
- [3] Cheswick, W. and Bellovin, S., *Firewalls and Internet Security: Repelling the Wily Hacker*, ISBN 0201633574, Addison-Wesley, 1994.
- [4] Ranum, M. and Avolio, F., “A Toolkit and Methods for Internet Firewalls,” Proceedings of the summer USENIX conference, 1994.
(<http://www.avolio.com/fwtk.html>)
- [5] Shimomura, T. and Markoff, J., *Takedown: The Pursuit and Capture of Kevin Mitnick, America's Most Wanted Computer Outlaw—By the Man Who Did It*, ISBN 0-7868-89136, Warner Books, 1996.
- [6] Stoll, C., *The Cuckoo's Egg: Tracking a Spy through the Maze of Computer Espionage*, ISBN 0671726889, Reprint edition, Pocket Books, 1995.
- [7] Meta Global Networking Strategies File 549, November 24, 1997.

FREDERICK M. AVOLIO is an independent security consultant. He has lectured and consulted on Internet gateways and firewalls, security, cryptography, and electronic mail configuration for both government and industry, working in the UNIX and TCP/IP communities since 1979. He is a top-rated speaker and contributor to NetWorld+Interop, USENIX, SANS, TISC, and other security-related forums. With Paul Vixie, Avolio wrote the book *Sendmail: Theory and Practice*, published by Digital Press. He has an undergraduate degree in Computer Science from the University of Dayton and a Master of Science from Indiana University. E-mail: fred@avolio.com

Was the Melissa Virus So Different?

by Barbara Y. Fraser, Lawrence R. Rogers, and Linda H. Pesante,
Software Engineering Institute, Carnegie Mellon University

Was the recent electronic mail-based *Melissa* virus so different from similar events in our noncyberspace lives that it merits special behavior? We don't think so. But recent events raise some interesting questions about where to draw the line in our concern about the safety of our mailbox contents.

We regularly receive samples in the mail and don't give them much thought. They run the gamut from laundry detergents to shampoos to cereals to pain relievers. How often do we rip open that sample box of sugar-coated cereal and chomp down a few handfuls as a snack? Do we question whether the labeling accurately reflects the contents of the package? And what about the shampoo samples in those convenient little bottles, just the right size for tossing into our travel bag for the next trip. We use the shampoo with no thought that it might really be hair dye that would turn our hair purple or green. Then there are the sample medications and herbal remedies. Do we use the sample, assuming that it is exactly what it seems to be, without verifying it in some way?

For many of us, these examples represent common behavior today. When we open the samples we find in our mailbox, we don't question whether someone intent on harming us has sent a product that appears to be something we would use and that seems to come from a trusted source. Rarely, if ever, would we call manufacturers and ask whether they had really sent the sample.

How different is this from our approach to the contents of our electronic mailbox? We urge people *never* to click on an attachment before verifying its contents—or at least not until they've verified that it came from the stated sender. Surely we must make these recommendations because of malicious code in electronic mail messages. But we may be asking people to behave differently in cyberspace than they typically do in their noncyberspace life.

What are we to do then? Responsible cyberspace behavior says to trust nothing and verify everything as completely as possible. This scenario would mean that attachments added to an electronic mail messages must be analyzed before being used. To be the most effective, analyzers must be kept up-to-date with the latest information. Even then, rapidly spreading viruses like *Melissa* can slip under our "radar" for a while. Tools that support authentication and integrity are another building block we should use to gain trust in information that we should otherwise consider untrustworthy.

In our noncomputer lives, how do we know that the medication sample that came in the mail actually came from the attributed vendor? How do we know that the sample was not changed after it left the manufacturing point? The best we can do is to call the manufacturer and exchange some information about the sample: product numbers, packaging color, descriptions of the sample, and so on. Still, we cannot be completely sure that the product is what the packing says it is. Similarly, how do we know that the electronic mail attachment actually came from the stated sender or that it was not changed in transit?

Here cyberspace has the edge over noncyberspace. Technologies are available that help us to verify the mail sender (authentication) and the validity of the message (integrity). Alas, none of the available technologies are multivendor, interoperable, or approved or endorsed by the Internet's standardization body. These technologies are an improvement over their noncyberspace counterparts, but they are not yet mature enough or widespread enough to be as effective as they ultimately will become. Unfortunately, we need that maturity now.

Returning to our original question: Was the Melissa virus so different? Our answer is *no*, it was not so different from the comparable free samples we receive in our noncyberspace lives. Unfortunately, those lives are fraught with the same kind of problems, yet we accept those risks with little concern for our well-being. The real answer is that both our cyberspace and noncyberspace lives need to change to reflect the challenges of our modern world.

About Melissa

The CERT CC began receiving reports of a new virus on Friday, March 26, 1999. The macro virus is activated when a user opens an infected document in Microsoft Word 97 or Word 2000 with macros enabled. The virus is then quickly spread by sending an infected document to the first 50 addresses in the victim's Microsoft Outlook address book. It also infects the **Normal.dot** template file, a situation which in turn causes other Word documents created using this template to be infected with the virus. If these newly infected documents are opened by a second user, the document, including the virus, will propagate, sending the document to 50 addresses in the second user's address book. The CERT CC handled over 300 reported incidents involving Melissa, affecting over 100,000 computers. This estimate is very conservative because it counts only those who contacted the CERT CC. It is believed that millions of host computers were infected.

References

- [1] <http://www.cert.org/advisories/CA-99-04-Melissa-Macro-Virus.html>
- [2] <http://www.melissavirus.com/>
- [3] <http://www.nai.com/valert>
- [4] <http://www.datafellows.com/news/pr/eng/19990327.htm>
- [5] http://www.mcafee.com/about/press_releases/pr040299.asp
- [6] http://www.cert.org/other_sources/viruses.html

To subscribe to CERT Advisories:

http://www.cert.org/contact_cert/certmaillist.html

BARBARA FRASER is a senior member of the technical staff at the Software Engineering Institute (SEI) located at Carnegie Mellon University. She is currently working in the Networked Systems Survivability Program of the SEI and the CERT® Coordination Center. Barbara leads the team that is currently developing an adaptive security management model for networked systems that will allow organizations to adapt to technology and organization changes while maintaining an appropriate level of security in their networked systems. Her professional interests are in developing tools and techniques for improving the survivability of technologies currently deployed in the Internet. Barbara has been involved with the CERT Coordination Center since 1990. She has developed and delivered many talks and courses on Internet security and security incident response, and has worked with many organizations to help them understand and address security issues as they relate to the Internet. Barbara is currently coteaching a graduate course, "The Economics of Information Security," for the Heinz School of Public Policy at Carnegie Mellon University. Barbara is active in the security area of the Internet Engineering Task Force (IETF) and was one of the authors of RFC 1281, "Guidelines for the Secure Operation of the Internet," and RFC 2196, "Site Security Handbook." She is currently a member of the Security Area Directorate and chairs two IETF working groups (GRIP and SSH). Prior to joining the SEI, Barbara was a senior engineer at Martin Marietta Corporation (now Lockheed Martin), where she led a team of software engineers in the development of aircraft simulator software. Barbara holds a bachelor's degree in biology and an M.S. degree in computer science. E-mail: byf@cert.org

LAWRENCE R. ROGERS is a senior member of the technical staff in the Networked Systems Survivability Program at the Software Engineering Institute (SEI). The CERT Coordination Center is also a part of this program. Larry's primary focus in this group is analyzing system and network vulnerabilities and helping to transition security technology into production use. His professional interests are in the areas of the administering systems in a secure fashion and software tools and techniques for creating new systems being deployed in the Internet. Before joining the SEI, Larry worked for ten years at Princeton University, first in the Department of Computer Science on the Massive Memory Machine project, and later at the Department of Computing and Information Technology (CIT). While at CIT, Larry directed and managed the UNIX Systems Group that was charged with administering the UNIX computing facilities used for undergraduate education and campus-wide services. Larry coauthored the book *Advanced Programmer's Guide to UNIX Systems V* with Rebecca Thomas and Jean Yates. Larry received a B.S. degree in Systems Analysis from Miami University in 1976 and an M.A. degree in Computer Engineering in 1978 from Case Western Reserve University. E-mail: lrr@cert.org

LINDA HUTZ PESANTE has been a member of the technical staff of the Software Engineering Institute (SEI) since 1987. She is currently the leader of the Information Services Team for the CERT Coordination Center and SEI Networked Systems Survivability Program. She also teaches communication skills in the Master of Software Engineering Program at Carnegie Mellon University. At the University, she is a member of the Institutional Review Board for the Protection of Human Subjects in Research. She holds a B.A. in English and M.A. in professional writing from Carnegie Mellon, and an M. Ed. from the University of Pittsburgh. She has published on the topics of technical communication, network security, and teaching writing in computer science and software engineering programs. E-mail: 1hp@cert.org

Book Review

OPSF *OSPF: Anatomy of an Internet Routing Protocol*, John T. Moy, Addison Wesley Longman, ISBN 0-201-63472-4, 1998.
<http://www.awl.com/cseng/titles/0-201-63472-4>

Audience

John Moy takes the somewhat difficult topic of Internet routing and presents an understandable and engaging tour of specific parts of routing and how this one instance interrelates with other parts of Internet routing. This book is not for the routing novice, although the first couple of chapters provide a quick overview and history of routing and one viewpoint on the distinctions between two architectural choices in routing protocol design, *Distance Vector* and *Link State*. This book is really targeted for people that have a basic understanding of what routing is and would like to gain an understanding of this particular tool in the Internet routing “toolbox.”

Organization

The second section goes into great detail on one implementation of the Link State architecture, *Open Shortest Path First Protocol* (OSPF). There is a companion volume which contains OSPF specific details and includes source code for building an OSPF service on FreeBSD systems. He covers some background in the design phases of OSPF, delineating why certain choices were made in the evolution of OSPF as we know it today and then starts into what I think of as the heart of the book, an understandable, brief discussion of OSPF design with packet formats. In this section of the book, the author takes a textbook approach and closes each chapter with a series of exercises which test understanding of the principles covered in each chapter. At the end of the section, the FAQ answers a number of questions which operators that are considering OSPF will ask.

The book then changes focus and examines the basics of routing in the context of multicast aware infrastructure. This is an area that is still very dynamic and several of the presumptions that John makes in this section may not be as relevant in today’s networking environment. However, he does demonstrate the ability of OSPF to support new features, in this case the variant called *Multicast OSPF* or MOSPF. A discussion of the integration of MOSPF into OSPF networks as well as MOSPF in *Distance Vector Multicast Routing Protocol* (DVMRP) networks points out how different routing protocols can work together. DVMRP forms the central core of the Multicast Backbone or *Mbone*. Both DVMRP and MOSPF lack policy features that many operators demand and so this section remains more of academic interest in understanding how multicast can work.

The fourth section covers configuration and management of OSPF in real networks. Of specific interest to me is the discussion on how OSPF can take advantage of authentication features to ensure the integrity of the routing protocol and the data it sends. Others may find that a discussion of tools for troubleshooting more interesting. A fair amount of the discussion in this section deals with the use of *Simple Network Management Protocol* (SNMP) as the tool for managing and configuring OSPF. Its not clear to me that operators of parts of the Internet are comfortable with this approach since SNMP has known vulnerabilities. Such techniques are useful for monitoring OPSF activities and may be used in private networks with a higher comfort level.

Protocol Review

The book closes with a review of popular routing protocols, both current and historic for unicast and multicast environments. John covers some basic ideas on protocol interactions when systems run more than one but does not cover the interactions between multicast and unicast protocols.

—Bill Manning, USC-ISI
manning@isi.edu

Would You Like to Review a Book for IPJ?

We receive numerous books on computer networking from all the major publishers. If you've got a specific book you are interested in reviewing, please contact us and we will make sure a copy is mailed to you. The book is yours to keep if you send us a review. We accept reviews of new titles, as well as some of the "networking classics." Contact us at ipj@cisco.com for more information.

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

ICANN Update

As mentioned in previous issues of IPJ, the *Internet Corporation for Assigned Names and Numbers* (ICANN) began operation in early November 1998. Recently, ICANN announced that five companies have been selected to participate in the initial testbed phase of the new competitive *Shared Registry System*. These five participants will be the first to implement the new system for competition in the market for **.com**, **.net**, and **.org** domain name registration services. Currently, registration services for these domains are provided by Network Solutions, Inc. (NSI), which has enjoyed an exclusive right to handle registrations under a 1993 Cooperative Agreement with the U.S. Government. The five registrars participating in the testbed are, in alphabetical order: America Online, CORE (*Internet Council of Registrars*), France Telecom/Oléane, Melbourne IT, and register.com.

Under the Cooperative Agreement between NSI and the U.S. Government, the competitive registrar testbed program began on April 26 and will last until June 24, 1999 (Phase I). Following the conclusion of Phase I, the Shared Registry System for the **.com**, **.net**, and **.org** domains will be opened on equal terms to all accredited registrars, meaning that any company that meets ICANN's standards for accreditation will be able to enter the market as a registrar and offer customers competitive domain name registration services in these domains.

Meanwhile, ICANN continues to work on the formation of several *supporting organizations*, namely the *Domain Name Supporting Organization* (DNSO), the *Address Supporting Organization* (ASO), and the *Protocol Supporting Organization* (PSO). More information is available at: www.icann.org

IETF and Related links

The *Internet Engineering Task Force* (IETF) is responsible for the development of standards for Internet technology. Membership to the IETF is open and you can participate in person or subscribe to the IETF mailing list. The IETF meets three times per year. For a list of future meetings and other IETF information see: <http://www.ietf.org>

SIGCOMM

If you want to learn about the latest developments on the research side of networking you should check out SIGCOMM, the Association for Computing Machinery's Special Interest Group on Communications. You can find out more about the group and their annual conference at: <http://www.acm.org/sigcomm/sigcomm99>

Send us your comments!

We look forward to hearing your comments and suggestions regarding anything you read in this publication. Send us e-mail at: ipj@cisco.com

This publication is distributed on an "as-is" basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Engineering
MCI WorldCom, USA

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Edward R. Kozel, Sr. VP, Corporate Development
Cisco Systems, Inc., USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the Cisco News
Publications Group, Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 1999 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-10/5
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

Bulk Rate Mail U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol Journal

September 1999

Volume 2, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Web Caching	2
Gigabit Ethernet	21
One Byte at a Time	26
Letter to the Editor	29
Book Reviews	30
Call for Papers	36
Fragments	37

FROM THE EDITOR

More and more of the data traffic on the Internet is due to World Wide Web activity. Given the often-complex graphics contents of Web pages, this traffic represents a significant amount of data and leads to an overall requirement for more bandwidth across the system. But building “bigger pipes” is not the only way to achieve better performance. Generally speaking, Web pages are relatively static objects that reside in *one* location and are accessed repeatedly by *many* users, often from “far away.” If the contents of the most frequently accessed pages can be stored by a proxy residing more “local” with respect to the end user, significant reductions in download delay can be accomplished. Since the Internet comprises many expensive international circuits, such local mirroring of content is also highly desirable from the point of view of the Internet Service Providers. Storing information in a proxy server is called *caching*, and it is the subject of our first article. Geoff Huston explains the motivation behind—and the different approaches to—caching.

The most popular Local-Area Network (LAN) technology is *Ethernet*. Invented in 1973 by Bob Metcalfe as a 3-Mbps technology, Ethernet has evolved to the now-familiar 10Base-T and 100Base-T standards. Standardized in 1998, *Gigabit Ethernet* is the subject of our second article. Bill Stallings gives an overview of the Gigabit Ethernet standards and their application in enterprise networks. There is already discussion about 10-Gigabit Ethernet and even 100-Gigabit Ethernet. We will keep you posted on these developments.

Some readers have suggested that we publish a few short articles on limited topics. In this issue we bring you the first in what we hope will become a series of articles under the general heading “One Byte at a Time.” The article is by Tom Thomas and he discusses *active* and *passive* modes of the File Transfer Protocol (FTP). If you have suggestions for future topics in this series, please contact us at ipj@cisco.com

The so-called “Millennium Bug” or “Y2K Problem” has been well reported in all the media. Our *Fragments* section gives some specific information relating to Y2K and the Internet.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Web Caching

by Geoff Huston, Telstra

Web browsing dominates today's Internet. More than two-thirds of the traffic on the Internet today is generated by the Web. In looking at how to improve the quality of service delivered by the Internet, a very productive way to start is examining the performance of Web transactions. It is here that Web caching can play a valuable role in improving service quality for a large range of Internet users.

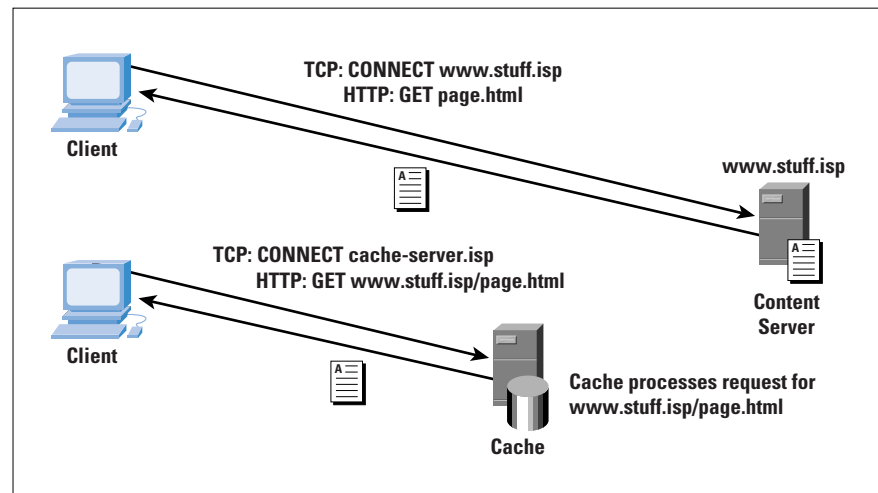
There are two types of Web caches—a *browser cache* and a *proxy cache*. A browser cache is part of all popular Web browsers. The browser keeps a local copy of all recently displayed pages, and when the user returns to one of these pages, the local copy is reused. By contrast, a proxy cache is a shared network device that can undertake Web transactions on behalf of a client, and, like the browser, the proxy cache stores the content. Subsequent requests for this content, by this or any other client of the cache, will trigger the cache to deliver the locally stored copy of the content, avoiding a repeat of the download from the original content source. In this article we look at proxy caches in further detail, particularly at the aspects of deployment of proxy caches in Internet Service Provider (ISP) networks.

What Is Proxy Web Caching?

When a browser wishes to retrieve a URL, it takes the host name component and translates that name to an IP address. A HTTP session is opened against that address, and the client requests the URL from the server.

When using a proxy cache, not much is altered in the transaction. The client opens a HTTP session with the proxy cache, and directs the URL request to the proxy cache instead (Figure 1).

Figure 1: A Proxy Web Transaction



If the cache contains the referenced URL it is checked for freshness by comparing with the “Expires:” date field of the content, if it exists, or by some locally defined freshness factor. Stale objects are revalidated with the server, and if the server revalidates the content, the object is remarked as fresh. Fresh objects are delivered to the client as a *cache hit*.

If the cache does not have a local copy of the URL, or the object is stale, this is a *cache miss*. In this case the cache acts as an agent for the client, opens its own session to the server named in the URL, and attempts a direct transfer to the cache.

The Pros and Cons of End-to-End Web Access

The original design principle of the Internet architecture is that of the end-to-end model^[2, 3]. Within this model the network is a passive instrument that undertakes a best effort to forward packets to the specified destination. Each packet generated by a host is assumed to be forwarded to the addressed destination, and any response to the datagram is assumed to come from that destination address.

The World Wide Web transaction protocol, the *Hypertext Transfer Protocol* (HTTP)^[4, 5], is constructed upon this model, where a client’s Web fetch causes a TCP session to be opened with the specified target host. The ensuing HTTP conversation identifies the requested data on the destination host, and this data is then passed back to the client. This delivery model is best expressed as a *just-in-time delivery model*, where the data is passed to the client on demand.

This delivery model has many significant advantages. The content server can modify the content, and all subsequent client requests are provided with the updated information, so that updates are immediately reflected in the delivered data. The content server is also able to track all content requests, allowing the content provider to track which particular content is being requested, the identity of each requestor, and how often each content item is referenced. The content provider can also differentiate between various clients, and, using some form of security model, the content provider can authenticate the client and deliver privileged information to certain clients. In this model the content provider can also differentiate between clients, delivering certain information to some clients, and *different* information to other clients of the content server.

Many web systems have been constructed based on the capability of this end-to-end delivery model. Continuously updating Web pages that use either *server push* or *client pull* to regularly update the content on the client’s display are used to display stock market prices, weather maps, or network management screens. Client identification can be used to create combined public and virtual private information servers, where a class of identified users can be directed to internal content environments, while other clients are passed to a default public content environment. Such systems form the basis of extranet environments, and can also be used to form part of a virtual private network.

Where information has a defined locality, this tool is very useful. Security and authentication is also used to provide services where the transaction requires some level of privacy. Electronic trading systems, credit card transactions, and related financial systems on the Web make use of such client authentication capabilities. The individual transaction can be encrypted using socket-level encryption,^[13] or the entire TCP session can be encrypted using an IP session-level encryption tool such as IP Security (IPSec).

For all these benefits available in an end-to-end model of Web content delivery, there are some balancing drawbacks. A server providing very popular content is placed under considerable stress, both in the number of simultaneous client connections active at any time and in the total volume of data being delivered from the server in the surrounding network. This load is expressed both as a server system load, and as load on the surrounding network. Improving the performance of such systems may entail improving the server throughput, increasing the number of servers through the use of server farms and a traffic manager, and improving the capacity of the local network to deliver the increased volume. However, all these measures may not address all the problems in maintaining quality of the content delivery. Modem-based client systems, and low-bandwidth wireless-based client systems are constrained by a combination of the restricted bandwidth of this last hop and the associated imposed end-to-end delay in conversing with the server. Improving the capacity of the server may not necessarily reduce the number of simultaneously active client connections. Reducing the delay between the client and the point of delivery of the content will improve the performance of content delivery.

In addition, the network itself may not be efficiently utilized. Web traffic does have considerable levels of duplication, where a set of clients request copies of the same content, and the network carries duplicates of the data to each client. For a network provider, where transmission capacity is a business cost, importing the content just once, and then passing local copies of this content to each client, is one method of improving the carriage efficiency of the network.

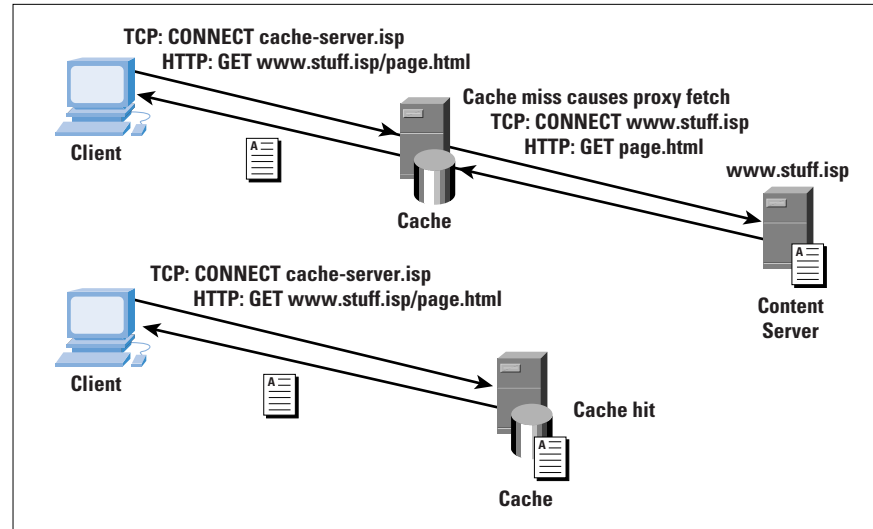
In terms of the ability to improve the service performance of delivery of content to a global network of clients, and in terms of the ability to improve the carriage efficiency of the network, caching of content makes some sense to the content provider, to the ISP, and to the end client.

The Pros and Cons of Web Proxy Caching

The same benefits of improved performance and reduced outbound traffic loads can be realized for World Wide Web traffic through the deployment of Web caches. Web caches are basically no different from any other form of caching. The client request is passed through a *cache agent*, which makes the request to the original source as a proxy for the client. The response of the server is retained in a local cache, and a copy

is passed to the client. If the same request is passed to the cache agent soon after the original request was serviced, the response can be generated from the cache without further reference to the original source. The operation of a Web cache is shown in Figure 2.

Figure 2: A Web Cache



Measurements of ISP traffic profiles indicate that some 70 percent of a typical ISP's traffic is Web-based traffic. An analysis of Web requests indicates that the typical level of similarity of requests (for the same object as one previously requested) can be as high as 50 percent of all Web-based traffic.

There are two hit-rate measures, a *page hit rate* and a *byte hit rate*. A page hit rate measures the proportion of individual HTTP requests that can be served from the cache, irrespective of the size of the page. A byte hit rate measures the ratio of the number of bytes delivered from the cache in hits against the number of bytes in misses. Experience to date has indicated that page hit rates of somewhere between 40 to 55 percent are achievable for a well-configured cache. In such circumstances the associated byte hit rate is between 20 and 35 percent. The major contributor to the hit rate is in image files.

For many ISPs, particularly those operating outside of North America, transmission costs dominate the cost profile of the ISP's operation. If the cache performed at even 60 percent of a theoretical maximum caching performance, the ISP could reduce its external traffic volume requirements by some 13 percent. When the costs of caching are compared to the costs of transmission, this difference can be a significant one in the cost base of the ISP's operation.

For example, if the average cost of transmission is \$150 per gigabyte, and the ISP has a typical carriage profile of purchasing 1000 gigabytes per month from an upstream ISP with a 70-percent Web traffic profile, then a cache operating at a 25-percent byte hit rate can save the ISP a recurrent expenditure of \$26,250 per month. If the cache costs \$100,000

as a capital expenditure and \$2000 per month in operational costs to support the service, then a business case analysis would see the cache activity return some \$18,000 per month to the business, net of annualized capital and operational expenditures.

The other benefit is to the client, where the reduced network delay between the client and the local cache results in an increase in speed of Web page delivery for cached content.

The average size of a Web transaction is some 16 data packets within the TCP flow. Within a TCP slow-start flow-control process, the first cycle will transmit one packet and wait for an ACK. The reception of the ACK will trigger transmission of two more packets in the second round-trip cycle, and then the sender will await two ACKs. Reception of these two ACKs will trigger a further four packets in the third cycle and eight in the next cycle, and the remaining single packet in the fifth cycle. Therefore, allowing for optimal behaviour of the TCP slow-start algorithm, this average Web transaction takes some five round-trip times. If a user is located some distance away from the Web page, and the round-trip time to the source is 300 ms, the propagation delay of the page load will be 1.5 seconds. In comparison, if the round-trip time to the local Web cache is 2 ms, then the propagation delay of the page load will be 10 ms. These latency figures assume an uncongested network in both cases. In this case, as long as the Web cache search can complete within 1 second, the cache will appear to be far faster to the user.

A slightly different analysis is possible when comparing the performance of a cache configured at the headend of a cable-IP system versus the performance of direct access. The difference in latency in this case is due to both the closer positioning of the cache to the user and the greatly increased effective bandwidth from the cache to the user. A cache download can operate at speeds of megabits per second, as compared to kilobits or tens of kilobits per second when using dialup modem or ISDN services. For a 100K image download, the dial user may experience a 60-second delay, and the same delivery from a local cache via cable-IP may take less than half a second.

The trade-off with caching is that of balancing the the cost of carriage capacity, both in terms of monetary cost of the carriage and the performance cost of the transaction time of the application, against the cost of the use of caching. For non-North American ISPs, in which there is typically a large cache hit rate against North American server locations, the benefits of widespread use of caching are quite substantial. For cable-IP operators, the benefits of local cache operation lie in the ability to exploit the benefits of the very-high-speed final hop from the headend to the end user. For other ISPs, the benefits of caching may be less dramatic, but nevertheless, there are tangible positive outcomes of caching in terms of performance and cost that can be exploited.

As with direct-access models, this approach also has drawbacks. We have already noted the various ways in which the end-to-end model of Web content delivery has been exploited to provide time-based content, client-based content, and secure delivery of content. Caches insert themselves within the end-to-end semantics of the original transaction model, and intercept the transaction by presenting a proxy of the original endpoint. The content delivered from the cache is the content based on the time the cache undertook its request to the server, and the content delivered from the server is based on the server's view of the identity of the cache, rather than the identity of the end client.

With cached content in operation, the cached-content server no longer has an accurate picture of the number of times an item of content is viewed, and by whom. The server cannot authenticate the client, nor can the server deliver any information that is based on the supposed identity of the client. Equally, the client has potential problems, because the client may not be aware that the content has been delivered by the proxy cache. The content may not properly reflect the client's identity, and the information may be based on the security trust model of the server to the cache, rather than the server to the end client, and again the client may not be aware of such a change in security domains. If the content is time-dependent, the content will reflect the time at which the cache retrieved the content, rather than the time the client made the request.

All of this tends to suggest that caching is not a universally applicable tool. Part of the challenge in deploying cache servers is to understand the models of cache deployment and Web content delivery, and ensure that the cache does not intrude in ways that distort the integrity of content delivered to the end user.

Web Cache Hits Versus Web Server Hits

One of the biggest tensions is the balance between the cache operator's desire to maximize the hit rate of the cache system and the desire of many Web page publishers to maintain an accurate count on the number of hits of the page and from where those hits occur. In most cases, it is the requests that are of interest here, rather than the control of delivery of the content. The Web publisher is not necessarily interested in absorbing the hits for Web content. Indeed, many Web publishers see value in distributing the load of content delivery of fixed-content material further out toward the client base, rather than the Web publisher bearing the cost of the distribution load from the local site.

Static pages, composed of plain text and images, are readily cached. As a consequence, the original page publisher may not obtain an accurate count of the number of times the page was displayed by users if the Web server's log was analysed. Some Web page designers place information in the Web page directives; this information directs the Web cache server not to reuse a cached page. The most common way of doing this is to set the "Expires:" Web page information header to the current date and

time, so the next time the page is referenced, a new fetch will be undertaken. One of the more common hacks to cache servers to attempt to improve the hit rate is to allow this directive to be ignored.

This server hit-count problem has plagued cache deployment for many years now. Although there are real requirements in the areas of authentication and security, time-based content, and client-based content that mandate certain types of content being flagged as non-cacheable, much of the data that is marked as non-cacheable has been marked in this way simply for the server to capture the identity of the client. Such “cache-busting” practices are unnecessarily wasteful of network resources, and can overload the content server. There is an Internet Proposed Standard extension to HTTP^[6] intended to provide a “Meter” header, where a cache can communicate demographic information relating to client “hits” back to the original content server. The extension also proposes usage limiting, where a server can provide content with a limit on the number of times the information can be used by the proxy cache before revalidating the content with the server.

Web-Caching Models

There are many models of how to invoke a proxy cache.

Explicit Caching

Some proxy cache systems are deployed as a user-invoked option, in which the user nominates a cache server to the browser as a proxy agent, and the browser then directs all Web requests to the proxy cache. At any stage, the user can instruct the browser to turn off the use of the proxy cache, and request the browser to undertake the transaction directly with the client. Modern browsers when configured with a proxy cache may also use the approach of attempting direct access when a request via a proxy cache results in a fetch error. In the proxy cache mode of operation, the destination address of the underlying transport session is then the address of the cache server, while the HTTP content of the transaction remains unaltered. Such caches can be deployed within a client’s local network, with the intent of minimizing the amount of traffic passed to the external provider ISP. Additionally, The ISP can operate such a voluntary cache for use by its clients. If the ISP operates in this mode, the benefits to the user in using the cache need to be clearly stated and understood by both the client and the ISP, and the client must be made aware of the location of the cache in configuring his or her local browser.

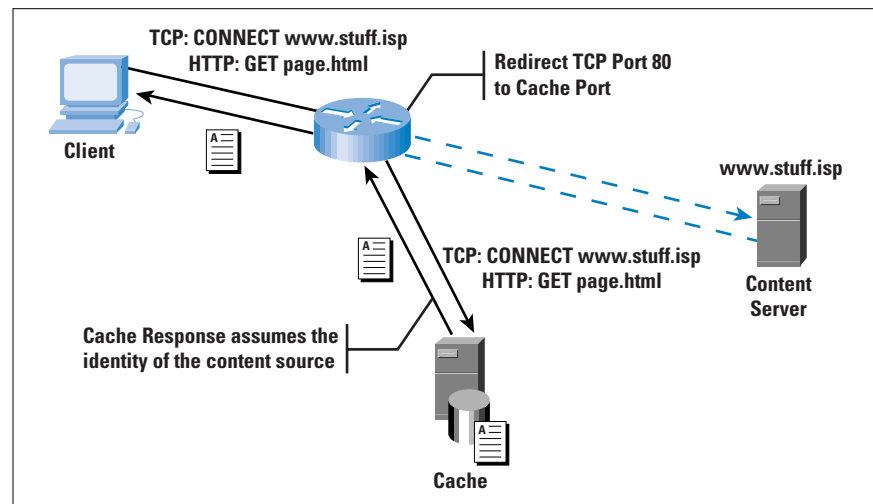
Forced Explicit Caching

Some ISPs, notably in the dialup service provider sector, operate in a highly cost-competitive market. In such a market service performance and service price are critical business factors, and the provider may choose to operate its network in a forced-cache mode. Here, all Web traffic on TCP port 80 (the port used by the HTTP Web transport protocol) is blocked from direct outbound access, and the ISP’s clients are forced to configure their browsers to use the provider’s cache for external Web access. This technique is commonly termed *forced caching*.

Transparent Caching

The use of a cache for all Web traffic also can be undertaken by the ISP, without the explicit configuration of the identity of the proxy cache into the user's browser. Irrespective of precisely how this setup is engineered, and there are numerous ways of engineering it, this technique is termed *transparent caching*. With transparent caching the user, and the user's browser, may not be explicitly aware that caching is being undertaken when processing the user's requests. Here the network has to intercept HTTP packets destined to remote Web servers, and present these packets to the proxy cache. Once the page is located, either as a cache hit or a cache miss, the cache must then respond to the original requestor by assuming the identity of the original destination (Figure 3).

Figure 3: Transparent Caching



It should be noted that no mechanism to date of explicit or transparent caching is completely transparent to both the Web client and the Web server. Where the Web server uses an end-to-end security access model the transparent cache may fail, because the cache will present its address as the source of the request, rather than that of a client. This scenario may result in a page-denied error to the cache request, whereas the client could have completed the transaction directly with the server. In those situations where the use of the cache is mandated, either through filters and a forcing function, or through transparent network redirection, there is no user-visible workaround to the error, and the level of user frustration with the entire cache service rises dramatically.

Under some circumstances it may be possible to work around transparent cache fetch errors. One approach is for a cache fetch error to trigger the cache subsystem to establish an HTTP session with the content server using the source address of the client, and then pass the original HTTP GET request to the server. The server's response is then passed to the client using a TCP bridge. (A TCP bridge is where the connecting device is required to translate the sequence numbers of the TCP headers between the two TCP sessions). Having the cache subsystem intercept the server's packets addressed to the client does require careful coordination with the cache router, and TCP bridging is also quite complex in its

operation, so such solutions tend to be somewhat unstable under load stress. An alternative approach is for the cache to pass a TCP RST back to the client, and instruct the cache router to insert a temporary entry in its redirection filter so that any subsequent TCP port 80 connection from the client to the server's address is not redirected to the cache.

If the sole benefit to the client is improved speed of response, then the ISP must understand that the performance of the Web cache systems must be continually tuned to be highly responsive to Web requests under all load conditions experienced by the ISP. Performance of cache hits must be maintained at a level consistently faster than the alternative of direct client access to the original client site. Performance of cache misses must be at a level that is not visibly slower than that of direct access to the original site. If the user's perception of performance of the cache drops, the benefit to the user also drops. In the case of user-selected caching, the users will turn off the cache option in their browser and return to a mode of direct access.

The business model of a cache is that the capital and operational costs associated with localizing traffic to the cache result in cost reductions to the ISP, when compared to the operation of a noncached network. These cost reductions can be passed on to all users through operation of the entire service at a lower price point or selectively passed on to those clients who make use of the cache through some form of cache-use tariff. The generic model of applying the cost reduction to the ISP's service tariff is certainly an advantage in a price-competitive marketplace. However, unless the performance of the cache is consistently very high, and the transparency of the cache is close to perfect, each individual user may attempt to use direct-access methods.

The alternate business model is to pass on the marginal cost savings to those clients who make use of the cache, and at a level that corresponds to the client's use of the cache and its effectiveness in operating at a high cache hit rate. If, for example, the ISP uses a charging model that includes a tariff component based on the amount of data delivered to the client during the accounting period, this tariff component could be adjusted by the amount of use the client made of the cache system and the relative operating efficiency of the cache in generating cache hits.

As an example, if traffic is tarified at \$100 per gigabyte as delivered to the customer, a discounted value can be derived for traffic delivered from the Web cache. If the average cache byte hit rate is 30 percent, then after factoring in the costs of capital equipment and operational support, the traffic from the cache could be tarified at \$80 per gigabyte. Here, the benefit of using the Web cache is passed directly to those clients who make use of the cache, who both enjoy lower tariffs in direct proportion to their use of the cache and derive superior performance through using the cache. The accounting for this marketing model is certainly a more involved process, involving additional accounting systems and processing to undertake an accurate per-client view of cache usage.

It is becoming increasingly evident that a robust business model associated with a model of discretionary use of a Web cache is that of access to a lower unit price of traffic. In this way, the user sees the incentive of immediate financial benefit in choosing to use the cache system. When the provider deploys transparent or forced caching, translating the benefits of caching into an overall reduced tariff structure for all clients is a more robust business model.

Web-Cache Systems

Cache systems can take a variety of forms. The original Web server from CERN, the original location of the development of Web software, allowed a mode of proxy behaviour. This cache server model was developed significantly in the Harvest Project, a research project at the University of Colorado. As an evolutionary path, the *Harvest* cache server is being further developed within the scope of the development of the *Squid* cache server software and the associated *Internet Caching Protocol* (ICP).

Currently numerous freely available proxy cache systems are available, such as Squid, and many systems are available commercially, such as the Cisco Systems *Cache Engine*. Some of these systems are software packages that operate on a conventional operating system platform, while some use a customized platform kernel, which is optimized for the demands of a cache-delivery environment.

Many of the characteristics of Web caching systems are relevant to the performance of the caching environment. The first is the *size* of the cache server. The relationship between the size of the cache and its hit rate is not a linear relationship. For typical patterns of Web use generated from a relatively large user population, a cache of 1 gigabyte or so will yield reasonable hit rates. Further increase of the cache size will yield incremental improvements in the cache hit rate, where the incremental rate is best described by a negative exponential relationship. Thus, caching systems with 10 gigabytes of storage do not produce performance characteristics markedly different from larger 100-gigabyte caching systems. No objective best size of cache system can be determined, because local environments differ, but every environment exhibits the law of diminishing returns, in which the addition of further cache capacity yields no tangible difference in the cache effectiveness. Large caches take some time, in the order of days or even weeks, to build up a sufficiently large repository of cached data to produce an improved cache hit rate. Generally, 10- to 100-gigabyte cache systems provide extremely effective cache performance, as long as the cache is allowed to stabilize for some weeks following startup. Memory demands in a cache also need to be carefully configured. The URL index of the storage system is stored in memory in most cache architectures in order to perform fast cache lookups, so that the more disk storage configured, the larger the memory requirements.

The next parameter is the *number of simultaneous cache requests* that the cache server can manage efficiently. Note that this metric is different from the number of requests per second that the cache server can manage. The number of simultaneous sessions that the cache server can support is related to the amount of resources allocated to the cache request and the total resource capacity of the box.

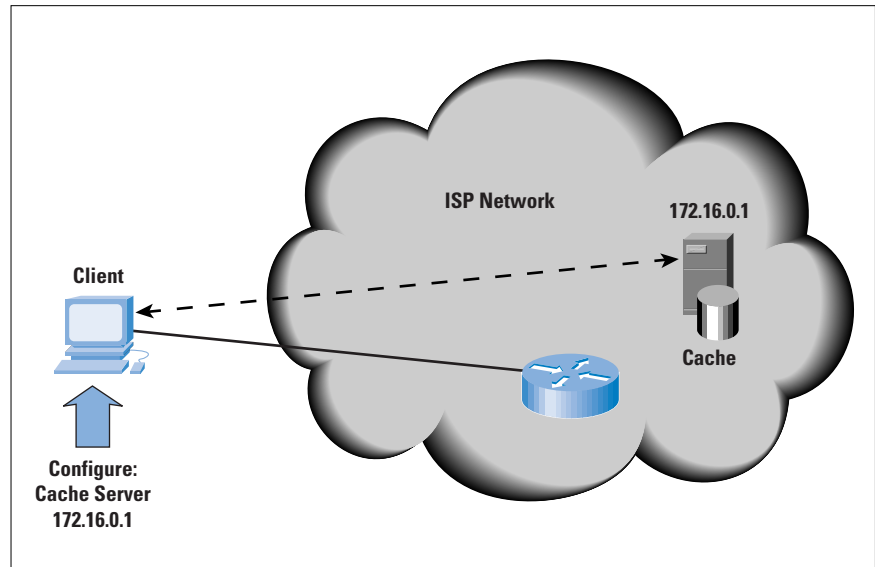
The environment of deployment is very relevant to the performance of the cache environment. The related metric to the number of simultaneous requests that can be managed is the average time to process a request. Combining these two metrics provides the number of requests per second that the cache system can process. The same unit will have a different performance metric of requests per second when deployed in different parts of the Internet. If the cache system is deployed with a satellite-based feed, then the average time to process a cache miss is considerably longer because of the higher latency of the satellite path. This scenario leaves the process of managing the original request open for a longer period, blocking other requests from using this process slot. If the same unit is deployed in a location where cache misses take fractions of a second to process, the process slot can be quickly reused. Each active client connection also consumes memory, and the client connection will remain open for as long as it takes to complete the Web transaction, either for a hit or a miss. The greater the mean round-trip delay for a miss, the greater the number of concurrent active sessions held in the cache. Similarly, the greater the number of low-speed modem or wireless-based clients, the greater the number of concurrent active sessions in the cache. Whether the client operates in transparent mode or in explicit proxy caching mode is also an important consideration. Browser clients use an explicit proxy cache with a persistent connection, while if the cache is a transparent cache, the cache will see clients bring up and drop HTTP connections each time the base URL changes. This session reestablishment, together with the additional Domain Name System (DNS) resolution load imposed on the client, can add up to half a second to the transparent cache response time as compared to the explicit cache response.

Web Cache Deployment Models

In this section we first examine scaling issues for explicitly referenced cache configurations, and then look at the changes to the model introduced through transparent caching.

The simplest deployment model of an explicit cache is that of deployment of a single cache system as a browser-selectable resource. This system can be deployed within an ISP's server environment with a TCP port-80 interface opened for client access. Such a deployment model is shown in Figure 4.

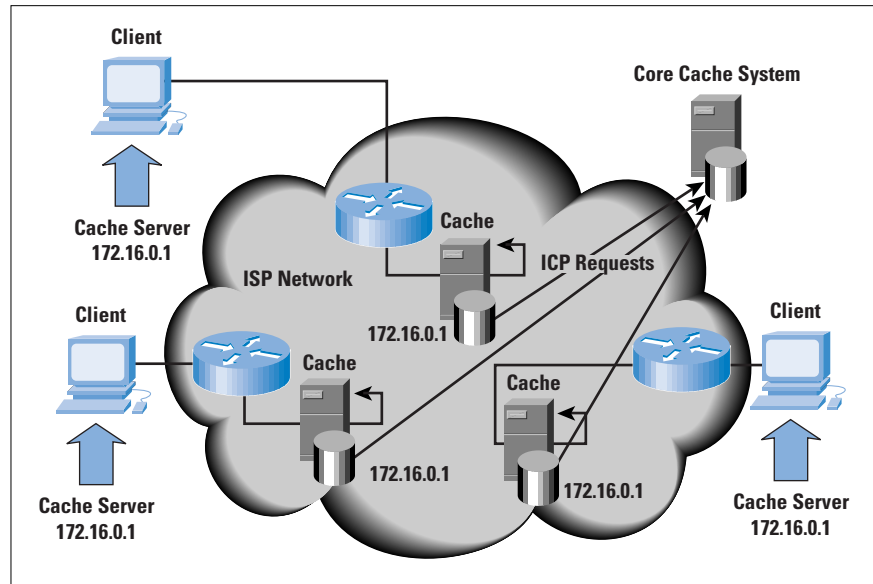
Figure 4: A Selectable
Web Cache



Single Web proxy cache systems can be placed under some significant load, and an overloaded and poorly performing cache is perhaps worse than no cache at all. However, scaling this deployment model can prove challenging. Where an ISP operates multiple access points, or points of presence (POPs), one scaling solution is to deploy a server at each POP and use the same IP address for each server. This solution allows the ISP to provide a consistent configuration to all clients and to augment capacity at any location seamlessly. If the cache itself is responsible for advertising the common IP address into the routing system, the caches can also act in a mutual backup role. Failure of a single server will shut down the local route advertisement. Traffic directed to this address will then be carried by the routing system to the next closest proxy cache. There may be some level of TCP session resets for sessions that were active on the failed unit, but in all other respects the switchover is seamless to the client base, and the recovery of an operational state among a set of such servers can be left to the routing system. This deployment model is indicated in Figure 5. Such servers can be configured as a set of local satellite systems to a larger caching core, using an *Internet Cache Protocol* (ICP) configuration to set up a caching hierarchy.

ICP is a lightweight message format for communicating between Web caches^[7]. The message format is a simple two- packet exchange, where a Web cache passes a URL query to another cache. The response is either a hit or a miss, indicating the presence of the URL object on the remote cache. On top of this protocol can be constructed cache hierarchies, to allow multiple neighboring caches to pool their resources effectively.

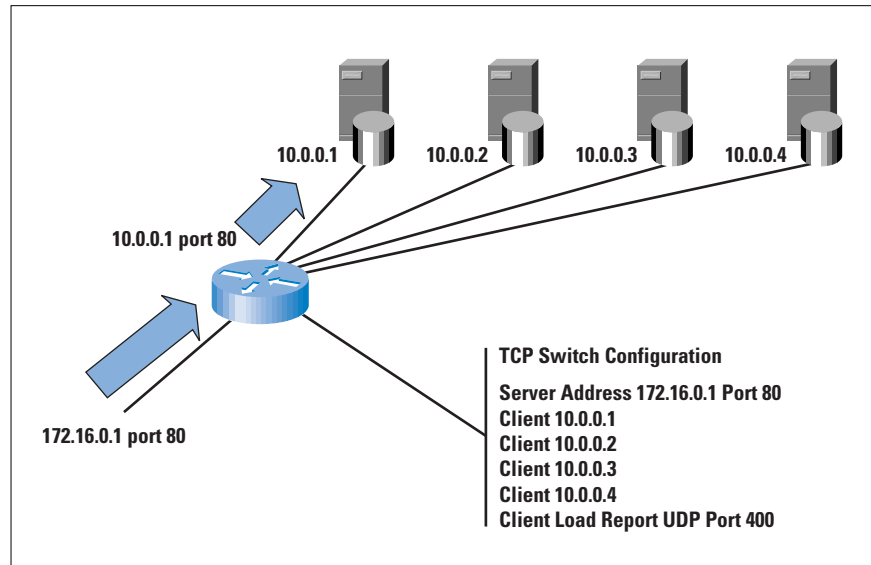
Figure 5: Replicated Web Caches



The proposed mode of configuration of caches is into a tree-structured hierarchy^[8]. In such a hierarchy every participating cache is organized with a connection of neighboring peers and an ICP parent. When a cache request cannot be serviced from the local cache, the cache first uses a set of local configuration rules to determine if the server is local. If so, the cache queries the server directly for the content. If the server is remote, the cache issues a set of simultaneous ICP queries to all its cache peers. If any peer responds with an ICP hit, the cache then requests the peer to provide the referenced content. If all peers respond negatively to the ICP query, or a two-second timeout elapses, the cache then requests the URL from its designated parent. The parent may use a peer referral, or the parent may refer the query to its parent, or perform a cache retrieval on behalf of the original request. The intent of this mode of operation is to use a lightweight query response protocol to allow a local collection of caches to pool their cached data. ICP has also been used with additional policy constraints, although the protocol itself is not capable of describing or carrying overly complex retrieval policies. Other intercache protocols are available, including the *Hyper Text Caching Protocol* (HTCP) and the *Cache Array Routing Protocol* (CARP), which offer functionality in terms of intercache cooperation similar to that of ICP^[9].

Another scaling measure is to alter the single server to multiple servers, using a TCP-based, load-sharing mechanism in the switching system to ensure that the servers are evenly loaded. This setup is shown in Figure 6. Such a simple load-sharing system may even the load on each server, but it will cause each server to act independently of its sibling servers. It is essential in such an environment to use ICP to coordinate the servers so that they will refer to each other before initiating a new fetch from the content server.

Figure 6: Load-Balancing Web Caches



In such a configuration each cache will contain content also held in neighboring caches. Although this scenario may allow some form of server load balancing, particularly when the servers continually communicate their current load conditions to the load-balancing switch, there is still some inefficiency in the cache farm operation through the potential replication of content on each of the component caches. One direction of scaling the cache servers is to take a collection of cache servers and allow each cache server to specialize in the content it holds. However, the outer TCP destination address does not help the server determine which URL is being requested. In an explicit cache configuration, the browser is directing the TCP session to the externally advertised TCP address of the server farm. The URL information is embedded within the HTTP payload. Some developments have been made in this area, where, with a combination of TCP spoofing and TCP session bridging, a server switch can select the appropriate cache for each HTTP-referenced URL, and then logically connect the client's TCP session to a TCP session to the selected cache to deliver the URL to the client.

Transparent caching presents some further deployment challenges. The functional requirement is to pass all Web requests through a proxy cache server without the explicit knowledge of the client. Two generic techniques exist to achieve this goal:

- *Inline caches:* The first of these approaches is to pass all traffic through a two-port cache server. All non-HTTP traffic is simply passed straight through the device without alteration. HTTP traffic is intercepted and passed to a cache module. The major concern with this approach is the introduction of a single point of failure with an active network element. Any failure of the cache may well prevent all further traffic from entering or leaving the served subnetwork.

- *Redirection caches:* A technique that does not place the cache as a critical point of potential failure is to use policy redirection within the router, redirecting all port-80 traffic to the attached cache. Normally such a policy redirect would infer that the cache is located one hop away from the router, so that such a redirection is normally a local solution. Redirection to a tunnelled interface does allow some greater flexibility in this setup, and the one cache farm could, in such an approach, service a collection of redirecting routers. The failure mode of this form of operation remains a concern, because the redirection mechanism in the router would not normally be aware of the operational status of the cache.

Transparent caches need to ensure that the full URL is inserted into the HTTP level request. When the browser assumes that the request is directed to the content server, the GET request may specify a URL relative to the server. In such cases, the transparent server will need to perform a DNS lookup of the destination IP address of the TCP session in order to reconstruct the complete URL.

Although the DNS lookup does have some performance implications to transparent caches, the major issue for transparent caches is to devise a fail-safe mechanism, so that if the cache server fails for any reason, the caching redirection is disabled. One solution is to use a redirection function within the router in conjunction with a keepalive-based Web cache management protocol. This scenario is the basis of the *Web Cache Coordination Protocol* (WCCP)^[10]. WCCP also adds the ability to load share across multiple cache servers through content distribution. Transparent caching assists in this task because the destination address in the IP packets can be used as the basis of the cache selector. The keepalive exchange between the router and the cache server system allows the router to cease redirecting Web traffic upon failure of the servers.

Alternative solutions rely on the cache itself participating in a local routing environment. The redirecting router uses policy-based redirection to forward all port-80 traffic to an address announced by the cache system at a high routing priority. The same address is also announced by the default path router at a low routing priority. Failure of the cache system will result in a withdrawal of the high-priority route, and while the redirection will remain in place on the router, the redirection will be in the direction of the default route.

Another challenge is to process cache misses at a speed comparable with normal noncached Web retrieval. A process of pulling the document into the cache and then serving the document to the original requestor does not meet that objective. The transparent cache has to feed the document to the requestor while simultaneously creating a stored copy for subsequent cache serving.

However, the largest challenge to the transparent cache is that it can serve only documents that are not dependent on the identity of the requestor being preserved. Web servers that use an end-to-end model of access, based on source address identification, or Web servers that attempt to present different documents to the client based on the client's source address, do not fit within the transparent caching model. There is much interest in solutions that allow a transparent cache to effectively shut down in the case of a Web retrieval error, and allow the original requestor the ability to conduct a HTTP conversation directly with the server in such situations. Although there is interest in a network-only solution, it appears at this stage that some level of assistance from the browser may be required. One model of operation is that a transparent cache records the network-level flow identification of a failed Web retrieval, and passes a retry signal back to the requesting browser, and also passes this flow identifier back to the redirector as a temporary filter entry. When the requestor retries the query, as per the signal from the cache, the redirecting router will refrain from redirecting the flow to the cache, and allow an end-to-end session to operate.

Accounting for Web Cache Use

These deployment systems allow for user-optional cache configuration. If the ISP wants to account for the use of the cache, then the cache server or the switch that feeds the cache server must play an active role in accounting collection.

If every network address is uniquely advertised to the ISP by a particular client, then the task of accounting for cache use can be performed using the logged records of the cache system itself. Because every IP address can be uniquely mapped to an ISP client, it is possible to also associate the volume of bytes delivered by the cache to the identified client.

Unfortunately, two factors make this supposition of address uniqueness somewhat weak. First, dialup address assignment implies that the association of an IP address to a client is held only within dialup accounting records in the first instance, and the binding is valid only between the times referenced in the start and stop records. This scenario can be configured into an accounting model by simultaneously processing the dial accounting records when attempting to associate a particular IP address at a particular time to a client.

The second factor is slightly more challenging. For an ISP that offers permanent access transit services, the potential exists that any particular IP address may not be uniquely routed. Normally, such multiple access environments are part of a Border Gateway Protocol (BGP)-based interaction with multiple clients. Knowing the IP address of the query agent is not enough. Ascertaining the next-hop Autonomous System (AS) number as well as the IP address is now necessary to determine the client using the cache.

The implication is that the accounting records now need to be generated on the router that is also the entry point to the cache. In addition, the router must participate in the interior BGP (iBGP) core mesh to maintain current AS path-selection choices. Given the considerable overheads that such an engineering design entails, an alternative approach is to restrict the cache accounting role to account for those cases where the cache client is readily identified. A common measure is that the lower tariff is available only to customers who are “singly homed” with the ISP. Not only is this a strong market incentive for customer loyalty, it also allows simple engineering solutions for cache accounting, because the lookup from the IP address in the cache log to a customer account is then relatively straightforward. Such measures allow a cache-use tariff to be very competitively positioned in the market.

As well as accounting issues, another component for the consideration of optional use of Web caches is that of the necessity of restricting the use of the cache to clients of the ISP. The motive for so doing is to ensure that the cache is available only to clients of the service and not to clients of peer ISPs. It may not be an issue worth the effort of solving, and the first questions ISPs should ask is, “To what extent does this happen, and what impact does it have on the operation of the Web cache systems?” In most cases, the accounting of cache usage may reveal that this issue is one of negligible proportions, and any effort expended in devising an engineering solution would far outweigh the loss to the ISP through such use of the service.

If the measurement of such usage is considered sufficient to warrant engineering solutions, then the mechanisms available to the ISP are to ensure that the Web cache access is filtered at the edges of the ISP network and to ensure that access is possible only by ISP clients, or that the address of the cache is not exported in the routing system to peer ISPs or upstream ISPs.

Further Deployment Challenges

It is highly likely that further development will occur with cache servers in the near future. Large-scale backbone IP networks that use OC-3c (155 Mbps) or OC-12c (622 Mbps) transport cores may carry tens of thousands of requests per second. Designing transparent caches that fit within a transport core at such a scale does present dramatic scaling issues in terms of cache system performance. This factor continues to elude many of today’s products available on the market. The generic architecture today is to use a cache network that attempts to place the cache systems closer to the access edge of the network, where the Web request volumes are within the scale of today’s cache systems.

Transparency of the cache remains an issue, and it is perhaps an area of further refinement within the specification of the underlying HTTP Web server protocol, as well as further refinement of the operation of Web browsers and transparent cache systems. A potential implementation within Web browsers may allow the user to state the acceptability of using a cache to complete a request, and allow noncache Web page retrieval attempts on cache failure, in the same way that the provider can use page expiration directives to direct a cache not to store the presented data.

References and Further Reading

- [1] Huston, G. *ISP Survival Guide*, ISBN 0-471-31499-4, John Wiley & Sons, November 1998.
A more comprehensive view of the technology, business, and strategy behind the Internet service sector.
- [2] Clark, D.D. "The Design Philosophy of the DARPA Internet Protocols," Proceedings of SIGCOMM 88, *ACM Computer Communications Review* (CCR), Volume 18, Number 4, August 1988, pp. 106–114 (reprinted in ACM CCR Volume 25, Number 1, January 1995, pp. 102–111).
The original paper describing the end-to-end design philosophy used within the Internet protocols.
- [3] Carpenter, B., Ed. "Architectural Principles of the Internet," RFC 1958, Informational RFC, June 1996.
A summary of the design principles underlying the current Internet architecture.
- [4] Berners-Lee, T., et al. "Hypertext Transfer Protocol—HTTP/1.0," RFC 1945, Informational RFC, May 1996.
The specification of Version 1.0 of the HTTP protocol.
- [5] Fielding, R., et al. "Hypertext Transfer Protocol—HTTP/1.1," RFC 2616, Draft Standard RFC, June 1999.
The specification of Version 1.1 of the HTTP protocol.
- [6] Mogul, J., and Leach, P. "Simple Hit-Metering and Usage-Limiting for HTTP," RFC 2227, Proposed Standard RFC, October 1997.
A proposed extension to HTTP to allow a content server to receive hit reports from a proxy cache.
- [7] Wessels, D., and Claffey, K. "Internet Cache Protocol (ICP), Version 2," RFC 2186, Informational RFC, September 1997.
A description of the ICP protocol.
- [8] Wessels, D., and Claffey, K. "Application of Internet Cache Protocol (ICP), Version 2," RFC 2187, Informational RFC, September 1997.
A description of the structure of cache hierarchies, and their ICP-based interaction.

- [9] Melve, I. "Inter Cache Communications Protocols," Internet Draft, Work in progress, **draft-melve-intercache-comproto-00.txt**, November 1998.
An overview of intercache communications protocols currently available, and a collection of references that describe these protocols in further detail.
- [10] Cieslak, M., and Foster, D. "Web Cache Coordination Protocol V1.0," Internet Draft, Work in progress, **draft-ietf-wrec-web-pro-00.txt**, June 1999.
A description of Version 1 of the WCCP protocol to support the operation of transparent caches. The protocol defines the interaction between a router and a neighboring cache system.
- [11] "Squid Internet Object Cache"—Resource Web page.
<http://squid.nlanr.net>
A very useful page of resources and references related to the Squid implementation of Web caching.
- [12] "Distribution of Stored Information on the Web," Online Tutorial, Ross, K., Institut Eurecom, October 1998. Available at:
<http://www.eurecom.fr/~ross/CacheTutorial/DistTutorial.html>
A good overview of proxy caching technologies, and also a good analysis of their efficiency of operation.
- [13] Stallings, W. "SSL: Foundation for Web Security," *The Internet Protocol Journal*, Volume 1, Number 1, June 1998.

[This article is based in part on material in The ISP Survival Guide, by Geoff Huston, ISBN 0-471-31499-4, published by Wiley in 1998^[1]. Used with permission].

GEOFF HUSTON holds a B.Sc and a M.Sc from the Australian National University. Closely involved with the development of the Internet for the past decade, particularly within Australia, he was responsible for the initial build of the Internet within the Australian academic and research sector. Huston is currently the Chief Technologist in the Internet area for Telstra. He is also an active member of the IETF, and is the chair of the Internet Society Board of Trustees. He is author of *The ISP Survival Guide*, ISBN 0-471-31499-4, and coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. Both books are published by John Wiley & Sons, E-mail: **gih@telstra.net**

Gigabit Ethernet

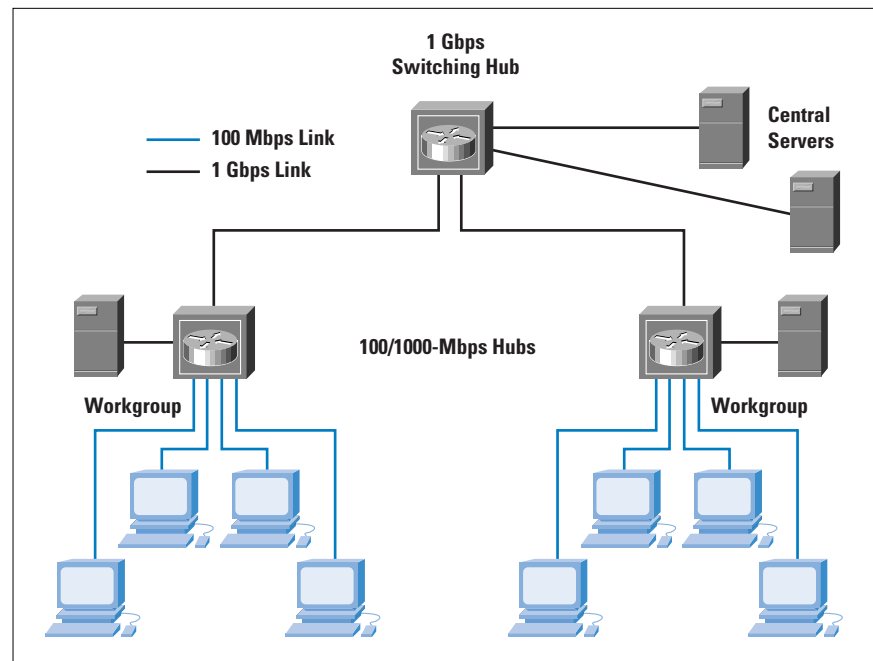
by William Stallings

In late 1995, the IEEE 802.3 committee formed a High-Speed Study Group to investigate means for conveying packets in Ethernet format at speeds in the gigabit-per-second range. A set of 1000-Mbps standards have now been issued.

The strategy for Gigabit Ethernet is the same as that for 100-Mbps Ethernet. While defining a new medium and transmission specification, Gigabit Ethernet retains the carrier sense multiple access collision detect (CSMA/CD) protocol and frame format of its 10- and 100-Mbps predecessors. So it is compatible with the slower Ethernets, providing a smooth migration path. As more organizations move to 100-Mbps Ethernet, putting huge traffic loads on backbone networks, demand for Gigabit Ethernet is intensifying.

Figure 1 shows a typical application of Gigabit Ethernet. A 1-Gbps LAN switch provides backbone connectivity for central servers and high-speed workgroup switches. Each workgroup LAN switch supports both 1-Gbps links, to connect to the backbone LAN switch and to support high-performance workgroup servers, and 100-Mbps links, to support high-performance workstations, servers, and 100-Mbps LAN switches.

Figure 1: Example Gigabit Ethernet Configuration



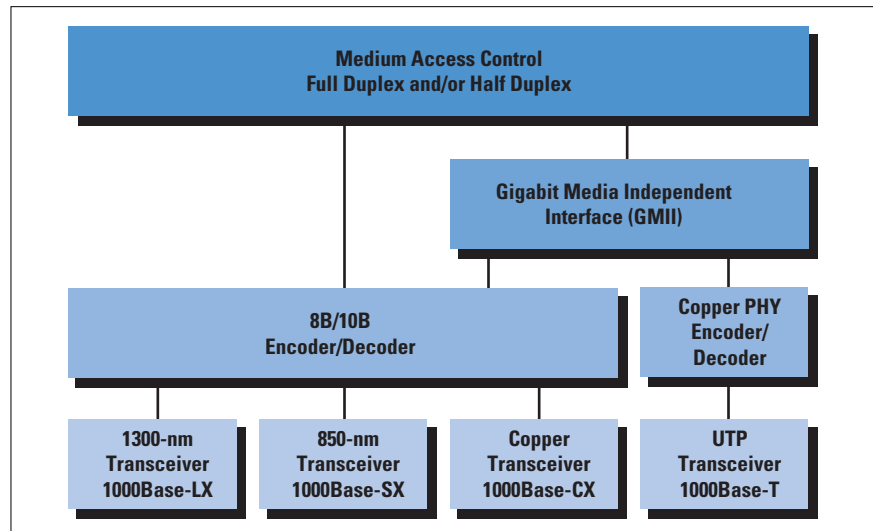
Protocol Architecture

Figure 2 shows the overall protocol architecture for Gigabit Ethernet. The Media Access Control (MAC) layer is an enhanced version of the basic 802.3 MAC algorithm. A separate gigabit medium-independent interface (GMII) has been defined and is optional for all the medium options except unshielded twisted-pair (UTP).

The GMII defines independent 8-bit-parallel transmit and receive synchronous data interfaces. It is intended as a chip-to-chip interface that lets system vendors mix MAC and physical sublayer (PHY) components from different manufacturers.

Two signal encoding schemes are defined at the physical layer. The 8B/10B scheme is used for optical fiber and shielded copper media, and the pulse amplitude modulation (PAM)-5 is used for UTP.

Figure 2: Gigabit Ethernet Layers



Media Access Layer

The 1000-Mbps specification calls for the same CSMA/CD frame format and MAC protocol as used in the 10- and 100-Mbps versions of IEEE 802.3. For traditional Ethernet hub operation, in which only one station can transmit at a time (half-duplex), the basic CSMA/CD scheme has two enhancements:

- *Carrier extension:* Carrier extension appends a set of special symbols to the end of short MAC frames so that the resulting block is at least 4096 bit-times in duration, up from the minimum 512 bit-times imposed at 10 and 100 Mbps. This extension makes the frame length of a transmission longer than the propagation time at 1 Gbps.
- *Frame bursting:* This feature allows for multiple short frames to be transmitted consecutively, up to a limit, without relinquishing control for CSMA/CD between frames. Frame bursting avoids the overhead of carrier extension when a single station has a number of small frames ready to send. extension when a single station has numerous small frames ready to send.

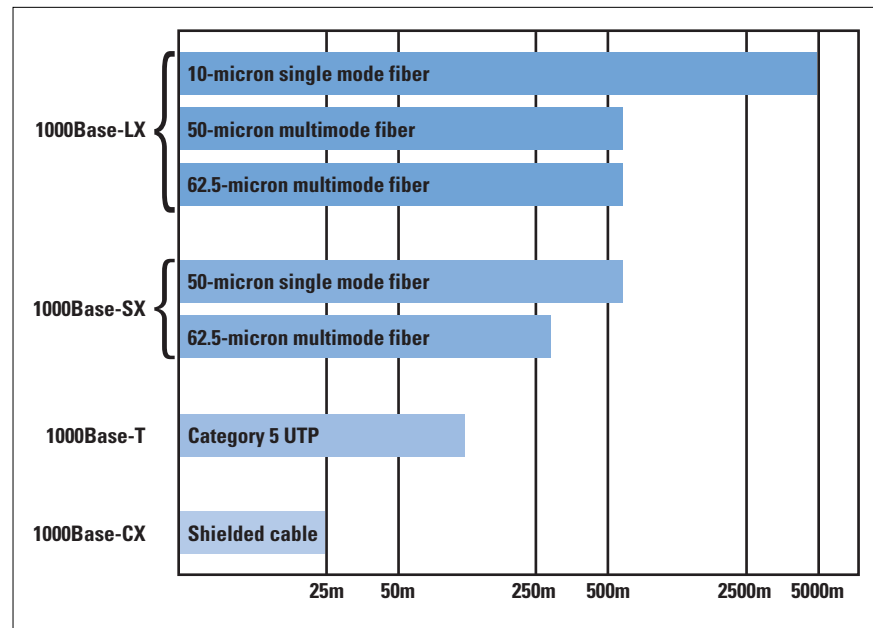
With a LAN switch (full-duplex operation), which provides dedicated rather than shared access to the medium, the carrier extension and frame bursting features are not needed. They are unnecessary because data transmission and reception at a station can occur simultaneously without interference and with no contention for a shared medium. All the gigabit products on the market use a switching technique, and so do not implement the carrier extension and frame bursting.

With a switching technique, full-duplex operation is employed, and the CSMA/CD protocol is not needed. The gigabit specification expands on the pause protocol that is defined for 100-Mbps Ethernet by allowing asymmetric flow control. Using the autonegotiation protocol, a device may indicate that it may send pause frames to its link partner but will not respond to pause frames from its partner.

Physical Layer

The current 1-Gbps specification for IEEE 802.3 includes the following physical-layer alternatives (Figure 3):

Figure 3: Gigabit Ethernet Media Options (log scale)



- **1000Base-LX:** This long-wavelength option supports duplex links of up to 550 m of 62.5- μ m or 50- μ m multimode fiber or up to 5 km of 10- μ m single-mode fiber. Wavelengths are in the range of 1270 to 1355 nm.
- **1000Base-SX:** This short-wavelength option supports duplex links of up to 275 m using 62.5- μ m multimode or up to 550 m using 50- μ m multimode fiber. Wavelengths are in the range of 770 to 860 nm.
- **1000Base-CX:** This option supports 1-Gbps links among devices located within a single room or equipment rack, using copper jumpers (specialized shielded twisted-pair cable that spans no more than 25 m). Each link is composed of a separate shielded twisted-pair running in each direction.
- **1000Base-T:** This option makes use of four pairs of Category 5 unshielded twisted-pair copper wires to support devices over a range of up to 100 m.

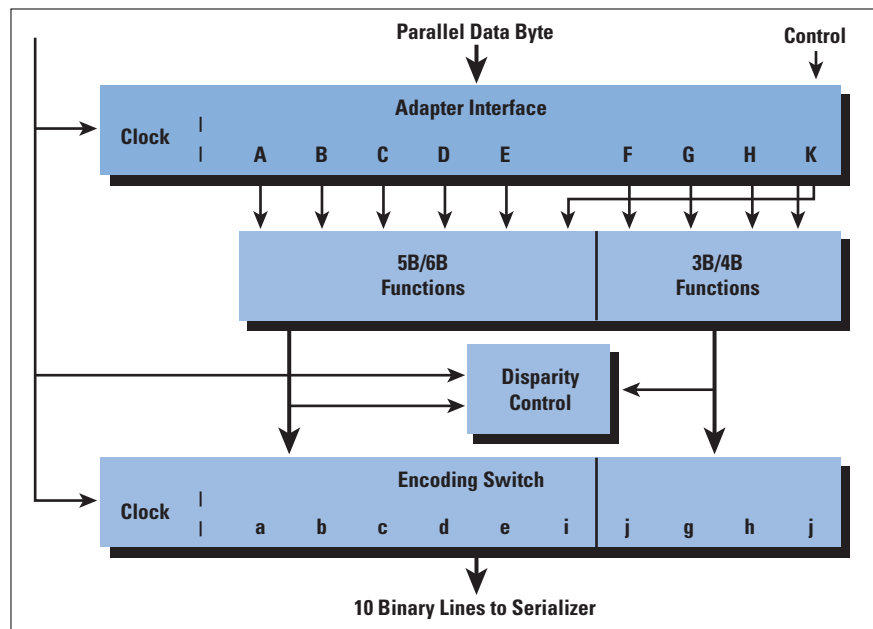
Digital Signal Encoding Techniques for Gigabit Ethernet

The encoding scheme used for all the Gigabit Ethernet options except twisted-pair is 8B/10B. This scheme is also used in Fibre Channel. With 8B/10B, each 8 bits of data is converted into 10 bits for transmission. The 8B/10B scheme was developed and patented by IBM for use in its 200-megabaud ESCON interconnect system.

- The developers of this code list the following advantages:
- It can be implemented with relatively simple and reliable transceivers at low cost.
- It is well balanced, with minimal deviation from the occurrence of an equal number of 1 and 0 bits across any sequence.
- It provides good transition density for easier clock recovery.
- It provides useful error-detection capability.

The 8B/10B code is an example of the more general $mBnB$ code, in which m binary source bits are mapped into n binary bits for transmission. Redundancy is built into the code to provide the desired transmission features by making $n > m$. Figure 4 illustrates the operation of this code. The 8B/10B code actually combines two other codes, a 5B/6B code and a 3B/4B code. The use of these two codes is simply an artifact that simplifies the definition of the mapping and the implementation; the mapping could have been defined directly as an 8B/10B code. In any case, a mapping is defined that maps each of the possible 8-bit source blocks into a 10-bit code block. There is also a function called *disparity control*. In essence, this function keeps track of the excess of zeros over ones or ones over zeros. An excess in either direction is referred to as a disparity. If there is a disparity, and if the current code block would add to that disparity, then the disparity control block complements the 10-bit code block. This complement has the effect of either eliminating the disparity or at least moving it in the opposite direction of the current disparity.

Figure 4: 8B/10B Encoding



The encoding mechanism also includes a control line input, K, which indicates whether the lines A through H are data or control bits. In the latter case, a special nondata 10-bit block is generated. A total of 12 of these nondata blocks are defined as valid in the standard. These blocks are used for synchronization and other control purposes.

For 1000Base-T, the encoding scheme used is PAM-5, over four twisted-pair links. Therefore, each link must provide a data rate of 250 Mbps. PAM-5 provides better bandwidth utilization than simple binary signaling by using five different signaling levels. Each signal element can represent two bits of information (using four signaling levels). In addition, a fifth signal level is used in a forward error correction scheme.

References

A good tutorial on Gigabit Ethernet is [1]. The Gigabit Ethernet Alliance is at <http://www.gigabit-ethernet.org>

- [1] Frazier, H., and Johnson, H. "Gigabit Ethernet: From 100 to 1,000 Mbps." *IEEE Internet Computing*, January/February 1999.
- [2] *Gigabit Ethernet: Technology and Applications for High-Speed LANs* by Rich Seifert, ISBN 0-201-18553-9, Addison-Wesley, 1998. (Reviewed in *The Internet Protocol Journal*, Volume 1, Number 2, September 1998.)

WILLIAM STALLINGS is a consultant, lecturer, and author of more than a dozen books on data communications and computer networking. He has a PhD in computer science from M.I.T. His latest book is *Data and Computer Communications, Sixth Edition* (Prentice Hall, 1999). His home in cyberspace is <http://www.shore.net/~ws> and he can be reached at ws@shore.net

One Byte at a Time: Is Your FTP Active or Passive?

by Thomas M. Thomas, NetCerts

What many people don't know is that the *File Transfer Protocol* (FTP) has multiple modes of operation that can dramatically affect its operation and, as a result, the security of your network. These modes of operation determine whether the FTP server or FTP client initiates the TCP connections that are used to send information from the server to the client. The FTP protocol supports two modes of operation, as follows:

- The first FTP mode of operation is known as *normal*, though it is often referred to as *active*. This mode of operation is typically the default.
- The second FTP mode of operation is known as *passive*.

In active (normal) FTP, the client opens a control connection on port 21 to the server, and whenever the client requests data from the server, the server opens a TCP session on port 20. In passive FTP, the client opens the data sessions, using a port number supplied by the server.

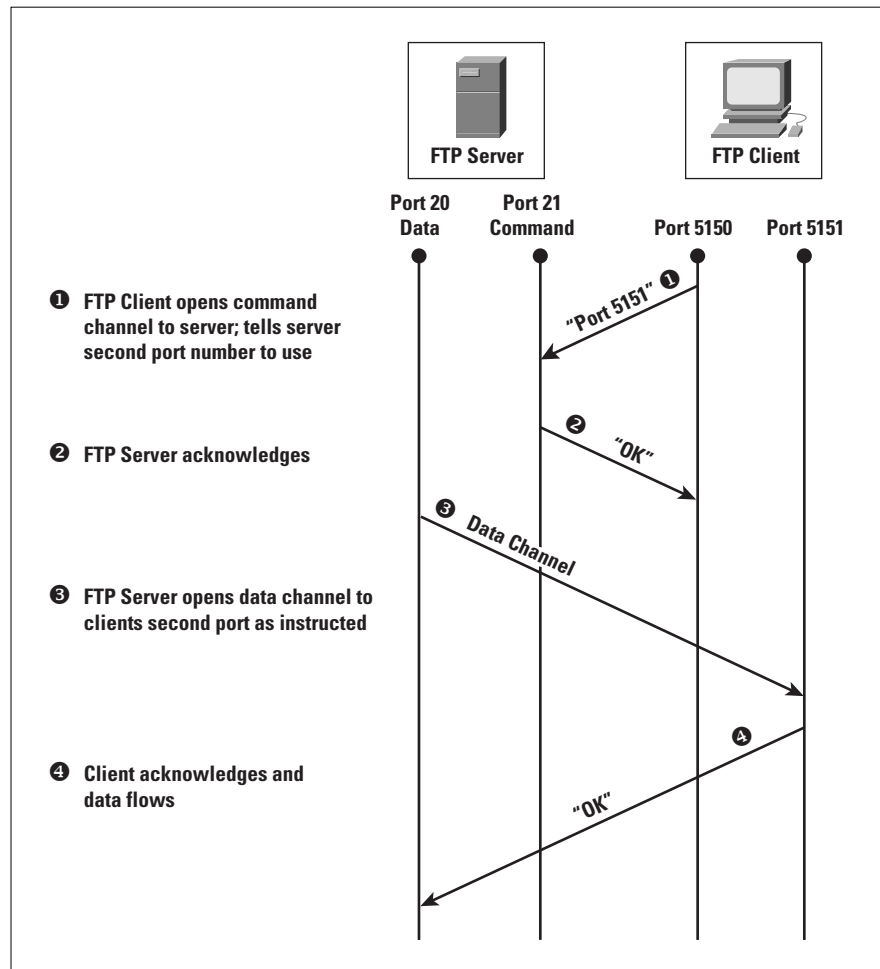
Active FTP Operation

The active mode of operation is less secure than the passive mode. This mode of operation complicates the construction of firewalls, because the firewall must anticipate the connection from the FTP server back to the client program. The steps of this mode of operation are discussed below and are shown in Figure 1.

- The client opens a control channel (port 21) to the server and tells the server the port number to respond on. This port number is a randomly determined port greater than 1023.
- The server receives this information and sends the client an acknowledgement "OK" (ack). The client and server exchange commands on this control connection.
- When the user requests a directory listing or initiates the sending or receiving of a file, the client software sends a "PORT" command that includes a port number > 1023 that the client wishes the server to use for the data connection.
- The server then opens a data connection from port 20 to the client's port number, as provided to it in the "PORT" command.

The client acknowledges and data flows.

Figure 1: Active-Mode
FTP Connection



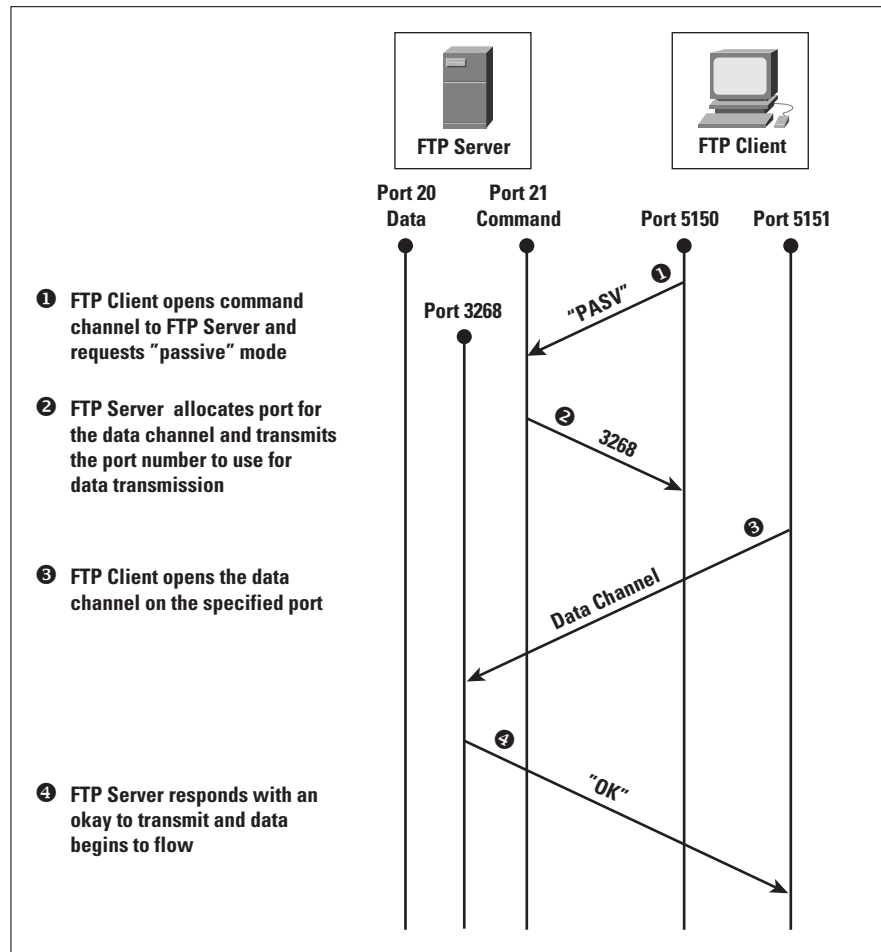
Passive FTP Operation

This mode of operation is assumed to be more secure because all the connections are being initiated from the client, so there is less chance that the connection will be compromised. The reason it is called passive is that the server performs a "*passive open*." The steps of this mode of operation are discussed below and are shown in Figure 2.

- In passive FTP, the client opens a control connection on port 21 to the server, and then requests passive mode through the use of the "PASV" command.
- The server agrees to this mode, and then selects a random port number (>1023). It supplies this port number to the client for data transfer.
- The client receives this information and opens a data channel to the server-assigned port.

The server receives the data and sends an "OK" (ack).

Figure 2: *Passive-Mode
FTP Connection*



References

- [1] *Internetworking With TCP/IP, Volume 1: Principles, Protocols, Architecture, Third Edition*, by Douglas E. Comer, ISBN 0-13-216987-8, Prentice Hall, 1995.
- [2] R. Braden, "Requirements for Internet hosts—application and support," RFC 1123, October 1989.
- [3] S. Bellovin, "Firewall-Friendly FTP," RFC 1579, February 1994.
- [4] P. Deutsch, A. Emtage, A. Marine, "How to Use Anonymous FTP," RFC 1635, May 1994.

THOMAS M. THOMAS II has recently founded his own company, NetCerts (www.netcerts.com), to assist network engineers working toward their Cisco Certifications. Before starting NetCerts, Tom worked as a Course Developer at Cisco Systems for the Worldwide Training division. He worked as part of a team on a new course on Multilayer Switching. He also wrote the book *OSPF Network Design Solutions* for Cisco Press. Tom has also worked as a senior network engineer and group leader of the Advanced Systems Solutions Engineering Team for MCI's Managed Network Services. In this capacity, he developed network maintenance Standard Operating Procedures and performed various in-house training duties. Before joining MCI, Tom worked as a technical team leader at AT&T Solutions, where he provided technical support and network management for Cisco routers over ATM and Frame Relay and configured various networking protocols. E-mail: tthomas@netcerts.com

Letter to the Editor

The article “Was the Melissa Virus so Different?” (*The Internet Protocol Journal*, Volume 2, Number 2, June 1999) by Barbara Y. Fraser et al makes an interesting comparison between events in our real and virtual lives, comparing e-mail borne viruses with commercial samples delivered to our physical mail boxes. While I think the comparison is a useful exercise, the authors fail to point out one of the fundamental differences between these two worlds.

An electronic message contains a finite amount of information: a careful sender can make sure his identity cannot be revealed. In contrast, a physical “message” (i.e., mail bomb, extortion letter, etc.) contains an essentially unlimited amount of information: from finger prints and material analysis to DNA traces, a potential perpetrator can never be certain that he can deny his involvement. For cyberspace crimes the chance to be caught is (and is perceived to be) much smaller. As a result, many virus authors have but the slimmest motive for their deed.

The fact that the Melissa author was quickly identified because of a hidden signature in Microsoft Word is little comfort. For reasons of privacy, this feature has been disabled: it was a bug, not a feature.

To extend the analogy: suppose a simple device would become available that can look up a person’s full ID based on a DNA trace (a few molecules) on any object touched or handled. Move the scanner over the door handle and you know who’s been visiting. The ramifications would be extensive. Most likely, the as-yet hypothetical device would be illegal except for police use.

—Ernst Lopes Cardozo, Aranea Consult BV
e.lopes.cardozo@aranea.nl

Send us your comments!

We look forward to hearing your comments and suggestions regarding anything you read in this publication. Send us e-mail at: ipj@cisco.com

Changes at the IPJ Web Site

Now you can find every issue of *The Internet Protocol Journal* in both PDF and HTML format at www.cisco.com/ipj. We are also pleased to announce that Nikkei Business Publications in Tokyo has provided an introduction to IPJ in Japanese, as well as translation of some of the titles from previous issues at: <http://nit.nikkei.co.jp/ipj.html>. We hope to set up similar links with other publications around the world.

Book Reviews

DHCP *DHCP—A Guide to Dynamic TCP/IP Network Configuration*, by Berry Kercheval, ISBN 0-13-099721-8, Prentice Hall PTR, 1998, http://www.prenhall.com/ptrbooks/ptr_0130997218.html

First, I should note that this book arrived at the perfect time for me: I am involved in adding *Dynamic Host Configuration Protocol* (DHCP) support to a software product and needed a quick, thorough understanding of DHCP that went into sufficient detail to support some key design decisions. The book provided me with exactly what I wanted. However, as to whether or not this is a book you should own or even want to read, that is a much more difficult question to answer.

Organization

The author begins with a chapter of general background information. Then, in a logical progression, he goes through an overview of DHCP and on to explicit details of both the client and server aspects of the protocol. In other sections he covers server administration, DHCP and IP Version 6 (IPv6), and the future of DHCP. He then briefly reviews a few available implementations. In supporting sections he covers the relationship between DHCP and the *Domain Name System* (DNS), specifically Dynamic DNS. In one chapter he discusses the relationship between directory services and DHCP, in particular, the *Lightweight Directory Access Protocol* (LDAP). He then concludes with three appendices: one lists DHCP vendors, another covers the available DHCP options, and a final appendix provides the DHCP RFCs, RFC 2131 and RFC 2132.

Presentation

Overall, the book is well planned and easy to read. The background information is clearly written and gives sufficient material to assure that even novice readers will not get left behind. The author clearly explains the origins of DHCP in BOOTP and the continuing relationship between the two protocols. He also provides many examples that help make the more difficult aspects of DHCP easier to grasp. The chapters tend to progress in a logical order, making absorption of the fairly technical subject almost easy.

The presentation, however, is somewhat marred by minor errors and omissions. None of these mistakes would confuse an expert, but they will make it harder for the novice to be sure what he or she is to understand. In one example, a client workstation on net 10.0.1.0 is offered, and selects, an address of 10.0.2.32. This scenario is, however, clearly unroutable, and the example only confuses the reader. The author also makes a good effort at defining terms the first time they are used, and then again in an extensive glossary. However, for some reason he never defines two key terms: *broadcast* and *multicast*. Since both techniques are core to understanding DHCP, this oversight is difficult to understand.

The chapters on DHCP are fairly exhaustive in their examination of the protocol from overview to minutiae. The roles of clients, servers, and relay agents are well described and documented with sample packets. Each packet field is thoroughly explained and easy to grasp. However, the sections of LDAP and Dynamic DNS could have been presented better. The reader is left with a glimpse of possible relationships between the protocols, but without enough information to really pull it all together. Notably missing is any mention of remote access and the *Remote Authentication Dial-In User Service* (RADIUS) protocol. DHCP and RADIUS perform similar functions in different situations, and there has been much discussion in the past year or two about use of DHCP to manage RADIUS IP address assignments.

Summary

This book sets out to accomplish a limited goal: informing the reader about the basics of DHCP. A couple of detours along the way provide useful information about related technologies (such as DNS and LDAP). The author makes no assumptions about the user's technical capability and level of knowledge. This is perhaps the book's major strength and its biggest weakness. Because of his assumptions about the reader's technical ability, a lot of space is devoted to giving background and reference information assuring that the reader has the necessary foundation to understand the more complex aspects of DHCP. If the background information and appendices (all of which are available on the net and consist mostly of the RFCs anyway) are removed from the book, little is left: without the appendices there are only 144 pages. Given that the book costs \$45, and that the 144 pages are essentially a guided explanation of the RFCs anyway, the technically competent reader might do just as well to download the RFCs and slog through them.

However, for the non-technical reader, or someone who just wants it all in one convenient volume, the author's approach is well worth the cost of the book and the (short) time required to read it. Explanations are clear and concise, terms are well defined, and everything the reader needs to grasp about the complexities of DHCP is right there, in a logical order.

—Richard Perlman, Lucent Technologies
perl@lucent.com

Information Warfare *Information Warfare and Security*, Dorothy E. Denning, ISBN 0-201-43303-6, Addison-Wesley, 1999, <http://www.awl.com/cseng/0-201-43303-6/>

It has been said that “information is power,” and they who control the information control the power. Whether the information is broadcast on the evening news, printed in a newspaper, etched on stone tablets, or published on a USENET newsgroup or Internet Web page, we rely on information in our daily lives, and trust that most of the information we receive and process is accurate.

“Information warfare.” What images does it conjure up for you? Propaganda wars via pamphlets dropped from airplanes, or “cyber-terrorists” versus the FBI on the Internet—or something else entirely? Dr. Denning covers all bases in this, her latest book. The “warfare” of the title is specifically the battle between the good guys and “information terrorists.”

This book is a textbook for a course by the same name at Georgetown University. No one, however, should be scared off by this knowledge. This book is incredibly approachable, intended for a broad audience. It is an introduction to information warfare, but really concentrates on computer- and network-based information. Anyone involved or interested in computer and network security would benefit from this book. Many sections are self-contained, so a reader can jump back and forth among the sections. All the sections are interesting and informative, and should be to both the highly technical reader as well as those for whom technology is peripheral to their jobs, but who require or desire deeper and broader knowledge of information warfare.

About the Author

Dorothy E. Denning is Professor of Computer Science at Georgetown University. She is a well-known expert in the areas of computer security and cryptography, and has been called as an expert witness to testify before the U.S. Congress. She is the author of over 100 papers on computer and Internet security, and has written three other books in addition to this one: *Cryptography and Data Security* (a coeditor with Peter Denning), *Rights and Responsibilities of Participants in Networked Communities*, with Herbert S. Lin, and *Internet Besieged: Countering Cyberspace Scofflaws*. She is also a frequent contributor to security-related publications.

Organization

Information Warfare and Security has three parts. Part 1 starts with a very exciting (and still timely) discussion of the role information warfare played in the Gulf War in the early 1990s. The tone and flavor of this opening chapter continues throughout the book. Randomly put your finger in the book and you will be able to start an enjoyable and interesting read (though I recommend reading beginning to end). Part 1 introduces basic concepts upon which the work is built. Chapters 2 and 3 present a taxonomy of information warfare, relating it to information security and assurance, and suggesting four arenas of activity: play, crime, individual rights, and national security. The author discusses goals, motivations, culture, and concerns. Included is the no-doubt apocryphal, but always fun, quote attributed to Secretary of State for War Henry Stimson, upon the 1929 “discovery” of the Black Chamber code-breaking operation: “Gentlemen, do not read one another’s mail.”

Part 2 focuses on offense. This section covers topics that, for the most part, will be new to many readers. The chapters cover open source material and privacy (and piracy of information), “social engineering,” and its kin. The threat from insiders—legitimate and those who have broken in, gets a thorough treatment. Eavesdropping also is examined, from cellular and pager intercepts, to the mysterious-to-most-people area of traffic analysis, to surveillance, packet-sniffing, and other electronic eavesdropping attacks.

Chapter 8 looks in detail at well-known computer hacking techniques and the tools that implement the attacks. Chapter 9 discusses identity theft, including forged e-mail and stolen accounts, IP-spoofing (stealing the identity of a computer), and Trojan Horse attacks. Finally, Part 2 ends with a chapter dedicated to computer viruses, both real and hoaxes.

Topics discussed in Part 3, “Defensive Information Warfare,” will be familiar to most readers who understand computer and network security. Chapter 11 not only describes cryptographic techniques for protecting information, but also covers *steganography*, or “the practice of hiding a message in such a manner that its very existence is concealed”—and anonymity. Chapter 12, “How to Tell a Fake,” deals with methods for determining identity or trustworthiness of entities or information. Chapter 13 talks about access control mechanisms, including firewalls, and intrusion detection. Covering vulnerability monitoring and analysis, risk analysis, risk management, and incident response, Chapter 14 possibly should have started Part 3. Devices, mechanisms, and methods should be deployed after an understanding of what is contained in this chapter. Part 3, and the book, end with a chapter dedicated to discussing the role of government in defensive information warfare. Also included are descriptions of recent (1990s) actions, laws, and initiatives of the U.S. Government in this area.

Throughout, the book is seasoned with stories—infowar stories, if you will—and background information, allowing the novice not only to understand, but also to enjoy learning what is contained within.

A Book for the Lecture Hall or Armchair

It is not surprising that *Information Warfare and Security* so thoroughly covers the space of information warfare theory, measures, and countermeasures, not because it weighs in at over 500 pages, but because it was written as a text for a course that had to cover all of this material. What may be surprising to readers unfamiliar with Dr. Denning is that such complete coverage could be done in such an easy-to-read way. I have no doubt that this book is and will continue to be useful and effective in the classroom. In addition, the reader studying for accreditation in a field requiring this knowledge, or the professional wanting to “brush up,” “fill in,” or just “kick back,” will find much here to commend itself.

—Frederick M. Avolio, Avolio Consulting
fred@avolio.com

Cryptonomicon *Cryptonomicon*, Neal Stephenson, ISBN 0-380-97346-4, Avon Books, 1999. <http://www.cryptonomicon.com/main.html>

It isn't often that you find reviews of works of fiction in these pages, but *Cryptonomicon* deserves special treatment. Neal Stephenson's latest work is a 918-page science fiction World War II thriller that I couldn't put down. You have to love a novel that has plot points that depend on the technical details of prime number theory, Pretty Good Privacy (PGP), public key infrastructure (PKI), Secure Shell (SSH), Global Positioning System (GPS), secure e-mail, and other Internet applications. Truly this is an epic novel of techno-epic proportions.

The story takes places during both World War II and modern times. The contemporary action revolves around an offshore data haven created by a Silicon Valley startup with the usual coterie of managers, venture capitalists, lawyers with class-action suits, marketeers, and nerds that you'll easily recognize. These entrepreneurs think nothing of flying across the Pacific to attend a meeting and then flying home to get in some quality family time.

The war setting revolves around a small group of code crackers who travel around the globe planting misinformation behind German and Japanese lines. The two groups are literally related: the modern generation is the progeny of the wartime crackers. Both groups are going after hidden caches of gold, among other things, buried near the Philippines.

Technology

There is much technology here for any self-respecting computer geek to digest. Think of Tom Clancy playing with the latest laptops and the Internet rather than with the latest guns. There is even an appendix describing the technical details of one of the crypto algorithms using synchronized decks of playing cards (a key plot point in the book). Stephenson blends in descriptions of undersea cable laying and salvage operations with the cracking of the *Enigma*^[1] codes and hunting down German submarines. At one point, the code-cracking wartime division has to change its numerical designation because it can be factored into two prime numbers—too obvious.

One of my favorite scenes happens early in the book, when the modern-day principals of the crypto firm are meeting some of their backers and potential clients for the first time. The firm's engineer (using the built-in pinhole camera of the laptop) programs his UNIX laptop to surreptitiously capture a photo of whoever is using the keyboard during a demo of the firm's crypto technology, but hides his program in a way that any UNIX hacker would appreciate. He then e-mails the collected digital photos to a friend to try to confirm their identity.

Balance

Unlike Clancy, this book has characters with some depth to them and doesn't overdo the technology. The relationship of the war and modern-day periods is nicely tied together in the end, and the familiarity of the modern-day business relationships is sometimes almost too painful to read.

—David Strom, publisher of *Web Informant*
david@strom.com

References

[1] See <http://www.nsa.gov:8080/museum/enigma.html>

Would You Like to Review a Book for IPJ?

We receive numerous books on computer networking from all the major publishers. If you've got a specific book you are interested in reviewing, please contact us and we will make sure a copy is mailed to you. The book is yours to keep if you send us a review. We accept reviews of new titles, as well as some of the "networking classics." Contact us at ipj@cisco.com for more information.

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

More ICANN News

The *Internet Corporation for Assigned Names and Numbers* (ICANN) recently announced that seven additional applicant companies have met its registrar accreditation criteria.

As accredited registrars, these seven companies will compete in the market for domain name registration services in the **.com**, **.net**, and **.org** domains. In addition, they will be able to participate the ongoing testbed program for the *Shared Registry System*, which allows multiple ICANN-accredited registrars to provide domain name registration services in these domains. Under an agreement announced August 6 by the U.S. Department of Commerce and Network Solutions, Inc. (NSI—the developer of the Shared Registry System), new registrars that have signed an accreditation agreement with ICANN will be eligible to join the initial five testbed registrars as participants in the testbed operation. The testbed phase is currently scheduled to conclude on September 10, 1999.

The seven new companies join the 57 companies that have already been accredited by ICANN starting in April, 1999. Until the initial introduction of competition in June, registration services in the **.com**, **.net**, and **.org** domains were provided solely by NSI under a 1992 Cooperative Agreement with the U.S. Government.

The additional seven companies named are: CommuniTech.Net, Inc. (United States), GANDI (France), iDirections, Inc. (United States), InterNeXt (France), ProBoard Technologies (United States), PSI-USA (United States), and Signature Domains, Inc. (United States). Further information about these companies will be made available on the ICANN Web site:

<http://www.icann.org/registrars/accreditation.html>

Under an October 6, 1998 amendment to the Cooperative Agreement between NSI and the U.S. Government, the process of opening the Internet Domain Name System's three largest domains to competition was launched with a testbed phase that began on April 26. Five companies were initially accredited to use the NSI Shared Registry System in a test operation designed to ensure that the introduction of competition occurs in a smooth, coordinated manner.

By qualifying to be accredited as registrars, the seven new registrars join the five original testbed registrars, as well as the 52 other companies that have already qualified for ICANN accreditation. The Shared Registry System testbed program has been expanded to extend to all accredited registrars that sign the standard testbed registrar agreements with NSI and meet technical certification requirements.

ICANN is a non-profit, international corporation formed in September 1998 to oversee a select set of Internet technical management functions currently managed by the U.S. Government, or by its contractors and volunteers. Specifically, ICANN is assuming responsibility for coordinating the management of the *Domain Name System* (DNS), the allocation of IP address space, the assignment of protocol parameters, and the management of the root server system. For more information, see <http://www.icann.org>. Here you will also find information about ICANN's upcoming public meetings.

INET 2000

INET 2000: The Internet Global Summit, is a special INET. Hosted by the Internet Society, the Summit will be held 18–21 July 2000, in Yokohama, Japan. The place, the date, and the fact that it is the 10th anniversary of this important event all mark it as an exceptional year.

To be considered as a speaker, panelist, tutorial instructor, or poster presenter, please see <http://www.isoc.org/inet2000/callforabstracts.shtml> for submission instructions and to read about this year's theme, "Global Distributed Knowledge for Everyone."

INET is the premier international event for Internet and internetworking professionals. Nowhere can such a broad cross-section of important movers of the Internet be found in one single location.

We look forward to receiving your abstract and seeing you in Japan!

—Jean-Claude Guedon and Jun Murai
Co-Chairs, INET 2000 Program Committee

Y2K and The Internet

As the countdown to the Year 2000 continues, a number of efforts are underway to ensure that the Internet continues to operate normally on January 1, 2000. Here we include some pointers to recent activities.

On July 30, 1999, the *President's Council on Year 2000 Conversion*, convened a roundtable meeting to examine the readiness of the Internet for the Year 2000 date change, and to coordinate efforts to maintain Internet performance and reliability during the transition to the new millennium. The roundtable brought together roughly 100 prominent organizations and individuals from different parts of the Internet community to discuss the Internet's Y2K readiness. Meeting participants included small and large ISPs, equipment vendors, root name server and domain registries, exchange points, network time servers, industry associations, and government officials. For more information see:

<http://www.y2k.gov/> and <http://www.mids.org/y2k/>

For small- and medium-sized businesses in the U.S. and in key trading partner countries, the U.S. Department of Commerce (DoC) is providing a strategic management tool to help battle the millennium bug. The *Y2K Self-Help Tool/CD-ROM* contains a software program that enables users to complete an inventory of assets that may be susceptible to Y2K problems, gauge the criticality of business processes, develop contingency plans and conduct remediation activities.

This CD-ROM contains a 10-minute discussion video, the software program for managing your Y2K process, a self-assessment checklist, contingency planning template, user guide and hotlinks to many helpful Y2K sites. It has been produced in several languages including English, Spanish, Mandarin Chinese, Japanese, French, Portuguese, Arabic and Russian. The software was developed by the DoC's National Institute of Standards and Technology Manufacturing Extension Partnership (MEP) in cooperation with the U.S. Department of Agriculture and the U.S. Small Business Administration.

To receive just the software, visit: www.nist.gov/y2k/software.htm and download *Conversion 2000: Y2K Jumpstart Kit*. To receive the complete CD-ROM with video and hotlinks, you can call 1-800-Y2K-7557 and ask for the Self-Help Tool in any of the languages listed above. If you are an association or organization interested in multiple copies of the CD-ROM for your members and staff, click on order form, print the form, complete the requested information, and fax it to 202-482-0077. Please note that there is a minimal charge for orders over 100 copies for duplication and shipping.

The *Internet Engineering Task Force* (IETF) has examined all of the protocol standards and related documents to identify any potential inherent Y2K problems in the Internet Protocol Suite. The resulting report, RFC 2626, "The Internet and the Millennium Problem (Year 2000)" can be found at <http://www.ietf.org/rfc/rfc2626.txt>

See also:

<http://www.apia.org>

<http://www.nety2k.org/>

<http://www.cert.org/y2k/indmessage.html>

<http://www.icann.org/committees/dns-root/y2k-statement.htm>

This publication is distributed on an "as-is" basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Engineering
MCI WorldCom, USA

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Edward R. Kozel, Sr. VP, Corporate Development
Cisco Systems, Inc., USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the Cisco News
Publications Group, Cisco Systems, Inc.
www.cisco.com*

*Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 1999 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-10/5
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

Bulk Rate Mail U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol Journal

December 1999

Volume 2, Number 4

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Internet Multicast Today	2
The Internet2 Project	20
One Byte at a Time	30
Book Review	33
Call for Papers	35
Fragments	36

FROM THE EDITOR

In June 1992 when I was editor and publisher of *ConneXions—The Interoperability Report*, we published an article entitled “First IETF Internet Audiocast.” Steve Casner and Steve Deering wrote: “The March Internet Engineering Task Force (IETF) meeting in San Diego was an exciting one for those interested in teleconferencing. In addition to several sessions on teleconferencing topics, we managed to pull off a ‘wild idea’ suggested by Allison Mankin from MITRE: live audio from the IETF site was ‘audiocast’ using IP multicast packet audio over the Internet to participants at 20 sites on three continents spanning 16 timezones.”

Multicast has come a long way since 1992. Today, every IETF meeting features several live streams of not only audio but also video and slide presentations. Multicast continues to be developed in the IETF, as protocols and tools are being revised and refined. In two articles, Jon Crowcroft and Mark Handley describe the technologies behind multicast. The first article, included in this issue, looks at the current state of multicast. The second article, to appear in a future issue of IPJ, will look at the problems that need to be solved before multicast can become a truly scalable service for the Internet.

Research into new, high-speed networking technologies and applications is taking place in many parts of the world. One example of such a research effort can be found in the Internet2 Project. Larry Dunn describes some of the technology and application development being conducted by Internet2 members.

Interest in *IP Version 6* (IPv6) is growing as organizations contemplate a world where millions of devices such as cellphones, PDAs, cable TV set-top boxes and so on are “Internet Ready.” The formation of the *IPv6 Forum* (www.ipv6forum.com) is some indication of this interest. We will look at a particular IPv4-to-IPv6 transition strategy in our next issue. In the meantime, Peter Salus takes a historical look at Internet addressing in our series “One Byte at a Time.”

And so we reach the end of 1999 and the end of Volume 2 of *The Internet Protocol Journal*. We wish you a pleasant holiday season and an uneventful transition to Y2K.

—Ole J. Jacobsen, Editor and Publisher

ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Internet Multicast Today

by Mark Handley, ACIRI and Jon Crowcroft, University College London

When you need to send data to many receivers simultaneously, you have two options: repeated transmission and broadcast. Repeated transmission may be acceptable if the cost is low enough and delivery can be spread out over time, as with junk mail or electronic mailing lists. Otherwise, a broadcast solution is required. With real-time multimedia, repeated delivery is feasible, but only at great expense to the sender, who must invest in large amounts of bandwidth. Similarly, traditional broadcast channels have been very expensive if they cover significant numbers of recipients or large geographic areas. However, the Internet offers an alternative solution: IP multicast effectively turns the Internet into a broadcast channel, but one that anyone can send to without having to spend huge amounts of money on transmitters and government licenses. It provides efficient, timely, and global many-to-many distribution of data, and as such may become the broadcast medium of choice in the future.

The Internet is a datagram network, meaning that anyone can send a packet to a destination without having to preestablish a path. Of course, the boxes along the way must have either precomputed a set of paths, or they must be relatively fast at calculating one as needed, and typically, the former approach is used. However, the sending host need not be aware of or participate in the complex route calculation; nor does it need to take part in a complex *signaling* or *call setup* protocol. It simply addresses the packet to the right place, and sends it. This procedure may be a more complex procedure if the sending or receiving systems need more than the default performance that a path or network might offer, but it is the *default* model.

Adding multicast to the Internet does not alter the basic model. A sending host can still simply send, but now there is a new form of address, the multicast or host group address. Unlike unicast addresses, hosts can dynamically subscribe to multicast addresses and by so doing cause multicast traffic to be delivered to them. Thus the IP multicast *service model* can be summarized:

- Senders send to a multicast address
- Receivers express an interest in a multicast address
- Routers conspire to deliver traffic from the senders to the receivers

Sending multicast traffic is no different from sending unicast traffic except that the destination address is slightly special. However, to receive multicast traffic, an interested host must tell its local router that it is interested in a particular multicast group address; the host accomplishes this task by using the *Internet Group Management Protocol* (IGMP).

Point-to-multipoint communication is nothing new. We are all used to the idea of broadcast TV and radio, where a shared medium (the radio frequency [RF] spectrum) is partitioned among users (transmitter or TV/radio station owners). It is a matter of regulation that there is typically only one unique sender of particular content on any given frequency, although other parts of the RF spectrum are given over to free use for multiparty communication (police radio, citizen band radio, and so on).

The Internet multicast *model*^[3] is very similar. The idea is to convert the mesh wide-area network that is the Internet (whether the public Internet, a private enterprise net, or intranet makes no difference to the model), into a shared resource for senders to send to multiple participants, or groups.

To make this group communication work for large-scale systems—in the sense of a large number of recipients for a particular group, or in the sense of a large number of senders to a large number of recipients, or in the sense of a large number of different groups—it is necessary, both for senders and for the routing functions to support delivery, to have a system that can be largely independent of the particular recipients at any one time. In other words, just as a TV or radio station does *not know* who is listening when, an Internet multicast sender does not know who might receive packets it sends. If this scenario sends out alarm bells about security, it shouldn't. A unicast sender has no assurance about who receives its packets either. Assurances about disclosure (privacy) and authenticity of sender/recipient are largely separate matters from simple packet delivery models. Security is a topic of much research and the focus for the recently formed *Internet Research Task Force* (IRTF) research group, *Secure Multicast Group* (SMuG).

The Internet multicast model is an extension of the datagram model; it uses the fact that the datagram is a self-contained communications unit that not only conveys data from source to destination, but also conveys the source and destination address information. In other words, in some senses, datagrams *signal* their own path, both with a source and a destination address in every packet.

By adding a range of addresses dedicated for sending to groups, and providing independence between the address allocation and the rights to send to a group, the analogy between RF spectrum and the Internet multicast space is maintained. Some mechanism, as yet unspecified, is used to dynamically choose which address to send to. Suffice it to say that for now, the idea is that somehow, elsewhere, the address used for a multicast session or group communication activity is chosen so that it does not clash with other uses or users, and is advertised to potential senders and receivers.

Unlike the RF spectrum, an IP packet to be multicast carries a unique source identifier, in that such packets are sent with the normal unicast IP address of the interface of the sending host.

It is also worth noting that an address that is being used to signify a group of entities must surely be a logical address (or in some senses a name) rather than a topological or topographical identifier. We shall see that this means there must be some service that maps such a logical identifier to a specific set of locations in the same way that a local unicast address must be mapped (or bound) to a specific location. In the multicast case, this mapping is distributed. Note also that multicast Internet addresses are in some sense “host group” addresses, in that they indicate a set of hosts to deliver to. In the Internet model, there is a further level of multiplexing, that of transport-level ports, and there is room for some overlap of functionality, since a host may receive packets sent to multiple multicast addresses on the same port, or multiple ports on the same multicast address.

This model raises numerous questions about address and group management, such as how these addresses are allocated. The area requiring most change, though, is in the domain of the routing. Somehow the routers must be able to build a distribution tree from the senders to all the receivers for each multicast group. The senders don’t know who the receivers are (they just send their data), and the receivers don’t know who the senders are (they just ask for traffic destined for the group address), so the routers have to do something without help from the hosts. We will examine this scenario in detail in the section “Multicast Routing.”

Roadmap

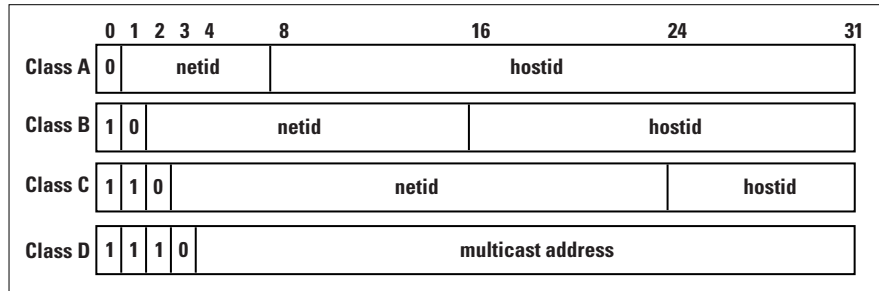
The functions that provide the Standard Internet Multicast Service can be separated into host and network components. The interface between these components is provided by IP multicast addressing and IGMP group membership functions, as well as standard IP packet transmission and reception. The network functions are principally concerned with multicast routing, while host functions also include higher-layer tasks such as the addition of reliability facilities in a transport-layer protocol. That’s the order in which we cover each of these functions in the rest of this article. At the end of the article we list the current status of *Internet Engineering Task Force* (IETF) specification for the various components.

Host Functions

As we stated above, host functionality is extended through the use of the IGMP protocol. Hosts and routers, which we will look at later, must be able to deal with new forms of addresses. When IP Version 4 addressing was first designed, it was divided into classes as shown in Figure 1.

Figure 1: Internet Address Classes

A	1.0.0.0	to 126.255.255.255
B	128.0.0.0	to 191.255.255.255
C	192.0.0.0	to 223.255.255.255
D	224.0.0.0	to 239.255.255.255



Originally Class A was intended for large networks, B for midsize networks, and C for small networks. Class D was later allocated for multicast addresses. Since then, classless addressing has been introduced to solve Internet scaling problems, and the rules for Classes A, B, and C no longer hold, but Class D is still reserved for multicast, so all IPv4 multicast addresses start with the high-order 4-bit “nibble”: 1110

In other words, from the 2^{32} possible addresses, 2^{28} are multicast, meaning that there can be up to about 270 million different groups, each with as many senders as can get unicast addresses! This number is many orders of magnitude more than the RF spectrum allows for typical analog frequency allocations.

For a host to support multicast, the host service interface to IP must be extended in three ways:

- A host must be able to join a group, meaning that it must be able to reprogram its network level, and possibly, consequentially, the lower levels, to be able to receive packets addressed to multicast group addresses.
- An application that has joined a multicast group and then sends to that group must be able to select whether it wants the host to loop-back the packets it sent so that it receives its own packets.
- A host should be able to limit the *scope* with which multicast messages are sent. The Internet Protocol contains a *Time-To-Live* (TTL) field, used originally to limit the lifetime of packets on the network, both for safety of upper layers, and for prevention of traffic overload during temporary routing loops. It is used in multicast to limit how “far” a packet can go from the source. We will see below how scoping can interact with routing.

When an application tells the host networking software to join a group, the host software checks to see if the host is a member of the group. If not, it makes a note of the fact, and sends out an IGMP membership report message. It also maps the IP address to a lower-level address and reprograms its network interface to accept packets sent to that address. There is a refinement here: a host can join “on an interface;” that is, hosts that have more than one network card can decide which one (or more than one) they wish to receive multicast packets via. The implication of the multicast model is that it is “pervasive,” so it is usually necessary to join on only one interface.

Taking a particular example to illustrate the IP-level to link-level mapping process, if a host joins an IP multicast group using an Ethernet interface, there is a mapping from the low 24 bits of the multicast address into the low 24 (out of 48) bits of the Ethernet address. Since this mapping is a many-to-one mapping, there may be multiple IP multicast groups occupying the same Ethernet address on a given wire, though it may be made unlikely by the address allocation scheme. An Ethernet LAN is a shared-medium network, thus local addressing of packets to an Ethernet group means that the packets are received by Ethernet hardware and delivered to the host software of *only* those hosts with members of the relevant IP group. Therefore, host software is generally saved the burden of filtering out irrelevant packets. Where there is an Ethernet address clash, software can filter the packets efficiently.

Operation of the IGMP protocol can be summarized as follows:

- When a host first joins a group, it programs its Ethernet interface to accept the relevant traffic, and it sends an IGMP Join message on its local network. This message informs any local routers that there is a receiver for this group now on this subnet.
- The local routers remember this information, and arrange for traffic destined for this address to be delivered to the subnet.
- After a while, the routers wonder if there is still any member on the subnet, and send an IGMP query message to the multicast group. If the host is still a member, it replies with a new message unless it hears someone else do so first. Multicast traffic continues to be delivered.
- Eventually the application finishes, and the host no longer wants the traffic. It reprograms its Ethernet interface to reject the traffic, but the packets are still sent until the router times the group out and sends a query to which no one responds. The router then stops delivering the traffic.

Thus joining a multicast group is quick, but leaving can be slow with IGMP Version 1. IGMP Version 2 reduces the leave latency by introducing a “Leave” message and a set of rules to prevent one receiver from disconnecting others when it leaves. IGMP Version 3 (not yet deployed) introduces the idea of *source-specific* joining and leaving, whereby a host can subscribe (or reject) traffic from individual senders rather than the group as a whole, at the expense of more complexity and extra state in routers.

Multicast Routing

Given the multicast service model described above, and the restrictions that senders and receivers don’t know each others’ location or anything about the topology, how do routers conspire to deliver traffic from the senders to the receivers?

We shall assume that if a sender and a receiver did know about each other, they could each send unicast packets to the other. In other words, there is a network with bidirectional paths and an underlying unicast routing mechanism already running. Given this network, there is a spectrum of possible solutions. At one extreme, we can flood data from the sender to all possible receivers and have the routers for networks where there are no receivers prune off their branches of the distribution tree. At the other extreme, we can communicate information in a multicast routing protocol conveying the location of all the receivers to the routers on the paths to all possible senders. Neither method is particularly desirable on a global scale, so the most interesting solutions tend to be hybrid solutions that lie between these extremes.

In the real world, there are many different multicast routing protocols, each with its own advantages and disadvantages. We shall explain each of the common ones briefly, because a working knowledge of their pros and cons helps us understand the practical limits to the uses of multicast.

Flood and Prune Protocols

Flood and Prune Protocols are more correctly known as *reverse-path multicast* algorithms. When a sender first starts sending, traffic is flooded out through the network. A router may receive the traffic along multiple paths on different interfaces, in which case it rejects any packet that arrives on any interface other than the one it would use to send a unicast packet back to the source. It then sends a copy of each packet out of each interface other than the one back to the source. In this way, each link in the whole network is traversed at most once in each direction, and the data is received by all routers in the network.

So far, this process describes *reverse-path broadcast*. Many parts of the network will be receiving traffic, even though there are no receivers there. These routers know they have no receivers (otherwise IGMP would have told them) and they can then send prune messages back toward the source to stop unnecessary traffic from flowing. Thus the delivery tree is pruned back to the minimal tree that reaches all the receivers. The final distribution tree is what would be formed by the union of shortest paths from each receiver to the sender, so this type of distribution tree is known as a *shortest-path tree* (strictly speaking, it's a reverse shortest path tree—typically the routers don't have enough information to build a true forward shortest-path tree).

Two commonly used multicast routing protocols fall in the class: the *Distance Vector Multicast Routing Protocol* (DVMRP)^[4] and *Protocol Independent Multicast Dense-Mode* (PIM-DM)^[5]. The primary difference between these protocols is that DVMRP computes its own routing table to determine the best path back to the source, whereas PIM Dense-Mode uses the routing table of the underlying unicast routing system, hence the term “Protocol Independent.”

It should be fairly obvious that sending traffic *everywhere* and getting people to tell you what they don't want is not a particularly scalable mechanism. Sites get traffic they don't want (albeit very briefly), and routers not on the delivery tree need to store prune state. For example, if a group has one member in the UK and two in France, routers in Australia still get some of the packets, and they need to hold prune state to prevent more packets from arriving! However, for groups where most places actually do have receivers (receivers are “densely” distributed), this sort of protocol works well. So although these protocols are poor choices for a global scheme, they might be appropriate within some organizations.

MOSPF

Multicast Open Shortest Path first (MOSPF^[12]) isn't really a category, but a specific instance of a protocol. MOSPF is the multicast extension to *Open Shortest Path First* (OSPF^[11]), which is a unicast link-state routing protocol.

Link-state routing protocols work by having each router send a routing message periodically listing its neighbors and how far away they are. These routing messages are flooded throughout the entire network, so every router can build up a map of the network. This map is then used to build forwarding tables (using a Dijkstra algorithm) so that the router can decide quickly which is the correct next hop for a particular packet.

Extending this concept to multicast is achieved simply by having each router also list in a routing message the groups for which it has local receivers. Thus given the map and the locations of the receivers, a router can also build a multicast forwarding table for each group.

MOSPF also suffers from poor scaling. With flood-and-prune protocols, data traffic is an *implicit* message about where there are senders, so routers need to store unwanted state where there are no receivers. With MOSPF, there are *explicit* messages about where all the receivers are, so routers need to store unwanted state where there are no senders. However, both types of protocol build very efficient distribution trees.

Center-Based Trees

Rather than flooding the data everywhere, or flooding the membership information everywhere, algorithms in the center-based trees category map the multicast group address to a particular unicast address of a router, and they build explicit distribution trees centered around this particular router. Three main problems need to be solved to get this approach to work:

- How is the mapping from group address to center address performed?
- How is the center location chosen so that the distribution trees are efficient?
- How is the tree actually constructed given the center address?

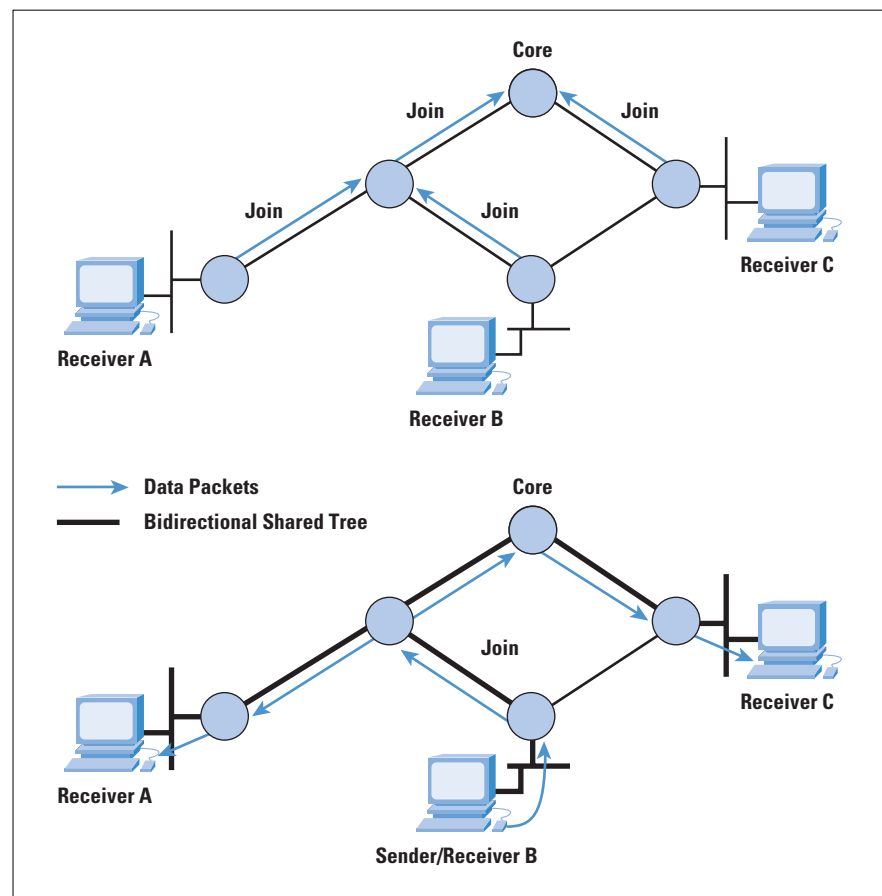
Different protocols have come up with different solutions to these problems. Three center-based tree protocols are worth exploring because they illustrate different approaches: *Core-Based Trees* (CBT), *PIM Sparse-Mode* (PIM-SM), and the *Border Gateway Multicast Protocol* (BGMP). However, we will leave discussion of BGMP until our second article because it is not currently deployed.

Core-Based Trees

Core-Based Trees (CBT^[1]) was the earliest center-based tree protocol, and it is the simplest.

When a receiver joins a multicast group, its local CBT router looks up the multicast address and obtains the address of the Core router for the group. It then sends a Join message for the group toward the Core. At each router on the way to the Core, forwarding state is instantiated for the group, and an acknowledgment is sent back to the previous router. In this way, a multicast tree is built, as shown in Figure 2.

Figure 2: Formation of a CBT Bidirectional Shared Tree



If a sender (that is, a group member) sends data to the group, the packets reach its local router, which forwards them to any of its neighbors that are on the multicast tree. Each router that receives a packet forwards it out of all its interfaces that are on the tree except the one the packet came from. The style of tree CBT builds is called a “bidirectional shared tree,” because the routing state is “bidirectional”—packets can

flow both up the tree toward the Core and down the tree away from the Core, depending on the location of the source, and packets are “shared” by all sources to the group. This scenario is in contrast to “unidirectional shared trees” built by PIM-SM as we shall see later.

IP multicast does not require senders to a group to be members of the group, so it is possible that a sender’s local router is not on the tree. In this case, the packet is forwarded to the next hop toward the Core. Eventually the packet will either reach a router that is on the tree, or it will reach the Core, and it is then distributed along the multicast tree.

CBT also allows multiple Core routers to be specified, adding a little redundancy in case the Core becomes unreachable. CBT never properly solved the problem of how to map a group address to the address of a Core. In addition, good Core placement is a difficult problem. Without good Core placement, CBT trees can be quite inefficient, and so CBT is unlikely to be used as a global multicast routing protocol.

However, within a limited domain, CBT is very efficient in terms of the amount of state that routers need to keep. Only routers on the distribution tree for a group keep forwarding state for that group, and no router needs to keep information about any source; thus CBT scales much better than flood-and-prune protocols, especially for sparse groups where only a small proportion of subnetworks have members.

PIM Sparse-Mode

The work on CBT encouraged others to try to improve on its limitations while keeping the good properties of shared trees, and *PIM Sparse-Mode*^[7] was one result. The equivalent of a CBT Core is called a *Rendezvous Point* (RP) in PIM, but it largely serves the same purpose.

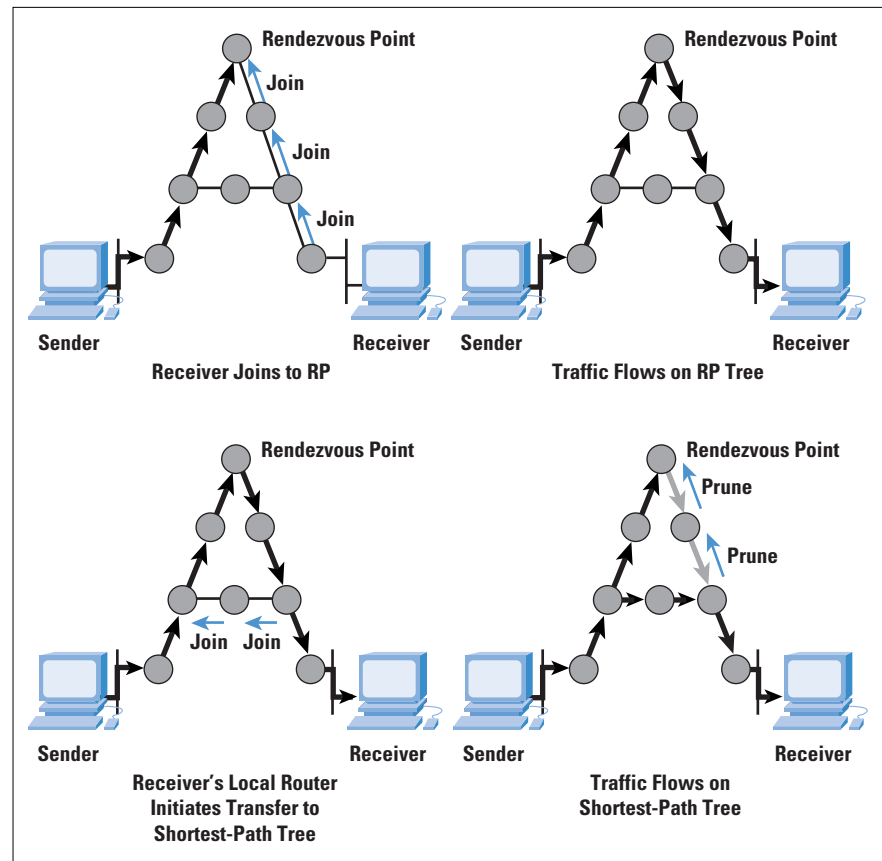
When a sender starts sending, whether it is a member or not, its local router receives the packets and maps the group address to the address of the RP. It then encapsulates each packet in another IP packet (imagine putting one letter inside another, differently addressed, envelope) and sends it unicast directly to the RP.

When a receiver joins the group, its local router initiates a Join message that travels hop-by-hop to the RP instantiating forwarding state for the group. However, this state is unidirectional state—it can be used only by packets flowing from the RP toward the receiver, and not for packets flowing back up the tree toward the RP. Data from senders is de-encapsulated at the RP and flows down the shared tree to all the receivers.

PIM-SM is an improvement on CBT in that discovery of senders and and tree building from senders to receivers are separate functions.

Thus PIM-SM unidirectional trees are not particularly good distribution trees, but they do start data flowing to the receivers. Once this data is flowing, the local router of a receiver can then initiate a transfer from the shared tree to a shortest-path tree by sending a source-specific Join message toward the source, as shown in Figure 3. When data starts to arrive along the shortest-path tree, a prune message can be sent back up the shared tree toward the source to avoid getting the traffic twice.

Figure 3: Formation of a PIM Sparse-Mode Tree



Unlike other shortest-path tree protocols such as DVMRP and PIM-DM, where prune state exists everywhere there are no receivers, with PIM-SM, source-specific state exists only on the shortest-path tree. Also, low-bandwidth sources such as those sending *Real-Time Control Protocol* (RTCP) receiver reports do not trigger the transfer to a shortest-path tree, a scenario that further helps scaling by eliminating unnecessary source-specific state.

Because PIM-SM can optimize its distribution trees after formation, it is less critically dependent on the RP location than CBT is on the Core location. Hence the primary requirement for choosing an RP is load balancing. To perform multicast-group-to-RP mapping, PIM-SM predistributes a list of candidates to be RPs to all routers. When a router needs to perform this mapping, it uses a special hash function to hash the group address into the list of candidate RPs to decide the actual RP to join.

Except in rare failure circumstances, all the routers within the domain will perform the same hash, and come up with the same choice of RP. The RP may or may not be in an optimal location, but this situation is offset by the ability to switch to a shortest-path tree.

The dependence on this hash function and the requirement to achieve convergence on a list of candidate RPs does, however, limit the scaling of PIM-SM. As a result, it is also best deployed within a domain, although the size of such a domain may be quite large.

Interdomain Multicast Routing

All the multicast routing schemes described so far suffer from scaling problems of one form or another:

- DVMRP and PIM-DM initially send data everywhere, and require routers to hold prune state to prevent this flooding from persisting.
- MOSPF requires all routers to know where all receivers are.
- PIM-SM needs predistribution of information about the set of RPs. Because traffic needs to flow to the RP, an RP cannot handle too many groups simultaneously, so many RPs are needed globally.

Thus each of these schemes is likely to be best deployed within a domain. How then does interdomain multicast routing take place?

Long-term solutions to this problem will be discussed in the second of these articles. In the meantime, the interim solution currently being deployed consists of multiprotocol extensions to the unicast *Border Gateway Protocol* (BGP) interdomain routing protocol, and a protocol called MSDP to glue PIM-SM domains together.

Multiprotocol BGP

For either technical or policy reasons, not all routers or peerings between Internet Service Providers (ISPs) are multicast capable. This situation complicates the use of PIM-SM for operation between domains because PIM assumes that the route obtained by unicast routing is good for multicast routing (strictly speaking, PIM assumes the reverse unicast path is good for forward-path multicast routing). If, in fact, the reverse unicast path is *not* good for forward-path multicast, then Join messages will often reach routers that do not support multicast, resulting in a lack of multicast connectivity. How then do we solve this problem?

BGP is the unicast interdomain routing protocol that is very widely used to connect unicast routing domains together. The multiprotocol extensions to BGP allow multiple routing tables to be maintained for different protocols. Thus with the *Multiprotocol Extensions for BGP-4* (MBGP)^[2], you can build one routing table for unicast-capable routes and one for multicast-capable routes using the same protocol. PIM can then use the multicast-capable routes to forward Join messages and can, therefore, detour around parts of the network that don't support multicast.

Multicast Source Discovery Protocol

In addition to the problem of designing a scalable mechanism for mapping multicast groups to RPs, attempts to use PIM-SM as an interdomain protocol are hindered by ISPs' desire not to be dependent on other ISPs' facilities. For example, consider a multicast group consisting of senders and receivers in two domains, A and B, run by two different ISPs. If the RP is in domain A, and there is some problem in domain A, then senders and receivers in domain B might still be unable to communicate with each other using multicast, even though they are in the same domain, because initial PIM register messages must go via the RP. ISPs do not want to be dependent on other ISPs for connectivity within their own domain, so it appears that using PIM-SM as an interdomain protocol would be unacceptable, even if there were no scalability problems.

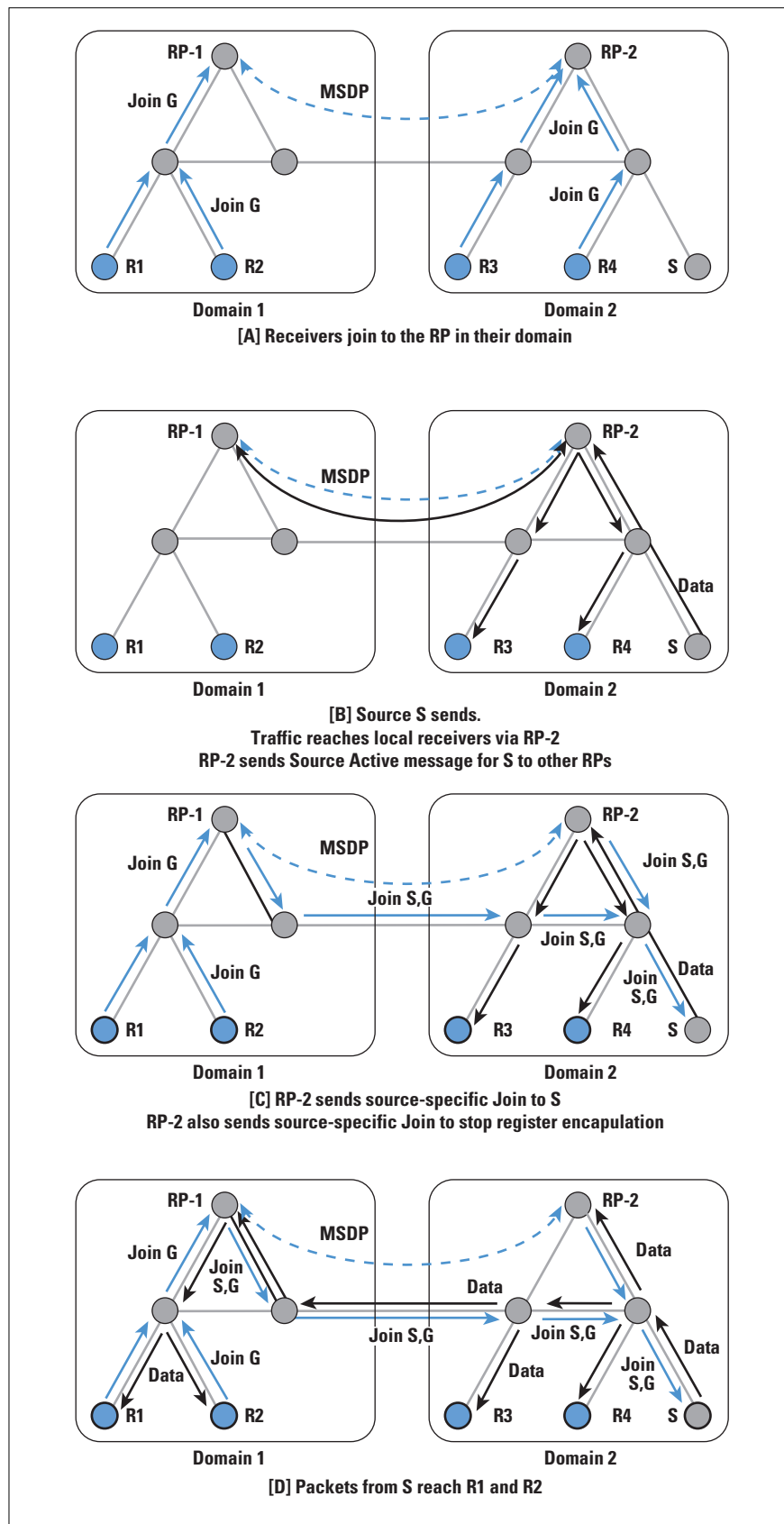
The *Multicast Source Discovery Protocol* (MSDP)^[8] is an attempt to work around this problem. It does not provide a long-term scalable solution, but does provide a solution that solves the ISP interdependence problem.

With MSDP, ISPs run PIM-SM within their own domain, and they have their own set of RPs for all groups within that domain. Additionally, the RPs within the domain are interconnected with each other and with RPs in neighboring domains using MSDP control connections to form a loose mesh.

The process is shown in Figure 4. Within domain 1, R1 and R2 send Join messages from group G to RP-1. Similarly, R3 and R4 send Join messages to RP-2. When S starts sending, its packets are encapsulated to RP-2 by its local router in the normal PIM-SM manner. RP-2 decapsulates the packets and forwards them down the group-shared tree within domain 2 to reach R3 and R4. In addition, it sends a *Source Active* message over the MSDP mesh to all other RPs. RPs like RP-1 that have active joiners for this group then send a source-specific Join back across the interdomain boundary toward S. Traffic is then delivered interdomain following the source-specific state laid down by the Join messages, and it is eventually delivered to R1 and R2.

MSDP uses the normal PIM-SM source-specific join mechanism interdomain following the MBGP multicast routes back to the source, but it sets up only a group-shared tree within each domain, avoiding the need to depend on remote RPs in different domains for the delivery of traffic between local members in a domain.

Figure 4: MSDP in Operation



As an interdomain routing protocol, however, MSDP has many shortcomings. In particular, every RP in every domain must be told about every source that starts sending, and a significant subset of the RPs must cache all this information so that receivers that join late can cause source-specific Joins to be sent by their local RP. Thus MSDP does not scale well if there are a large number of senders worldwide.

In addition, to ensure that the first few packets sent by a source do not get lost, they must be encapsulated and sent alongside the *Source Active* message to all the RPs that might possibly have receivers. If they are not encapsulated, then sources that send only a few packets every few minutes might never get any data through to receivers because the source-specific state has timed out after each time they send.

In summary, MSDP is not a scalable long-term solution to interdomain multicast routing. However, it does solve a real short-term problem faced by ISPs, and so it is currently seeing significant deployment.

Multicast Address Allocation

A local protocol for requesting multicast addresses from multicast address allocation servers has recently been standardized. This protocol is called *Multicast Address Dynamic Client Allocation Protocol*, or MADCAP^[10]. It is a relatively simple request-response protocol loosely modeled after the *Dynamic Host Configuration Protocol* (DHCP)^[6].

MADCAP is intended to be used with interdomain protocols that perform dynamic allocation of parts of the multicast address space between domains, but because these protocols are not yet deployed, they will be discussed in the second of these articles.

As an interim solution for interdomain address allocation, a simple static mechanism has been defined. This mechanism involves embedding the *Autonomous System* (AS) number of the domain as the middle 16 bits of a multicast address. Thus the domain with AS number 16007 would get multicast addresses in the range 233.64.7.0 to 233.64.7.255 (64 and 7 being the upper and lower bytes, respectively, of 16007). Known as *glop addressing*, this mechanism is experimental. It may be superseded by a dynamic mechanism in the longer term.

Multicast Scoping

When applications operate in the global Multicast backbone (MBone), it is clear that not all groups should have global scope. Not only is this constraint especially important for performance reasons with flood and prune multicast routing protocols, but it also is true with other routing protocols for application security reasons and because multicast addresses are a scarce resource. Being able to constrain the scope of a session allows the same multicast address to be in use at more than one place as long as the scopes of the sessions do not overlap. This is analogous to the same radio frequency being used by two radio stations operating far apart from one another—each will only be heard locally.

Multicast scoping can currently be performed in two ways, known as *TTL Scoping* and *Administrative Scoping*. Currently TTL scoping is most widely used, with only a very few sites making use of administrative scoping.

TTL Scoping

When an IP packet is sent, an IP header field called *Time To Live* (TTL) is set to a value between zero and 255. Every time a router forwards the packet, it decrements the TTL field in the packet header, and if the value reaches zero, the packet is dropped. The IP specification also states that the TTL should be decremented if a packet is queued for more than a certain amount of time, but this decrement is rarely implemented these days. With unicast, the TTL is normally set to a fixed value by the sending host (64 and 255 are commonly used) and is intended to prevent packets from looping forever.

With IP multicast, the TTL field can be used to constrain how far a multicast packet can travel across the MBone by carefully choosing the value put into packets as they are sent. However, because the relationship between hop count and suitable scope regions is poor at best, the basic TTL mechanism is supplemented by configured thresholds on multicast tunnels and multicast-capable links. Where such a threshold is configured, the router will decrement the TTL, as with unicast packets, but then will drop the packet if the TTL is less than the configured threshold. When these thresholds are chosen consistently at all of the borders to a region, they allow a host within that region to send traffic with a TTL less than the threshold, and to know that the traffic will not escape that region.

An example is the multicast tunnels and links to and from Europe, which are all configured with a TTL threshold of 64. Any site within Europe that wishes to send traffic that does not escape Europe can send with a TTL of less than 64 and be sure that its traffic does not escape.

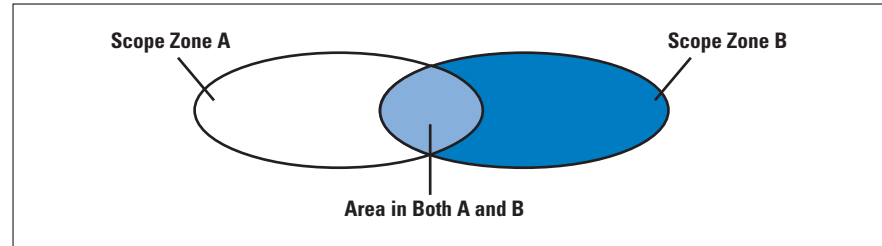
However, there are also likely to be thresholds configured within a particular scope zone—for example, most European countries use a threshold of 48 on international links within Europe, and because TTL is still decremented each time the packet is forwarded, it is good practice to send European traffic with a TTL of 63, a scenario that allows the packet to travel 15 hops before it would fail to cross a European international link.

Administrative Scoping

In some circumstances it is difficult to consistently choose TTL thresholds to perform the desired scoping. In particular, it is impossible to configure overlapping scope regions as shown in Figure 5, and TTL scoping has numerous other problems, so more recently, administrative scoping has been added to the multicast forwarding code in *mrouted* and in most router implementations.

Administrative scoping allows the configuration of a boundary by specifying a range of multicast addresses that will not be forwarded across that boundary in either direction.

Figure 5: Overlapping Scope Zones possible with Administrative Scoping



Scoping Deployment

Administrative scoping is much more flexible than TTL scoping, but it has many disadvantages. In particular, it is not possible to tell from the address of a packet where it will go unless all the scope zones that the sender is within are known. Also, because administrative boundaries are bidirectional, one scope zone nested within or overlapping another must have totally separate address ranges. This makes address allocation difficult from an administrative point of view, because the ranges ought to be allocated on a top-down basis (largest zone first) in a network where there is no appropriate top-level allocation authority. Finally, it is easy to misconfigure a boundary by omitting or incorrectly configuring one of the routers. With TTL scoping it is likely that in many cases a more distant threshold will perform a similar task, lessening the consequences, but with administrative scoping, there is less likelihood that this scenario will occur.

For these reasons, administrative scoping has been viewed by many network administrators as a speciality solution to difficult configuration problems, rather than as a replacement for TTL scoping, and the Mbone still very much relies on TTL scoping. However, this situation is set to change as a protocol for automatically discovering scope zones (and scope zone misconfigurations) starts to be deployed. This protocol is called the *Multicast Zone Announcement Protocol* (MZAP)^[9], and it will shortly become an IETF Proposed Standard. Eventually the use of configured TTL scopes to restrict traffic will cease to be used as a primary scoping mechanism.

Summary

In this article we have looked at the various routing systems that are used to devise delivery trees over which multimedia data can be sent for the purposes of group communication, and at address allocation and scoping mechanisms for this traffic.

After ten years of experimentation, IP multicast is not currently a ubiquitous service on the public Internet, but significant deployment has taken place on private intranets. The existing multicast routing and address allocation mechanisms work well at the scale of domains. However, as we have seen, there are still significant technical problems

concerning scaling to be overcome before multicast can be a ubiquitous interdomain service. In addition to the routing problems, we also still lack deployed congestion control mechanisms for multicast traffic, which are essential if multicast applications are to be safely deployed.

Despite these issues, IP multicast still shows great promise for many applications. Solutions have been devised to many of the remaining problems, although they have not yet been deployed. In the second of these articles, we will look at the proposed solutions for scalable interdomain routing and address allocation. We will also touch on multicast congestion control and the solutions that are currently emerging from the research community.

Document Status

A list of IETF specifications for the protocols discussed in this article is given below. We include the status for each document as of this writing (November 1999). For more information, check the IETF Web pages at www.ietf.org

Document	Status
IGMP v1	IETF Standard (RFC 1112)
IGMP v2	IETF Proposed Standard (RFC 2236)
IGMP v3	IETF work in progress
DVMRP	IETF Experimental Standard (RFC 1075)
PIM-Dense Mode	IETF work in progress
Multicast OSPF	IETF Proposed Standard (RFC 1584)
Core Based Trees	IETF Experimental Standard (RFC 2201)
PIM Sparse-Mode	IETF Experimental Standard (RFC 2362)
Multiprotocol BGP	IETF Proposed Standard (RFC 2283)
MSDP	IETF work in progress
MADCAP	IETF Proposed Standard (RFC 2730)
Glop Addressing	IETF work in progress

References

- [1] Ballardie, A., "Core Based Trees (CBT version 2) Multicast Routing," RFC 2189, September 1997.
- [2] Bates, T., Chandra, R., Katz, D., and Rekhter, Y., "Multiprotocol Extensions for BGP-4," RFC 2283, February 1998.
- [3] Deering, S., "Host Extensions for IP Multicasting," RFC 1112, August 1989.
- [4] Deering, S., Partridge, C., and Waitzman, D., "Distance Vector Multicast Routing Protocol," RFC 1075, November 1988.

- [5] Deering, S., Estrin, D., Farinacci, D., Jacobson, V., Helmy, A., Meyer, D., and Wei, L., "Protocol Independent Multicast Version 2 Dense Mode Specification," Internet Draft, work in progress.
- [6] Droms, R., "Dynamic Host Configuration Protocol," RFC 1531, October 1993.
- [7] Estrin, D., Farinacci, D., Helmy, A., Thaler, D., Deering, S., Handley, M., Jacobson, V., Liu, C., Sharma, P., and Wei, L., "Protocol Independent Multicast-Sparse Mode (PIM-SM): Protocol Specification," RFC 2362, June 1998.
- [8] Farinacci D. et al. "Multicast Source Discovery Protocol (MSDP)," Internet Draft, work in progress, June 1998.
- [9] Handley, M., Thaler, D., and Kermode, R., "Multicast-Scope Zone Announcement Protocol (MZAP)," Internet Draft, work in progress.
- [10] Hanna, S., Patel, M., and Shah, M., "Multicast Address Dynamic Client Allocation Protocol (MADCAP)," RFC 2730, December 1999.
- [11] Moy, J., "OSPF Version 2," RFC 2328, April 1998.
- [12] Moy, J., "Multicast Extensions to OSPF," RFC 1584, March 1994.
- [13] Miller, C. K., "Reliable Multicast Protocols and Applications," *The Internet Protocol Journal*, Volume 1, No. 2, September 1998.

JON CROWCROFT is a professor of networked systems in the Department of Computer Science, University College London, where he is responsible for a number of European and U.S. funded research projects in Multi-media Communications. He has been working in these areas for over 18 years. He graduated in Physics from Trinity College, Cambridge University, in 1979, and gained his MSc in Computing in 1981, and PhD in 1993. He is a member of the ACM, the British Computer Society, and is a Fellow of the IEE and a senior member of the IEEE. He is a member of the Internet Architecture Board (IAB) and was general chair for the ACM SIGCOMM from 1995 to 1999. He is also on the editorial team for the ACM/IEEE *Transactions on Networks*. With Mark Handley, he is the co-author of *WWW: Beneath the Surf* (UCL Press); he also authored *Open Distributed Systems* (UCL Press/Artech House), and with Mark Handley and Ian Wakeman, a third book, *Internetworking Multimedia* (Morgan Kaufmann Publishers), published in October 1999.

E-mail: **J.Crowcroft@cs.ucl.ac.uk**

MARK HANDLEY received his BSc in Computer Science with Electrical Engineering from University College London in 1988 and his PhD from UCL in 1997. For his PhD he studied multicast-based multimedia conferencing systems, and was technical director of the European Union funded MICE and MERCI multimedia conferencing projects. After two years working for the University of Southern California's Information Sciences Institute, he moved to Berkeley to join the new AT&T Center for Internet Research at ICSI (ACIRI). Most of his work is in the areas of scalable multimedia conferencing systems, reliable multicast protocols, multicast routing and address allocation, and network simulation and visualisation. He is co-chair of the IETF Multiparty Multimedia Session Control working group and the IRTF Reliable Multicast Research Group.

E-mail: **mjh@aciri.org**

[This article is based in part on material in *Internetworking Multimedia* by Jon Crowcroft, Mark Handley, and Ian Wakeman, ISBN 1-55860-584-3, published by Morgan Kaufmann in 1999. Used with permission].

The Internet2 Project

by Larry Dunn, Cisco Systems

Communication, connectivity, education, entertainment, e-commerce—across a broad spectrum of activities, the commodity Internet has made a strong impact on the way we live, work, and play. Nevertheless, many classes of applications do not yet run well, and some don't run at all, over the commodity net. As new applications are developed in disciplines from medicine to engineering to the arts and sciences, their success increasingly depends on an ability to use networks effectively. In research and education collaborations all over the world, efforts are under way to make use of new network technologies and develop network services that will facilitate these advanced applications. One such effort in the United States is called the *Internet2 Project*^[1].

The Internet2 Project was started in 1996 by 34 U.S. research universities. It has since grown to over 140 universities, and includes several corporate members and international partners. This article examines network technology used in Internet2, and looks at some of the engineering challenges involved in facilitating applications being developed by Internet2 members.

Background

In 1995, the U.S. National Science Foundation (NSF) funded a program to create the *very-high-performance Backbone Network Service* (vBNS)^[2]. The NSF provided funding to MCI, who interconnected five U.S. supercomputer centers and 3 *Network Access Points* (NAPs), where it was envisioned that supercomputer clients and other vBNS users would connect.

By 1996, congestion stemming from academic traffic to the commodity Internet had seriously congested the NAPs; it was accordingly recognized that clients of the supercomputer centers might be better served if the Research Universities, where Principal Investigators often resided, were themselves *directly* connected to the vBNS. So in 1996, the NSF accepted proposals as part of the *High-Performance Connections* (HPC) program^[3]. Schools applying for an HPC grant might receive \$350,000 over a 2-year period, provided their proposals met various criteria, including meritorious research that would benefit from the high-performance connection, a solid network plan, intention to investigate capabilities enabled by such a connection, commitment to share results with the community, matching funds from the University, and so on.

In October 1996, representatives from 34 universities met, and concluded that, while not all the schools had projects involving “meritorious research” that would meet the NSF criteria, they all *did* have a critical interest in deploying the kind of applications that such high-performance connections could enable.

Thus, to facilitate development and deployment of applications that would further the research and education mission of member universities, the Internet2 Project was formed.

From the beginning, the stated intention was to enable applications that could not run, or could not run well, on the “Commodity Internet.” Networks would be utilized or constructed only so as to facilitate this applications-enabling goal, and results/methods would be applied to the broader community as rapidly as possible.

Applications Focus

The list of applications being used or developed by Internet2 members is extensive. Several fall in the category of “meritorious research” as mentioned in the NSF HPC criteria. Examples include: remote instrument control (for instance, telescopes, microscopes), high-performance distributed computation, and large-scale database navigation. Other applications that further the education mission of member universities include tools to facilitate multisite collaboration, and asynchronous learning. Many examples in areas from science, engineering, art, language, music, and more can be found at the Internet2 applications Web site^[4]. In addition to individual applications, a couple of broad initiatives have a relationship with Internet2, including *The Internet2 Digital Video Initiative*, housed at the *International Center for Advanced Internet Research* (iCAIR)^[5], and the *Internet2 Distributed Storage Infrastructure Initiative* (I2-DSI)^[6].

The above applications share several challenging requirements, many of which translate to resource commitments that must be met by the network in an end-to-end fashion, including bandwidth and jitter. Additionally, the applications can become scalable only if more-mature middleware and control-plane infrastructure is developed. Necessary components include features such as *Authentication, Authorization, and Accounting* (AAA), scheduling, and coordination of resources managed by multiple administrative domains.

One compelling example of the network challenges present in a virtual collaborative environment is exemplified by a CAVE (*Cave Automated VR Environment*). See [7] for more details, but in brief, a single CAVE is a (10 x 10 x 10)-foot cube, with one wall removed. Users enter through the open wall, and using lightweight stereo-three-dimensional (3D) glasses, and a radio frequency (RF) mouse, can interact with an immersive environment created by rear-screen and direct projection on multiple walls and the floor. As an example, the interconnection of multiple CAVEs allows design teams in remote locations to jointly experience the operating “feel” of a new vehicle, and to dynamically adjust, design, or control parameters to see how the modified vehicle behaves.

The developers of CAVE software at Argonne National Labs have noted that the data flows in a CAVE consist of at least: control, text, audio, video, tracking, database, simulation, haptic, and rendering flows. Additionally, they have estimated the latency, jitter, and bandwidth requirements for these flows. Some of the flows represent a challenge in a single resource dimension, others have strict requirements in multiple resource dimensions.

Backbone Networks

At this time, Internet2 members may connect to either of two backbone networks, or both.

The vBNS is operated by MCI/Worldcom. It consists primarily of an IP-over-ATM network. Most schools connect at DS3 or OC-3c via ATM to a vBNS ATM switch. Interior vBNS links are OC-12c ATM. The schools peer with a Layer 3 router; a router is attached to each of the vBNS ATM switches. The vBNS routers are logically connected to each other via a full mesh of *Unspecified Bit Rate* (UBR) *Virtual Circuits* (VCs). The ATM switches are connected to each other via a second layer of ATM switches, which are part of MCI's commercial Hyperstream offering. While schools pass the vast majority of their traffic via peering at Layer 3 with the nearest vBNS border router, other services are available, including the option to establish *Variable Bit Rate* (VBR) VCs as needed, and the possibility to place some of the ATM-attached hosts of the school directly in a vBNS Classical IP *Logical IP Subnet* (LIS). This setup allows such hosts to send bytes directly to other ATM-attached hosts, bypassing the routers of both the school and the vBNS. The vBNS also carries native IP multicast traffic among members. In addition, the vBNS has a native IPv6 offering, which is achieved by deploying routers that run IPv6, and provisioning VCs to schools also running IPv6. The vBNS has also begun to offer an *IP-over-Synchronous Optical Network* (SONET) service, the first instance of which is an OC-48 *Packet-over-SONET* (POS) link from Northern to Southern California. Because the nominal partnership arrangement with the NSF expires in the year 2000, the vBNS has established a new network offering [called *Next Generation Network* (NGN)], to which schools and other entities may connect if the vBNS/NSF partnership is not renewed.

Measurement Tools in vBNS

One of the outcomes of the vBNS program has been the development of a variety of high-performance measurement tools. One such tool, called *OC-3mon* (and now, *OC-xMon*), was developed to allow passive capture (using optical splitters) of ATM cell and IP header information, to facilitate high-speed flow characterization. More detail is available at the vBNS Web site^[2]. Recently, further development of OC-xMon has been undertaken by the *Cooperative Association for Internet Data Analysis* (CAIDA)^[8]. CAIDA has perhaps the best collection of high-performance public-domain measurement and analysis tools in the world, and its Web site is definitely worth browsing.

The second backbone network to which Internet2 members can connect is called *Abilene*^[9]. Abilene was constructed by the *University Corporation for Advanced Internet Development* (UCAID) in collaboration with three industrial partners and Indiana University (IU). Partner contributions include fiber capacity from Qwest, SONET gear from Nortel, and routers from Cisco. The Abilene Network Operations Center (NOC) is staffed and operated by Indiana University. The network uses OC-48c POS interior links that initially connect ten routers in a partial mesh (a few interior links started as OC-12c, but are being upgraded). Abilene participants can connect at OC-3c or OC-12c, using either POS or ATM. See [10] for details on the router hardware architecture. For an insightful look at a research project that shows how this architecture can scale, see the second link in^[10] and also see Stanford Professor Nick McKeown's Tiny Tera homepage at^[11].

Measurement Tools in Abilene

It's worth spending a bit of time at the Abilene NOC Web site^[12]. One of the interesting tools developed there is the "Abilene Weather Map"^[13]. Abilene NOC has indicated that it will make source code for this tool available to Internet2 schools.

Gigapop Technology Survey

Internet2 schools can connect to either vBNS or Abilene directly. However, it is also common for several schools to converge their links at a "gigapop." This gigapop then connects to Abilene and/or vBNS, and possibly to commodity Internet Service Providers (ISPs) (to carry the "Commodity Internet" traffic of the school). Additionally, non-Internet2 schools, libraries, K-12, and state government networks also often converge at gigapops. Non-Internet2 schools typically don't forward traffic over Abilene or vBNS. But the common meeting point allows local exchange of local traffic, often affords larger aggregate commodity Internet connectivity for the gigapop participants, and allows direct access to other services that might be offered at the gigapop (Web caching, and so on).

The connectivity architecture used at gigapops varies widely. Detailed documentation for several gigapops can be found at^[14]. Some Gigapops are "Layer 2," meaning that each participant is responsible for exchanging routes and traffic among themselves directly. More often, gigapops are "Layer 3," meaning that the gigapop provides a router with which gigapop participants peer. The gigapop router then typically exchanges traffic with vBNS and/or Abilene, and possible commodity ISPs.

Some gigapops are implemented at a single site (for instance, *Metropolitan Research and Education Network* [MREN], *Southern Crossroads* [SoX]), while others are "distributed gigapops," meaning gigapop equipment exists at multiple locations (for instance, the *California Research and Educational Network* [CalREN-2], and *The Great Plains Network* [GPN]). Following are a couple of specific gigapop examples.

MREN

The MREN^[15] is built on a Layer 2 gigapop near Chicago that joins schools and research facilities from Illinois and several states in the Midwest. MREN members typically connect with OC-3c ATM links. Since MREN is a Layer 2 gigapop, the border router of each member peers directly with the border routers of other members. Additionally, each member's border router might peer with the Chicago-area vBNS or Abilene border router. vBNS and Abilene routers (as well as several other national research and international networks) peer here. Physically, the facility is built upon the Network Access Point (NAP) facility provided by Ameritech Advanced Data Services (AADS)^[16]. Routers typically peer with each other via ATM UBR *Permanent Virtual Paths* (PVPs), although other arrangements are possible.

CENIC/CalREN-2

The Corporation for Education Network Initiatives in California (CENIC)^[17] has constructed CalREN-2. The CalREN-2 distributed gigapop is interesting in several respects. First, as the name implies, it represents a distributed gigapop. In this case, three separate SONET ring facilities provide connectivity for Northern, Central (Los Angeles area), and Southern California schools. These three regions are linked to each other, and also to external networks.

Second, in each ring, there are two sets of OC-12c connections to each adjacent school. CalREN-2 has currently utilized these connections to construct both a ring of ATM connectivity, and a separate, parallel ring of POS connectivity. As a result, CalREN-2 is uniquely positioned to experiment simultaneously with both ATM and POS connectivity, performance, and QoS characteristics.

Third, to take the Northern schools as an example, the ring structure allows for a variety of Layer 3 topologies to be explored. For example, in a ring with these size and bandwidth characteristics, what are the trade-offs on application-level performance of inducing more hops while keeping the per-hop bandwidth high, versus dividing the bandwidth into smaller slices but creating a partial mesh that reduces the average Layer 3 hop count?

Engineering Challenges

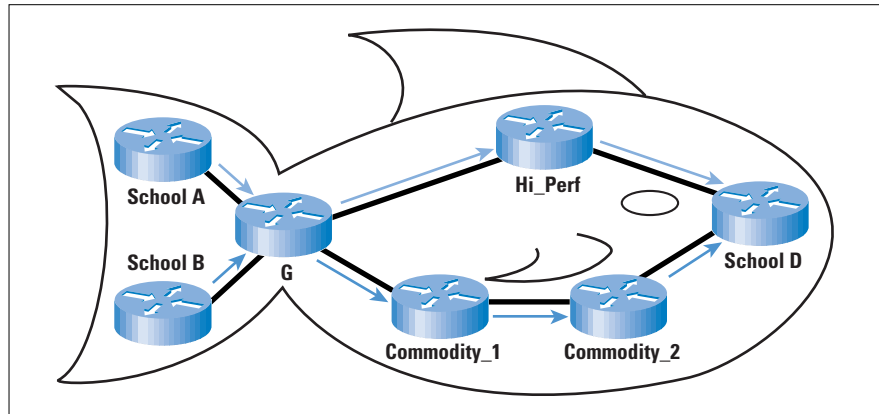
This section looks at some of the engineering challenges present in Internet2. They revolve around enabling applications with new network services, implementing appropriate policy, and doing all of this at high speed. Specifically, we'll look at Explicit Routing, Multicast, and Quality of Service.

Explicit Routing—The Fish Problem

The condition that several schools often converge at a gigapop, combined with the constraint that sometimes the funders of high-performance connectivity require that only the funded schools are allowed to use the high-performance connection, gives rise to a need for

“explicit routing” at the gigapop. The gigapop can forward packets through either a high-performance connection, or through the commodity Internet. Usually, for a single destination, traditional routing would have the gigapop use the “best” path to forward all packets to a particular destination. But when multiple policies must be implemented at the gigapop, the gigapop router must be able to “override” normal routing and forward packets on a path that’s not the “best.” A concrete example is shown in Figure 1.

Figure 1: The “Fish Problem”



Consider packets from schools A and B, both headed for destination D. Assume both schools are connected to gigapop G, and that G has two paths to D; one along *G-Hi_perf-D*, and the other along *G-Commodity_1-Commodity_2-D*. Further assume that A is allowed to use either path (and would prefer *G-Hi_perf-D*), but that B is prohibited from using *G-Hi_perf-D*. This scenario describes the “Explicit Routing Problem,” and since it is often drawn in a shape resembling a fish, is also known as the “Fish Problem.” The essence is that a routing decision at G must be made on some other criteria than just the destination IP address.

A couple of solutions to the fish problem have been used in the past, but they tend to have problems with either speed or scalability. For example, “policy routing,” which usually includes a method to look at both source and destination address, has historically shown low performance. Inserting ATM switches and using virtual circuits has been used in some cases, but this solution has scaling problems and requires extra equipment. Today, many Internet2 gigapops use a separate router per policy. In the case of needing two policies in the example above, this means two routers. This solution is expensive, but does have high performance.

One promising idea is to implement enough of the “policy routing” process in hardware to allow high-speed *source+destination+other_bits* lookups. While straightforward in concept, some point out that even with line-rate source-address routing capability, the method is flawed because it requires significant manual configuration, and is prone to creating black holes for traffic upon link failure. Proponents suggest that these shortcomings can be overcome.

Another promising mechanism is becoming available as a result of work done to facilitate *Multiprotocol Label Switching* (MPLS) in routers and switches. The idea here is that one of the underlying pieces of technology required for MPLS is “multi-FIB” (multiple *Forwarding Information Bases*). Instead of the traditional “single-FIB,” which always uses “the best” route to a destination, multi-FIB allows multiple forwarding tables to exist in a single router. This setup will allow a gigapop to implement multiple policies in one router, rather than the “one box per policy” that several gigapops have used previously. Note that in the case of a gigapop with a single router on which all members converge, multiple policies can be achieved with multi-FIB without actually using MPLS *Label-Switched Paths* (LSPs). For more complex gigapops, where members themselves may converge high-performance-eligible and ineligible traffic before forwarding on a single link to the gigapop, one might consider using simple LSPs to present the gigapop with traffic that is predifferentiated.

Multicast

Many of the applications in Internet2 schools use multicast. In addition to flows for videoconferencing or distance learning that use MPEG-1 (or slower) rates, a wide variety of applications require high-performance, scalable multicast. Examples include high-resolution immersive environments, collaborative real-time medical image diagnosis, and high-fidelity conferencing or distance learning (for instant, digital video camera rates of 30 mbps). When the Internet2 project began, many schools were on the *Multicast Backbone* (Mbone), and used *Distance Vector Multicast Routing Protocol* (DVMRP) tunnels to participate in multicast. Over the past year, one of the strong areas of collaboration between the Internet2 schools and the vendor community has been to develop and implement a migration strategy that allows Internet2 backbones, gigapops, and schools to move toward high-performance, scalable, native multicast support.

At the Internet2 conference in San Francisco in September 1998, the vBNS backbone was exposed to unprecedented levels of multicast stress. In a somewhat painful, but worthwhile, learning experience, it was concluded that *Protocol Independent Multicast-Dense Mode* (PIM-DM) did not scale well in highly meshed, high bitrate backbones. As a result, the vBNS has shifted to *PIM-Sparse Mode* (PIM-SM), and Abilene is being constructed with PIM-SM.

The current set of multicast components being applied in Internet2 (and leading ISPs) include: PIM-SM, *Multicast Border Gateway Protocol* (MBGP), and the *Multicast Source Discovery Protocol* (MSDP). MBGP allows distribution of routing information such that unicast and multicast routing can use noncongruent topologies.

MSDP allows independent domains to exchange information about multicast sources without creating interdomain Rendezvous Point (RP) dependencies. As they become standardized, it is expected that the *Border Gateway Multicast Protocol* (BGMP) and *Multicast Address Set Claim* (MASC) will be added to this infrastructure set.

Quality of Service

An area of broad interest in the Internet2 community centers on *Quality of Service* (QoS). The heart of QoS involves establishing strategies through which applications can be assured access to appropriate network resources when required. Typical examples of resources include end-to-end bandwidth, latency, or jitter. Of course, collateral issues and dimensions abound, including end-to-end vs. segment-only QoS; signaled vs. static provisioning; amount of state required by various approaches; level of granularity, precision, and strength of QoS “guarantee;” AAA issues; and reliability and recovery dynamics.

In an effort to start small, but make concrete progress, the Internet2 QoS working group^[18] has launched an experiment called the *Qbone*^[19]. Participants include backbone networks, gigapops, and individual schools and research labs worldwide. The Qbone will focus on deploying and using components developed by the Internet Engineering Task Force’s (IETF) *Differentiated Services* working group (Diffserv)^[20].

The initial Qbone plan is to deploy an approximation to the Expedited Forwarding (EF)^[21] forwarding behavior. The Qbone will start by statically allocating a small amount of EF bandwidth across boundaries between Autonomous Systems (ASs) to allow small EF flows among arbitrary combinations of schools/labs. Large flows, in these early stages, will have to be handled manually (much as they are today). In later stages the plan is to use *Bandwidth Brokers* (BBs) currently under development^[22] to aid in the automation of adjusting resource commitments between ASs (using interdomain BBs), and to aid in accepting application resource requests (using intradomain BBs, combined with policy servers and AAA mechanisms). The precise mechanics for BB interaction, trade-offs among signaling frequency, amount of state, scalability, and so on are certainly topics of research, but that’s part of what makes Qbone participation fun!

Summary

There is no single application or technology that makes Internet2 unique or exciting. Rather, the effort required to enable new applications that have strong bandwidth, latency, jitter, and coordination requirements has resulted in an infusion of energy from a variety of disciplines. Internet2 requires stretching existing technologies (ATM, POS, multicast, measurement), nurturing developing technologies (Quality of Service, explicit routing, Dense Wave-Division Multiplexing [DWDM], mobility), and participating in the invention of new technologies (all-optical infrastructures, extending AAA, and other resource allocation and

scheduling middleware). Internet2 requires attention to maturing components in backbone, gigapop, and campus environments in order to deliver on the promise of speedy transference of lessons learned to the commodity Internet. The effort so far has resulted in demonstration of truly stunning, impactful, and useful applications. It is the convergence of effort and rapid rate of change that makes Internet2 a challenging and rewarding endeavor.

Other Initiatives

Although this article has focused on aspects of Internet2 in the United States, there are many advanced Internet activities around the world. A partial list includes:

<http://www.dante.net/ten-155.html> (Europe)
<http://www.ukerna.ac.uk> (UK)
<http://www.dfn.de> (Germany)
<http://www.renater.fr> (France)
<http://www.surfnet.nl> (The Netherlands)
<http://apan.or.kr> (Asia/Pacific)
<http://www.singaren.net.sg> (Singapore)
<http://www.canet3.net> (Canada)
<http://www.cudi.edu.mx> (Mexico)
<http://www.reuna.cl> (Chile)
<http://www.ngi.gov> (U.S. Federal)
<http://www.startap.net> (International peering)

A more complete list of advanced Internet initiatives is maintained at:

<http://www.cisco.com/aia>

References

- [1] <http://www.internet2.edu>
- [2] <http://www.vbns.net>
- [3] See latest press release at:
<http://www.nsf.gov/od/lpa/news/press/99/pr9915.htm>
...and updated program announcement at:
<http://www.nsf.gov/pubs/1998/nsf98102/nsf98102.txt>
- [4] <http://apps.internet2.edu>
- [5] <http://i2dv.nwu.icaair.org/> and <http://www.icaair.org/>
- [6] <http://dsi.internet2.edu/>
- [7] <http://evlweb.eecs.uic.edu/pape/CAVE>
...has a great introduction to CAVE technology.
Also see the *Electronic Visualization Laboratory* homepage at:
<http://www.evl.uic.edu/EVL/index.html>
- [8] <http://www.caida.org>
- [9] <http://www.ucaid.org>, and Abilene specifics at:
<http://www.internet2.edu/abilene>

Abilene router details are at:

[10] <http://www.cisco.com/warp/public/cc/cisco/mkt/core/12000/index.shtml>

...and Nick McKeown's paper is at:

<http://www.cisco.com/warp/public/cc/cisco/mkt/core/12000/tech/fastwp.pdf>

[11] <http://tiny-tera.stanford.edu/tiny-tera/index.html>

[12] <http://www.abilene.iu.edu>

[13] <http://hydra.uits.iu.edu/~abilene/traffic>

[14] Following are several gigapop sites:

California's CENIC/CalREN2: <http://www.cenic.org>,

The Pacific/Northwest gigapop: <http://www.pnw-gigapop.net>

The Great Plains Network: <http://www.greatplains.net>

The Southern Crossroads, with members from Southeastern Universities Research Association: <http://www.sox.net>

MidAtlantic Crossroads: <http://www.networkvirginia.net/MAX>

MREN: <http://www.mren.org>

WestNet: <http://www.scd.ucar.edu/nets/Projects/Westnet>

North Carolina Gigapop: <http://www.ncni.net>

The Texas Gigapop: <http://noc.gigapop.gen.tx.us>

Northern Crossroads: <http://www.nox.org>

Philadelphia area Magpi: <http://www.magpi.net>

Pittsburgh-based NCNE: <http://www.ncne.net>

New York: <http://www.nysernet.org>

[15] <http://www.mren.org>

[16] <http://www.aads.net>, and <http://nap.aads.net/main.html>

[17] <http://www.cenic.org>

[18] <http://www.internet2.edu/qos/wg>

[19] <http://www.internet2.edu/qos/qbone>

[20] <http://www.ietf.org/html.charters/diffserv-charter.html>

[21] <http://www.ietf.org/rfc/rfc2598.txt>

[22] <http://www.merit.edu/working.groups/i2-qbone-bb>

LARRY DUNN is the Technology Development Manager in the Advanced Internet Initiatives Division at Cisco Systems. He serves on the Internet2 Quality of Service and Routing working groups. After receiving his PhD from the University of Minnesota (Electrical Engineering '92), he served as Director of Networking there, and subsequently as Director of Strategic Markets and Applications (Education) for FORE Systems. He periodically teaches Advanced Networking courses at the University of Minnesota. Research interests include test vector generation for combinational logic, network design and analysis, and Quality of Service techniques and deployment strategies.

E-mail: ldunn@cisco.com

One Byte at a Time: Internet Addressing

by Peter H. Salus

The source of all knowledge where the Internet is concerned is the set of *Requests for Comments* (RFCs). Because there are now well over 2,700 RFCs, however, only a few people track history, evolution, and outright paradigm shift.

Each node on the Internet—router or end system (often called “host” or “server”)—has a unique identifier attached to it; this identifier is its *address*. Any packet sent between nodes must use the destination address to tell the intervening routers where it should go.

In RFC 1 (April 1969), Steve Crocker laid out a scheme that allotted five bits to address space: enough for 32 addresses. By September 1969, when *Interface Message Processor* (IMP) No. 1 was installed in Kleinrock’s lab at UCLA, this number had grown to six bits (63 addresses). By 1972, it had become apparent that this number would be insufficient, and the address space was enlarged to eight bits (255 addresses). In fact, the *Advanced Research Projects Agency Network* (ARPANET) hit only 63 hosts in January 1976. This number was, however, already a lot in terms of the **HOSTS.TXT** tables that were distributed to every site. By August 1983, there were 213 hosts, and the eight-bit address barrier was being pushed.

Cerf’s original version of TCP (RFC 675; December 1974) and Postel’s of IP (RFC 760; January 1980) increased this “address space” to 32 bits, but the structure of the ARPANET was “flat,” that is, the hierarchical distributed name-to-address database we are familiar with only came about with Mills’ conceptualization of the *Domain Name System* (DNS) (RFC 799; September 1981), and its implementation by Paul Mockapetris (RFCs 882 and 883; November 1983).

Address Classes

The Internet Protocol uses a 32-bit addressing scheme and originally four classes of networks: A, B, C, D. (See Figure 1 on page 5). There are only 128 Class A networks, but each can have 16,777,216 unique host identifiers. Next, there are 16,384 Class B networks, with 65,535 unique identifiers; 2,097,192 Class C networks, with 255 hosts; and over 268 million Class D multicast groups. (A fifth class, Class E, is reserved and not available for general use).

Address Depletion

Using the 32-bit IP addressing scheme allowed for about 4 billion hosts on 16.7 million networks. Although this number of various kinds of addresses seemed like a lot, the expansion of the use of the Internet over the past decade has been explosive, and the original address classes did not allow for a flexible address assignment based on an organization’s particular need.

In August 1990 during the Vancouver *Internet Engineering Task Force* (IETF) meeting, Frank Solensky, Phill Gross, and Sue Hares projected that the current rate of assignment would exhaust the Class B space by March 1994.

CIDR

Classless Inter-Domain Routing or CIDR (RFCs 1518 and 1519; September 1993) was introduced to improve both routing scalability and address space utilization in the Internet. By eliminating the notion of “network classes,” CIDR allows for a better match between address requirements and address allocation. This results in expansion of the scope of hierarchical routing, which in turn improves scaling properties of the Internet routing system. CIDR has proven to be the palliative that has enabled the Internet to continue functioning while growth continues.

Even with this palliative, it was predicted in 1994 that, using the current allocation statistics, the Internet will exhaust the IPv4 address space between 2005 and 2011. With five more years of experience, which has also brought greater uncertainty as to gross numbers, we can push these dates out a bit, but exhaustion will come eventually.

Another factor that has slowed down the address depletion rate is the use of *Network Address Translation* (NAT). NAT devices allows an organization to have one external (“public”) address and many private (net 10 is often used) addresses internally. Since the internal addresses are not “seen” from the outside, they do not need to be globally unique. This approach has downsides (some protocols weren’t designed with NATs in mind), but from the address depletion point of view, it is a win. RFC 1597 describes “Address Allocation for Private Internets.”

If you are interested in current Internet addressing, an excellent book is available: *TCP/IP Addressing*, by Buck Graham, AP Professional, 1997. Graham does an excellent job on addressing, routing, and the various bizzarries involved in optimal routing, efficient use of address space, and making network management less onerous. This book is, however, not intended to be for elementary instruction; Graham primarily speaks to the professional market.

IPng aka IPv6

In the summer of 1994, the IETF set up an Internet Protocol next generation (IPng) task force, cochaired by Scott Bradner and Allison Mankin. (IPng later became known as IPv6 for “IP version 6”). Recommendations from that task force were released in October 1994 for discussion at the December 1994 IETF meeting. The basic goal was to have something in place before 2000, so that the time limit would not be pushed.

Unfortunately, as Bradner and Mankin stated in their recommendation: “Some people pointed out that this type of projection makes an assumption of no paradigm shifts in IP usage. If someone were to develop a new ‘killer application,’ (for example, cable-TV set top boxes), the resultant rise in the demand for IP addresses could make this an over-estimate of the time available.”

IPv6 provides for 128-bit addressing. This number is gigantic: larger than the estimated total number of molecules in the universe.

Books

Two noteworthy books are available on IPv6 itself: Christian Huitema’s *IPv6: The New Internet Protocol* (ISBN 0-13-241936-X, Prentice Hall, 1996) and Scott Bradner and Allison Mankin’s anthology *IPng* (ISBN 0-201-63395-7, Addison-Wesley, 1996), which provides an explanation of the task force’s process and explicates the services that are provided for (as, for example, ATM support). These books are both dated, but they are the best available now. Keeping up with what’s going on is easy, thanks to the IETF’s Web site <http://www.ietf.org>.

An excellent business and technical case for IPv6 is found in the Internet Architecture Board draft by Steve King and several colleagues (**draft-iab-case-for-ipv6-05.txt**). Other works in progress deal with the adjustments to Open Shortest Path First (OSPF), multicasting, mobility, and so on.

Transition

The period from 1981 through 1983—the time of conversion to DNS—was painful to all concerned. Over the past 15 years we have learned a lot, but the switch from IPv4 to IPv6 may be yet more painful. The drafts tell the tale of those who are striving to make things easier.

There has been much discussion about various kinds of transition mechanisms, and some of these may be less painful (more automated) than we might at first think. Remember, this pain is not because of the innate difficulty, but veering a ship that carries fewer than 250 passengers is far easier than veering a ship that carries 60 million. Some members of the community think that the pain may not justify the gain. The author is not one of them. It has been nearly 20 years since TCP/IP was made official, yet there are still UUCP networks.

In the author’s opinion, IPv6 will be here in a few years, if not sooner.

Reference

- [1] Fink, R., “IPv6—What and Where It Is,” *The Internet Protocol Journal*, Volume 2, No. 1, March 1999.

PETER H. SALUS is the author of *A Quarter Century of UNIX* (1994) and *Casting the Net: From ARPANET to Internet and Beyond* (1995). He is the Editor in Chief of *The Handbook of Programming Languages* (1998). His e-mail address is: peter@pedant.com

Book Review

An Engineering Approach to Computer Networking

An Engineering Approach to Computer Networking: ATM Networks, the Internet and the Telephone Network, Srinivasan Keshav, ISBN 0-201-63442-2, Addison-Wesley, 1997, <http://www.awl.com/cseng/titles/0-201-63442-2/>

The rapid convergence of telephone and data networks brings with it a collision of two diverse approaches to fundamental network design. This “New World,” as it is often called, requires us to understand both the analog-to-digital evolution of the voice network, with its redundant search for faultless reliability, and the persistent tolerance of the data network. Mirroring the industry trend, this book explores the three major networking technologies: ATM, the Internet, and telephone networks, with the idea that the design of any modern network requires consideration of the influence of at least two of the three technologies.

This book is a textbook. Keshav himself declares in the preface that “textbooks, almost by definition, tend to be boring,” and the reader will recall this subtle warning shortly into Chapter 2. This is definitely a book for those who have at least an intermediate knowledge of data networking and a need to understand the component parts of network implementations. Keshav takes a true engineering approach, in that he attempts to teach the building blocks of the major networking technologies—and this approach is what makes the book one of my all-time favorites. By examining the component parts and why they are required, Keshav leaves you prepared to engineer a network that meets any number of diverse criteria.

Organization

The book is organized into three sections. Section 1 gives an introduction to the future of data and voice networks and then introduces three of the major networking technologies. This section also gives an overview of the historic construction of networks, along with some fundamental definitions of some of the engineering principles by which networks function. As early as Chapter 1, Keshav explores the engineering philosophy behind common network technologies, illustrating the theories that underlie their design. My favorite example is his suggestion that the telephone network was engineered to be intelligent because its endpoints, the telephones, are simply dumb. While this sounds obvious, it provides a fundamental perspective on the design of the system that proves invaluable to understanding the origin of the various “components” of the network.

Section 2 begins with a short but requisite review of protocol layering and, after a brief discussion of common design constraints, begins to dissect the major components required of almost any network implementation. Chapter 8 is a fairly comprehensive review of switching and, as the book's title suggests, the chapter is full of comparative anatomy. Read this chapter for its valuable insight into why various switching mechanisms have emerged and for its comparison of how various switching functions are handled on three major networking technologies. Chapter 9 deals with scheduling network resources, with an excellent comparison of the variety of scheduling mechanisms and their effect on connections and packets. It covers policy considerations that are also required of scheduling disciplines, giving the reader a set of strategies for network design. Chapter 11 covers routing of packets as well as routing in the telephone network. In my opinion, this discussion alone makes this book a required part of any networking professional's library. Admittedly, there are books that better explain routing in both of these environments, but because of the proximity of the topics, this presentation helps the reader to understand the mechanics of both systems in a way that provides insight into the inherent issues posed by both technologies.

Section 3 pulls together the various component functions discussed in Section 2 and explains some of their implementation in the form of protocols. Section 3 is a short section, probably not intended as a thorough survey of networking protocols. Keshav documents an excellent set of references for Section 3, however, and leaves it up to the reader to pursue those that are relevant to his or her professional development.

Required Reading

An Engineering Approach to Computer Networking is definitely an A+ book, and should be required reading for anyone interested in the inner workings of data and voice networks. Although the author expects the reader to absorb quite a bit in every chapter, the time spent is well invested. The book is a refreshing alternative in that it provides an answer to the question of "why" the network works rather than being another treatise on "how" the network works.

—Jim LeValley, Cisco Press
levalley@cisco.com

Would You Like to Review a Book for IPJ?

We receive numerous books on computer networking from all the major publishers. If you've got a specific book you are interested in reviewing, please contact us and we will make sure a copy is mailed to you. The book is yours to keep if you send us a review. We accept reviews of new titles, as well as some of the "networking classics." Contact us at ipj@cisco.com for more information.

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

Fragments

Internet Policy Institute Launched

On November 9th, 1999 a group of distinguished Internet visionaries and scholars announced the creation of the *Internet Policy Institute*, the nation's first independent, nonpartisan think tank devoted exclusively to providing research and hard data on the Internet and society. The group also announced its first research project and an initiative aimed at educating the presidential contenders.

The creation of the new think tank was announced by Jim Barksdale, former CEO of Netscape, Vint Cerf, Senior Vice President of Internet Architecture of MCI WorldCom, Esther Dyson, author and Chairman of EDventure Holdings, Inc., Mario Morino, Chairman of The Morino Institute, and Kimberly Jenkins, President of the Internet Policy Institute.

The new, nonprofit think tank will employ well-known experts and scholars to research subjects ranging from the role of the Internet in privacy to the Internet's impact on taxation and health care.

"The Internet is surrounded by noise, hype, rumors, marketing, IPOs and the hopes of starry-eyed start-ups, but there is very little hard data on which policymakers can base critical decisions that will determine the future of the new medium and how it affects society," said Barksdale, co-chairman of the Internet Policy Institute's Board of Directors. Wayne Clough, President of Georgia Tech, is his co-chairman.

"The speed at which society has adopted the Internet is unprecedented," said Cerf, who was Chairman and founding president of the Internet Society, as well as one of the designers of the TCP/IP protocol. "If, as we expect, half the world will be online within the next four years, we must make sure that the policy decisions we make now are based on solid, well-researched data."

The Institute announced its first research project, to be undertaken in collaboration with The Brookings Institution, on "The Economic Pay-off from the Internet Revolution." The research will be led by Alice Rivlin, former vice chair of the Federal Reserve System's Board of Directors and former Office of Management and Budget director, now with the Brookings Institution, and Robert E. Litan, Vice President and Director of Economic Studies at The Brookings Institution and former associate director of the Office of Management and Budget. The research will produce the first comprehensive, systematic economic study by an independent research group of the subject.

The nature and extent of the impact is of special importance to macroeconomic policy—specifically monetary policy—to the extent that the Net is having or will have a material and sustained impact on the growth rate of productivity. The impact the Net has on specific industries, and the way it affects barriers to entry, has important implications for antitrust and regulatory policy.

Exactly one year before the next presidential election, the Internet Policy Institute also announced its first publications project, “Briefing the President: What the Next President of the United States Needs to Know About the Internet and Its Transformative Impact on Society.” The Institute also released the introduction to the project by Barksdale, while Cerf outlined the contents of the next paper, “What is the Internet (and What Makes It Work)” that will be released December 1. Over the course of the coming months, the Institute will release 13 papers to be presented in briefings to all the leading presidential contenders and later compiled into a book.

“We didn’t know five years ago the direction that the Internet would take,” Barksdale said. “I’ll bet that five years from now, we’ll be surprised by its new directions. We need to assure that an honest, objective approach is taken on Internet issues, to prevent decision making that hinders the potential of this amazing medium,” he said. For more information see: <http://www.internetpolicy.org>

APRICOT 2000

The *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) will be held in at the Intercontinental Hotel in Seoul, Korea from February 28th to March 2nd, 2000. APRICOT provides a forum for key Internet builders in the region to learn from their peers and other leaders in the Internet community from around the world. The week-long summit consists of seminars, workshops, tutorials, conference sessions, and birds-of-a-feather sessions—all with the goal of spreading and sharing the knowledge required to operate the Internet within the Asia Pacific region. For more information see:

<http://www.apricot.net>

More on Web Caching

If you enjoyed the article on Web Caching in our September 1999 issue, you might find the following paper of interest: “A Survey of Web Caching Schemes for the Internet,” by Jia Wang. You can find this article in the October 1999 issue of ACM SIGCOMM’s *Computer Communications Review* (Volume 29, Number 5). The paper is also available on line in either PostScript or PDF format:

<http://www.acm.org/sigcomm/ccr/archive/1999/oct99/ccr9910-jia-wang.html>

ICANN Update

On September 28, 1999, the United States Department of Commerce, Network Solutions, Inc. (NSI), and The Internet Corporation for Assigned Names and Numbers (ICANN) announced a series of agreements they had tentatively reached to resolve outstanding differences among the three parties. On November 4, 1999, based on public comment in writing and at a public forum held at the 1999 ICANN annual meeting, the ICANN Board approved revised versions of these agreements. The agreements were signed by the three parties on November 10, 1999. The full text of the agreements can be found on the ICANN Web site at www.icann.org. Here we include some highlights:

- NSI will operate the registry for the **.com**, **.net**, and **.org** top-level domains according to requirements stated in the agreement and developed in the future through the ICANN consensus-based process. All accredited registrars will have equal access to this registry.
- A revised registrar accreditation agreement between ICANN and registrars was adopted. To continue to register names with the **.com**, **.net**, and **.org** registry operated by NSI after November 30, 1999, registrars must have entered a new Registrar License and Agreement with NSI and the revised ICANN accreditation agreement.
- A revised NSI-Registrar License and Agreement was created under which competitive ICANN-accredited registrars are permitted to place and renew registrations in the registry.
- An amendment was made to Cooperative Agreement #NCR 92-18742 originally entered between NSI and the National Science Foundation (NSF) in 1992. On October 7, 1998, NSI and the United States Department of Commerce (which by then had assumed the NSF's role as lead agency of the U.S. Government) entered an Amendment 11 to that Cooperative Agreement under which NSI agreed to implement a shared registration system in which competitive registrars would enter registrations into the **.com**, **.net**, and **.org** registry on an equitable basis. Amendment 19 solidifies those arrangements and provides that in operating the registry NSI will abide by consensus policies adopted in the ICANN process.

At the annual meeting in early November, nine new directors joined the ICANN Board of Directors. They are Robert Blokzijl, Ken Fockler and Pindar Wong named by the The Address Supporting Organization (ASO); Amadeu Abril i Abril, Jonathan Cohen and Alejandro Pisanty named by the Domain Name Supporting Organization (DNSO); Jean-François Abramatic, Vinton G. Cerf and Philip Davidson named by the Protocol Supporting Organization (PSO).

The newly expanded ICANN Board will take on a major challenge in 2000 in its consideration of contending proposals for the future of Top Level Domains. After years of vociferous argument, the DNS community is no closer than it ever has been to a consensus on whether new name registries should be created, and if so, with what structure and registration rules.

Interplanetary Internet Special Interest Group Formed

The Internet Society (ISOC) recently announced the formation of the Interplanetary Internet Special Interest Group (IPNSIG). The IPNSIG exists to allow public participation in the evolution of the Interplanetary Internet. The technical research into how the Earth's Internet may be extended into interplanetary space has been underway for several years as part of an international communications standardization body known as the Consultative Committee on Space Data Systems (CCSDS). (See <http://www.ccsds.org/>)

The CCSDS organization is primarily concerned with communications standardization for scientific satellites, with a primary focus on the needs of near-term missions. In order to extend this horizon out several decades, and to begin to involve the terrestrial internet research and engineering communities, a special Interplanetary Internet Study was proposed and subsequently funded in the United States.

The Interplanetary Internet Study is funded by the Defense Advanced Research Projects Agency's Next Generation Internet Initiative, and presently consists of a core team of researchers from the NASA Jet Propulsion Laboratory, MITRE Corporation, SPARTA, Global Science & Technology and consulting researchers from The University of Southern California Information Sciences Institute, University of California Los Angeles and the California Institute of Technology. The primary goal of the study is to investigate how terrestrial internet protocols and techniques may be extended and/or used as-is in the exploration of deep space. The study team has also founded the IPNSIG and has formed the core of an Interplanetary Internet Research Group under the sponsorship of the Internet Research Task Force (IRTF).

The NASA IPN Study Team will act as liaison between the satellite and space communities and the ISOC/IRTF communities. The NASA IPN Study Team will assist with requirements and understanding of the deep space environment and missions, while the primary research on new or modified protocols will be conducted by the IRTF. In addition, the NASA Study Team will also act as liaison with the CCSDS.

The NASA Study Team will also enable simulated and actual opportunities to test protocols and the use of internet techniques in the space environment. For more information, visit: ipn.jpl.nasa.gov/

This publication is distributed on an "as-is" basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Engineering
MCI WorldCom, USA

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Edward R. Kozel, Member of The Board of Directors
Cisco Systems, Inc., USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the Cisco News
Publications Group, Cisco Systems, Inc.
www.cisco.com*

*Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 1999 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-10/5
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

Bulk Rate Mail
U.S. Postage
PAID
Cisco Systems, Inc.

The Internet Protocol Journal

March 2000

Volume 3, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

F R O M T H E E D I T O R

In This Issue

From the Editor	1
Routing IPv6 over IPv4.....	2
IP Security	11
QoS—Fact or Fiction?	27
Book Review.....	35
Call for Papers	38
Fragments	39

Work on a new version of the Internet Protocol, known as IPv6, has been under way for several years in the IETF. There is still some debate about when and how IPv6 will be deployed. Proponents of IPv6 argue that the demand for new IP addresses will continue to rise to a point where we will simply run out of available IPv4 addresses and that we should, therefore, start deploying IPv6 *today*. Opponents argue that such a protocol transition will be too costly and painful for most organizations. They also argue that careful address management and the use of *Network Address Translation* (NAT) will allow continued use of the IPv4 address space for a very long time. Regardless of the timeframe, a major factor in the deployment of IPv6 is an appropriate transition strategy that allows existing IPv4 systems to communicate with new IPv6 systems. A transition mechanism, known as “6to4,” is described in our first article by Brian Carpenter, Keith Moore, and Bob Fink.

In previous editions of this journal, we have looked at various security technologies for use in the Internet. Security mechanisms have been added at every layer of the protocol stack, and IP itself is no exception. IP Security, commonly known as “IPSec,” is being deployed in many public and private networks. In our second article, William Stallings describes the main features of IPSec and looks at how IPSec can be used to build Virtual Private Networks.

Our final article is a critical look at *Quality of Service* (QoS) in the Internet. The need to provide different priorities to different kinds of traffic in a network is well understood and the technical community has been hard at work developing numerous systems to address this need. Geoff Huston looks at the prospects of deploying QoS solutions that will operate across the Internet as a whole.

The Y2K transition has been described as a “nonevent” by many. However, the lessons learned and the collaborative coordination efforts that were put in place for this transition can hopefully be used in the future. A colleague of mine had to call a plumber to his house on New Year’s Eve. When he tried to pay for the repair with a credit card which had “00” as the expiration year, the plumber insisted that this meant the card was invalid. So while most systems were “Y2K compliant,” this particular plumber was clearly not. Do you have a Y2K story to share? Drop us a line at ipj@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

Connecting IPv6 Routing Domains Over the IPv4 Internet

by *Brian E. Carpenter, IBM & iCAIR*
Keith Moore, University of Tennessee
Bob Fink, Energy Sciences Network

A next-generation Internet Protocol^[1], known first as IPng and then as IPv6, has been under development by the *Internet Engineering Task Force* (IETF) for several years to replace the current Internet Protocol known as IPv4. The reasons behind the need for IPv6 are not covered here, but interested readers are encouraged to read “The Case for IPv6”^[2] for this background.

Of major importance during the development of IPv6 has been how to do the transition away from IPv4, and towards IPv6. The work on transition strategies, tools, and mechanisms has been part of the basic IPv6 design effort from the beginning. The current transition efforts, taking place at the *IETF IPng Transition Working Group* (ngtrans)^[3], will continue until it is clear that the transition will be successful.

These transition design efforts resulted in a basic Transition Mechanisms specification for IPv6 hosts and routers^[4] that specifies the use of a Dual IP layer providing complete support for both IPv4 and IPv6 in hosts and routers, and IPv6-over-IPv4 *tunneling*, encapsulating IPv6 packets within IPv4 headers to carry them over IPv4 routing infrastructures.

These concepts are heavily relied on for transition from the traditional IPv4-based Internet as we know it today, to an IPv6-based Internet. It is expected that IPv4 and IPv6 will coexist for many years during this transition.

Of great concern to transition strategy planners is how to provide connectivity between IPv6-enabled end-user sites (also known as *routing domains*) when they do not yet have a reasonable (or any) choice of *Internet Service Provider* (ISP) that provides native IPv6 transport services. One way to provide IPv6 connectivity between end-user sites (when native IPv6 service does not exist) is to use IPv6-over-IPv4 encapsulation (tunneling) between them, similar to the technique currently used in the 6bone^[5] IPv6 testbed network. This requires complexity for both end-user sites, and the networks providing the tunneling service (for instance, the 6bone backbone ISPs), in creating, managing, and operating manually configured tunnels.

The “6to4” transition mechanism, “Connection of IPv6 Domains via IPv4 Clouds without Explicit Tunnels”^[6], provides a solution to the complexity problem of using manually configured tunnels by specifying a unique routing prefix for each end-user site that carries an IPv4 tunnel endpoint address.

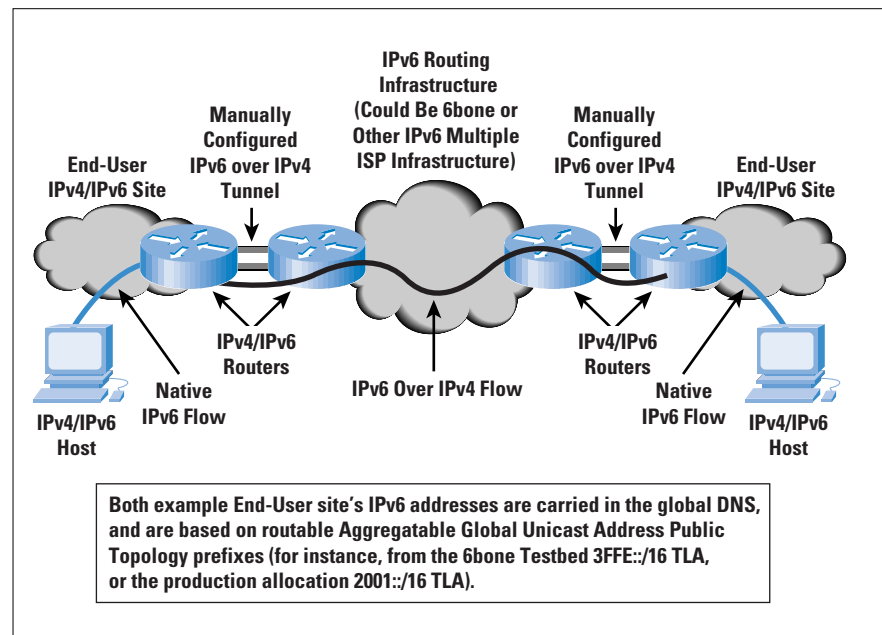
It should also be noted that each end-user site with as little as a single IPv4 address has a unique, routable, IPv6 site routing prefix thanks to the 6to4 transition mechanism.

Connecting IPv6 Routing Domains

When end-user site networks enable IPv6 in their local host and router systems, but have no native IPv6 Internet service, connectivity to other IPv6 routing domains across a worldwide Internet must be accomplished another way, or the value of a connected Internet is lost. Prior to the 6to4 transition mechanism, a site's network staff would have to rely on the manual configuration of IPv6-over-IPv4 tunnels to accomplish this connectivity.

This connectivity could be accomplished by arranging tunnels directly with each IPv6 site to which connectivity is needed, but more typically is done by arranging a tunnel into a larger IPv6 routing infrastructure that could guarantee connectivity to all IPv6 end-user site networks. (See Figure 1.) The 6bone IPv6 testbed was the first IPv6 routing infrastructure to provide worldwide IPv6 connectivity (starting in 1996), while more recently (late 1999) networks providing production IPv6 Internet service have also interconnected to provide this connectivity. In fact, the 6bone and production IPv6 routing infrastructures are well interconnected to guarantee worldwide IPv6 connectivity.

Figure 1: Configured Tunnel Overview



However, even given a solid, reliable, worldwide IPv6 routing infrastructure (similar to the IPv4-based Internet today), if an end-user site does not have a reasonable (or any) local choice for native IPv6 Internet service, a tunnel must be used.

The 6to4 mechanism addresses many of the practical difficulties with manually configured tunneling:

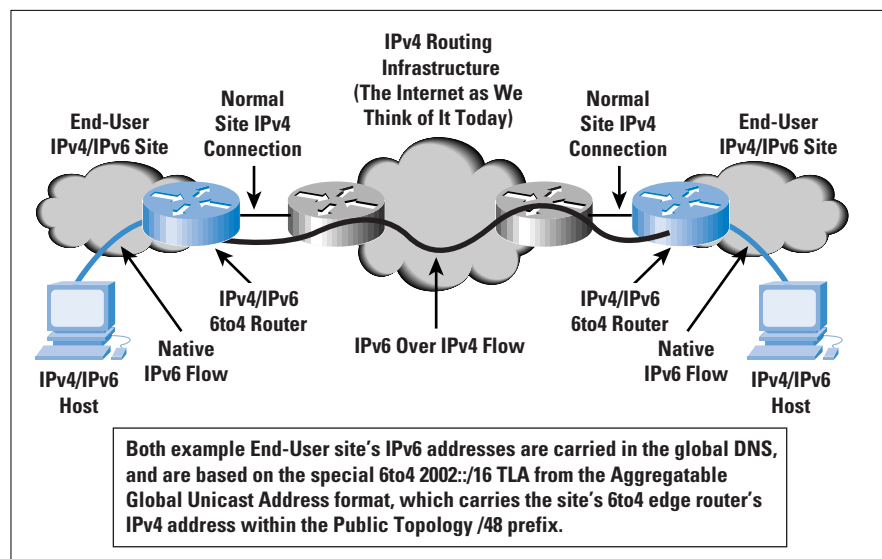
- The end-user site network staff must choose an IPv6 Internet service to tunnel to. This entails a process of at least three parts:
 - Finding candidate networks when the site's choice of IPv4 service does not provide IPv6 service (either tunneling or native),
 - Determining which ones are the best IPv4 path to use so that an IPv6-over-IPv4 tunnel doesn't inadvertently follow a very unreliable or low-performance path,
 - Making arrangements with the desired IPv6 service provider for tunneling service, a scenario that may at times be difficult if the selected provider is not willing to provide the service, or if for other administrative/cost reasons it is difficult to establish a business relationship.

Clearly it is easiest to use the site's own service provider, but in the early days of IPv6 transition this will often not be an option.

- An IPv6-over-IPv4 tunnel must be built to the selected provider, and a peering relationship must be established with the selected provider. This requires establishing a technical relationship with the provider and working through the various low-level details of how to configure tunnels between two routers, including answering the following questions:
 - Are the site and provider routers compatible early on in this process?
 - What peering protocol will be used (presumably an IPv6-capable version of the *Border Gateway Protocol Version 4* [BGP4]), and are the versions compatible and well debugged?
 - Have all the technical tunnel configuration issues between the site and provider been addressed?

Again, it is clearly easiest to perform all these steps if they are taken with the site's own IPv4 service provider.

Figure 2: 6to4 Tunnel Overview



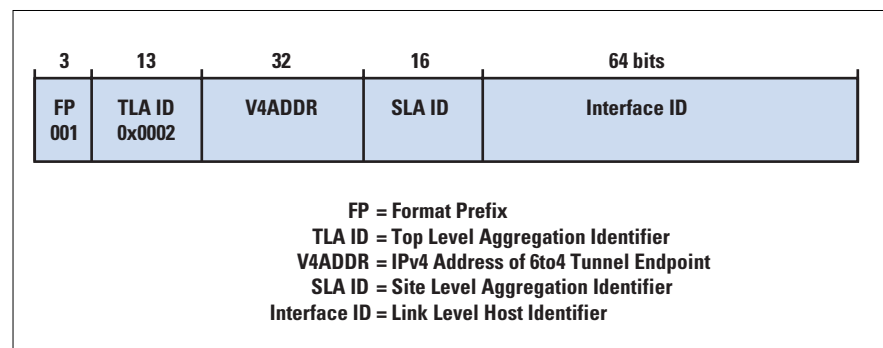
6to4 Eliminates Complex Tunnel Management

The 6to4 transition mechanism provides a solution to the complexity problem of building manually configured tunnels to an ISP by advertising a site's IPv4 tunnel endpoint (to be used for a dynamic tunnel) in a special external routing prefix for that site. Thus one site trying to reach another will discover the 6to4 tunnel endpoint from a *Domain Name System* (DNS) name to address lookup and use a dynamically built tunnel from site to site for the communication. (See Figure 2.) The tunnels are transient in that there is no state maintained for them, lasting only as long as a specific transaction uses the path. A 6to4 tunnel also bypasses the need to establish a tunnel to a wide-area IPv6 routing infrastructure, such as the 6bone.

The specification of a 48-bit external routing prefix in the IPv6 *Aggregatable Global Unicast Address Format* (AGGR)^[7] (see Figure 3) that provides just enough space to hold the 32 bits required for the 32-bit IPv4 tunnel endpoint address (called V4ADDR in Figure 3) makes this setup possible.

Thus, this prefix has exactly the same format as normal prefixes assigned according to the AGGR. Within the subscriber site it can be used exactly like any other valid IPv6 prefix, for instance, for automated address assignment and discovery according to the normal IPv6 mechanisms for this.

Figure 3: 6to4 Prefix Format



The Simplest Use of 6to4

The simplest scenario for 6to4 is when several sites start to use IPv6 alongside IPv4, and have no native IPv6 ISP service available. Thus each site identifies a router to run dual stack (that is, IPv4 and IPv6 together) and 6to4 tunneling, ensuring that this router has a globally routable IPv4 address (that is, not in private IPv4 address space).

It is assumed that this new 6to4 router is reachable by IPv6-capable hosts within the site. Although the various ways in which these hosts may be reached are not discussed in detail here, they include using IPv6-enabled site IPv4 routers, operating special IPv6-only routers in parallel with site IPv4 routers, using the “6over4” mechanism^[8], and employing other tunneling methods.

A new 6to4 site advertises the 6to4 prefix to its site via the *Neighbor Discovery* (ND) protocol^[9], which will cause IPv6 hosts at this site to have their DNS name/address entries to include the 6to4 prefix for the site in them.

In operation, when one IPv6-enabled host at a 6to4 site tries to access an IPv6-enabled host by domain name at another 6to4 site, the DNS will return both an IPv4 and an IPv6 IP address for that host, indicating that it is reachable by both IPv4 and IPv6. The requesting host selects the IPv6 address, which will have a 6to4 prefix, and sends a packet off to its nearest router, eventually reaching its site boundary router, which we assume has 6to4 service as well.

Sending and Receiving Rules for 6to4 Routers

When the requesting site's 6to4 router sees that it must send a packet to another site (that is, there is a nonlocal destination), and that the next hop destination prefix contains the special 6to4 *Top Level Aggregation* (TLA) value of 2002::/16, the IPv6 packet is encapsulated in an IPv4 packet using an IPv4 protocol type of 41, as defined in the *Transition Mechanisms RFC*^[4]. The source IPv4 address will be the one in the requesting site's 6to4 prefix (which is the IPv4 address of the outgoing interface to the Internet on the 6to4 router, and contained in the source 6to4 prefix of the IPv6 packet), and the destination IPv4 address will be the one in the next hop destination 6to4 prefix of the IPv6 packet.

When the destination site's 6to4 router receives the IPv4 packet, and recognizes that it has an IPv4 protocol type of 41, IPv4 security checks are made and the IPv4 header is removed, leaving the original IPv6 packet for local forwarding.

The sending rule above is the only modification to IPv6 forwarding, because the receiving rule was already specified for the basic IPv6 Transition Mechanism mentioned earlier^[4]. Along with advertisement of the 6to4 prefix by appropriate entries in the DNS, any number of sites can interoperate without manual tunnel configuration.

It is not necessary to operate an exterior routing protocol (for instance, BGP4+) for 6to4 simple scenarios because the IPv4 exterior routing protocol is handling this function. Also, no new entries in IPv4 routing tables result from the use of 6to4.

The Return Path and Source Address Selection

Packets must flow in both directions to be useful; thus it is essential that IPv6 packets sent use a packet with a 6to4 prefix as a source address when talking to a site with a 6to4 prefix; in other words, the destination must have a 6to4 prefix. In the simple example given above, this is not an issue because both sites have only IPv4 connectivity, so they have 6to4 prefixes for their site to communicate with. DNS lookups for host systems at these sites will return only one IPv6 address, which will be the one with a 6to4 prefix. Source address selection is thus not an issue.

As we will soon see, source address selection is an issue for more complex 6to4 usage scenarios; therefore, some source address selection algorithm is necessary in IPv6 hosts. The exact form and method of the algorithm to use is under active study at the IETF IPv6 (ipng) working group^[10], and an algorithm is likely to be chosen in early 2000. Meanwhile, for the purposes of understanding 6to4, it is sufficient to realize that when a 6to4 connected sending site is sending to a destination site using that site's 6to4 prefix, the sending host must guarantee that the source IPv6 address uses the sending site's 6to4 prefix.

More Complex 6to4 Usage Scenarios

Several more interesting 6to4 usage scenarios exist when a site has both 6to4 connectivity and native IPv6 connectivity. The simplest of these is when such a site is trying to reach another site that has only 6to4 connectivity, in which case the source address selection algorithm mentioned above is essential to ensure that the site's 6to4 IPv6 address is chosen. No destination selection is required because there is only one choice, that is, 6to4.

Similarly, when a site that has only 6to4 connectivity tries to reach a site with both 6to4 and native IPv6 connectivity, some host rule for choosing among multiple destination addresses must result in the 6to4 address being chosen, because only a local 6to4 IPv6 source address is available. Of course source selection is not an issue in this case because there is only the 6to4 IPv6 address to use.

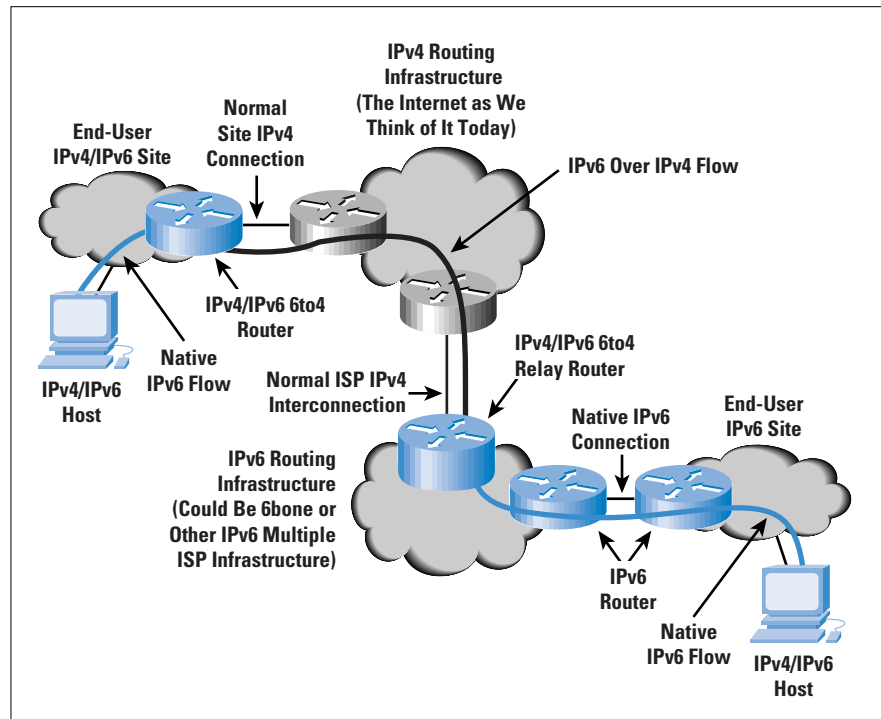
Another variation of these scenarios is when a site with 6to4 and native IPv6 connectivity is trying to reach another site that has only native IPv6 connectivity, making a source address selection algorithm essential to make sure the site's native IPv6 address is chosen. No destination selection is required, because there is only one choice, that is, the native IPv6 address.

Similarly, when a site that has only native IPv6 connectivity tries to reach a site with 6to4 and native IPv6 connectivity, a host rule is essential for choosing among multiple addresses to ensure that a native IPv6 address is chosen, because only a local native IPv6 source address is available. Again, source selection is not an issue in this case because only the native IPv6 address can be used.

An interesting choice develops in the situation when both sites have 6to4 and native IPv6 connectivity as both 6to4-to-6to4 and native IPv6-to-native-IPv6 connections are a possibility. Current thinking as of the writing of this article is to prefer the native IPv6 connection.

The 6to4 Relay

The most interesting, and most complex, 6to4 scenario is that of sites with only 6to4 connectivity communicating with sites with only native IPv6 connectivity. This is accomplished by the use of a 6to4 relay that supports both 6to4 and native IPv6 connectivity (Figure 4). The 6to4 relay is nothing more than an IPv4/IPv6 dual-stack router.

Figure 4: The 6to4
Relay

The 6to4 relay advertises a route to 2002::/16 for itself into the native IPv6 infrastructure it is attached to. The native IPv6 network operators must filter out and discard any 6to4 (2002:...) prefix advertisements longer than /16. In addition, the 6to4 relay may advertise into its 6to4 connection whatever native IPv6 routes its policies allow, which the 6to4 router at the 6to4-only site picks up with either a BGP4+ peering session, or with a default route, to the 6to4 relay.

Thus the 6to4-only site will try to send a packet to the native IPv6-only site by forwarding an encapsulated (tunneled) IPv6 packet to the 6to4 relay, which removes the IPv4 header (decapsulates) and forwards the packet on to the IPv6-only site.

Potentially, multiple 6to4 relays are needed, one for each separate IPv6 routing realm (collection of IPv6 routing ISPs). In practice, it is expected that all native IPv6 ISP services will be interconnected even if the use of inter-IPv6-ISP manually configured tunnels are required to do so. This is currently the case as of early 2000, because all 6bone 3FFE::/16 TLA networks and all production 2001::/16 subTLA networks are interconnected with each other.

It is expected that native IPv6 service providers will choose to operate 6to4 relays as a simple extension of their service. There are no special rules or exceptions to 6to4 as described here for this to happen because the 6to4 relay is simply operated as part of an end-user site that belongs to the IPv6 ISP.

Other Issues

Several other 6to4 issues are presented below for completeness.

- The IPv6 *Maximum Transmission Unit* (MTU) size could prove too large for some intermediate IPv4 link when a 6to4 tunnel is in use, thus IPv4 fragmentation will occur. Though undesirable, fragmentation is not disastrous, so the IPv4 “Do Not Fragment” bit should not be set in the IPv4 packet carrying the 6to4 tunnel.
- How sites move IPv6 packets internal to a site is not important to the 6to4 process. For illustrative purposes in this article, it is generally assumed that native IPv6 transmission exists within a site. This may not be strictly true because “6over4,” manual tunnels, and other methods of moving IPv6 packets could be in use. Nonetheless, it is not important to the 6to4 processes described here.
- Security issues with the 6to4 mechanism are not discussed here. The reader is referred to the current 6to4 draft for an explanation of these issues^[6].
- 6to4 sites with IPv6 connectivity must not inject their 6to4 prefix into the IPv6 routing infrastructure via the native IPv6 connection.
- It is not possible to assume the general availability of wide-area IPv4 multicast, so the 6to4 mechanism must assume only unicast capability in its underlying IPv4 carrier network. However, it is expected that IPv6 multicast packets may be sent to, or sourced from, a 6to4 router in the IPv4 encapsulated form, as described above. When IPv6 multicast is supported, an IPv6 multicast routing protocol must be used.
- The use of IPv6 Anycast is compatible with 6to4 prefixes.
- 6to4 for hosts only, as opposed to sites, is possible and will likely be developed in the future. However, details of this feature are not discussed in this article.
- The 6to4 mechanism is unaffected by the presence of a firewall at the border router.
- When using IPv4 *Network Address Translation* (NAT), 6to4 mechanisms remain valid, and the NAT device includes a fully functional IPv6 router with the 6to4 mechanism included. Combining 6to4 and NAT in this way offers the advantages of NAT for IPv4 use, and the additional address space of IPv6.
- There is no significant impact to either IPv4 or IPv6 routing table size caused by the proper implementation of 6to4.

Summarizing 6to4

The 6to4 mechanism allows isolated IPv6 routing domains to communicate with other IPv6 routing domains, even in the total absence of native IPv6 service providers. It is a powerful IPv6 transition tool that will allow both traditional IPv4-based Internet end-user sites and new IPv6-only Internet sites to utilize IPv6 and operate successfully over the existing IPv4-based Internet routing infrastructure.

For Further Reading

- [0] Fink, R., “IPv6—What and Where It Is,” *The Internet Protocol Journal*, Volume 2, No. 1, March 1999.
- [1] IPng and IPv6 information, including formal specifications, can be found at: <http://playground.sun.com/pub/ipng/html>
- [2] “The Case for IPv6,” an Internet Draft of the IAB, can be found at: <http://www.6bone.net/misc/case-for-ipv6.html>
- [3] IETF IPv6 Transition Working Group (ngtrans) information, including status of all its current projects, can be found at: <http://www.6bone.net/ngtrans/>
- [4] “Transition Mechanisms for IPv6 Hosts and Routers,” RFC 1933, can be found at: <http://www.ietf.org/rfc/rfc1933.txt>
- [5] The 6bone IPv6 Testbed Network is explained at: <http://www.6bone.net>
- [6] “Connection of IPv6 Domains via IPv4 Clouds without Explicit Tunnels” (“6to4”), an Internet Draft of the IETF ngtrans WG, can be found at: <http://www.6bone.net/misc/6to4.txt>
- [7] “IPv6 Aggregatable Global Unicast Address Format,” RFC 2374, can be found at: <http://www.ietf.org/rfc/rfc2374.txt>
- [8] “Transmission of IPv6 Packets over IPv4 Domains without Explicit Tunnels” (“6over4”), RFC 2529, can be found at: <http://www.ietf.org/rfc/rfc2529.txt>
- [9] “Neighbor Discovery for IP Version 6 (IPv6),” RFC 2461, can be found at: <http://www.ietf.org/rfc/rfc2461.txt>
- [10] IETF IPv6 Working Group (ipngwg) information, can be found at: <http://www.ietf.org/html.charters/ipngwg-charter.html>

BRIAN E. CARPENTER is a network researcher with the IBM Internet Division at iCAIR in Evanston Illinois. He is currently the Chair of the Internet Architecture Board (IAB) of the IETF. You can reach him at: brian@icair.org

KEITH MOORE is a network researcher at the Innovative Computing Laboratory of the Computer Science Department at the University of Tennessee. He is currently a Co-Director of the IETF Applications Area in the Internet Engineering Steering Group. You can reach him at: moore@cs.utk.edu

ROBERT FINK is a network researcher with the U.S. Dept. of Energy’s Energy Sciences Network (ESnet) at the Lawrence Berkeley National Laboratory. He is currently a co-chair of the IETF ngtrans (IPng Transition) Working Group, and leads the 6bone project. You can reach him at: fink@es.net

IP Security

by William Stallings

In 1994, the *Internet Architecture Board* (IAB) issued a report entitled “Security in the Internet Architecture” (RFC 1636). The report stated the general consensus that the Internet needs more and better security, and it identified key areas for security mechanisms. Among these were the need to secure the network infrastructure from unauthorized monitoring and control of network traffic and the need to secure end-user-to-end-user traffic using authentication and encryption mechanisms.

These concerns are fully justified. As confirmation, the 1998 annual report from the *Computer Emergency Response Team* (CERT) lists over 1,300 reported security incidents affecting nearly 20,000 sites. The most serious types of attacks included IP spoofing, in which intruders create packets with false IP addresses and exploit applications that use authentication based on IP address; and various forms of eavesdropping and packet sniffing, in which attackers read transmitted information, including logon information and database contents.

In response to these issues, the IAB included authentication and encryption as necessary security features in the next-generation IP, which has been issued as IPv6. Fortunately, these security capabilities were designed to be usable both with the current IP (IPv4) and IPv6, meaning that vendors can begin offering these features now, and many vendors do now have some *IP Security Protocol* (IPSec) capability in their products.

Applications of IPSec

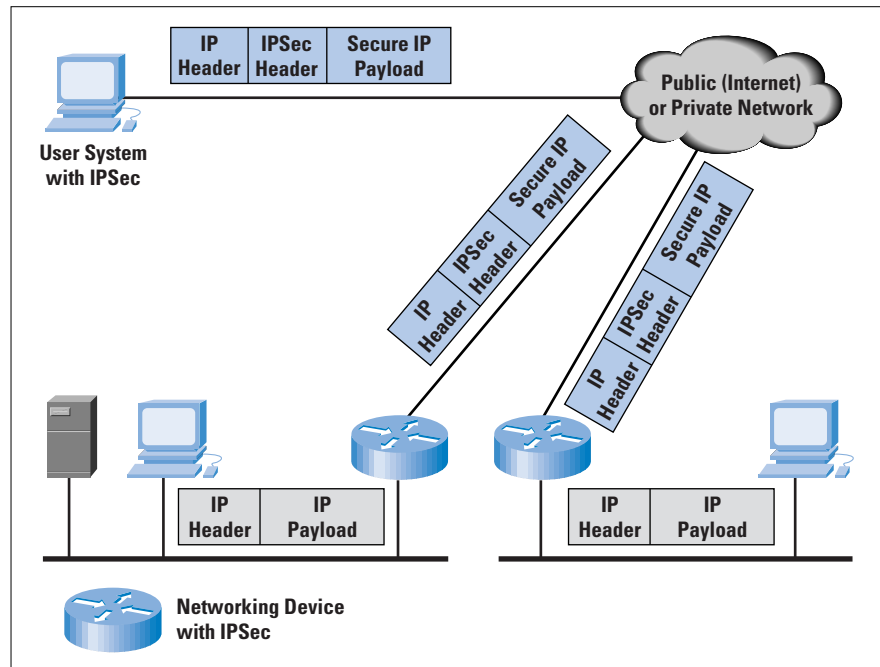
The Internet community has developed application-specific security mechanisms in numerous application areas, including electronic mail (*Privacy Enhanced Mail*, *Pretty Good Privacy* [PGP]), network management (*Simple Network Management Protocol Version 3* [SNMPv3]), Web access (*Secure HTTP*, *Secure Sockets Layer* [SSL]), and others. However, users have some security concerns that cut across protocol layers. For example, an enterprise can run a secure, private TCP/IP network by disallowing links to untrusted sites, encrypting packets that leave the premises, and authenticating packets that enter the premises. By implementing security at the IP level, an organization can ensure secure networking not only for applications that have security mechanisms but also for the many security-ignorant applications.

IPSec provides the capability to secure communications across a LAN, across private and public WANs, and across the Internet. Examples of its use include:

- Secure branch office connectivity over the Internet: A company can build a secure virtual private network over the Internet or over a public WAN. This enables a business to rely heavily on the Internet and reduce its need for private networks, saving costs and network management overhead.
- Secure remote access over the Internet: An end user whose system is equipped with IP security protocols can make a local call to an *Internet Service Provider* (ISP) and gain secure access to a company network. This reduces the cost of toll charges for traveling employees and telecommuters.
- Establishment of extranet and intranet connectivity with partners: IPSec can be used to secure communication with other organizations, ensuring authentication and confidentiality and providing a key exchange mechanism.
- Enhancement of electronic commerce security: Most efforts to date to secure electronic commerce on the Internet have relied upon securing Web traffic with SSL since that is commonly found in Web browsers and is easy to set up and run. There are new proposals that may utilize IPSec for electronic commerce.

The principal feature of IPSec that enables it to support these varied applications is that it can encrypt or authenticate *all* traffic at the IP level. Thus, all distributed applications, including remote logon, client/server, e-mail, file transfer, Web access, and so on, can be secured. Figure 1 shows a typical scenario of IPSec usage. An organization maintains LANs at dispersed locations. Traffic on each LAN does not need any special protection, but the devices on the LAN can be protected from the untrusted network with firewalls. Since we live in a distributed and mobile world, the people who need to access the services on each of the LANs may be at sites across the Internet. These people can use IPSec protocols to protect their access. These protocols can operate in networking devices, such as a router or firewall that connects each LAN to the outside world, or they may operate directly on the workstation or server. In the diagram, the user workstation can establish an IPSec tunnel with the network devices to protect all the subsequent sessions. After this tunnel is established, the workstation can have many different sessions with the devices behind these IPSec gateways. The packets going across the Internet will be protected by IPSec but will be delivered onto each LAN as a normal IP packet.

Figure 1: An IP Security Scenario



Benefits of IPSec

The benefits of IPSec include:

- When IPSec is implemented in a firewall or router, it provides strong security that can be applied to all traffic crossing the perimeter. Traffic within a company or workgroup does not incur the overhead of security-related processing.
- IPSec is below the transport layer (TCP, UDP), so is transparent to applications. There is no need to change software on a user or server system when IPSec is implemented in the firewall or router. Even if IPSec is implemented in end systems, upper layer software, including applications, is not affected.
- IPSec can be transparent to end users. There is no need to train users on security mechanisms, issue keying material on a per-user basis, or revoke keying material when users leave the organization.
- IPSec can provide security for individual users if needed. This feature is useful for offsite workers and also for setting up a secure virtual subnetwork within an organization for sensitive applications.

Is IPSec the Right Choice?

There are already numerous products that implement IPSec, but it is not necessarily the security solution of choice for a network administrator. Christian Huitema, who at the time of the development of the initial IP-Sec documents was the head of the IAB, reports that the debates over how to provide Internet-based security were among the most heated that he ever observed. One issue concerns whether security is being provided at the right protocol layer. To provide security at the IP level, it is necessary for IPSec to be a part of the network code deployed on all participating platforms, including Windows NT, UNIX, and Macintosh systems. Unless a desired feature is available on all the deployed platforms, a given application may not be able to use that feature.

On the other hand, if the application, such as a Web browser/server combination, incorporates the function, the developer can guarantee that the features are available on all platforms for which the application is available. A related point is that many Internet applications are now being released with embedded security features. For example, Netscape and Internet Explorer support SSL, which protects Web traffic. Also, many vendors are planning to support *Secure Electronic Transaction* (SET), which protects credit-card transactions over the Internet. However, for a virtual private network, a network-level facility is needed, and this is what IPSec provides.

The Scope of IPSec

IPSec provides three main facilities: an authentication-only function, referred to as *Authentication Header* (AH), a combined authentication/encryption function called *Encapsulating Security Payload* (ESP), and a key exchange function. For virtual private networks, both authentication and encryption are generally desired, because it is important both to (1) assure that unauthorized users do not penetrate the virtual private network and (2) assure that eavesdroppers on the Internet cannot read messages sent over the virtual private network. Because both features are generally desirable, most implementations are likely to use ESP rather than AH. The key exchange function allows for manual exchange of keys as well as an automated scheme.

The IPSec specification is quite complex and covers numerous documents. The most important of these, issued in November 1998, are RFCs 2401, 2402, 2406, and 2408.

Security Associations

A key concept that appears in both the authentication and confidentiality mechanisms for IP is the *Security Association* (SA). An association is a one-way relationship between a sender and a receiver that affords security services to the traffic carried on it. If a peer relationship is needed, for two-way secure exchange, then two security associations are required. Security services are afforded to an SA for the use of AH or ESP, but not both. A security association is uniquely identified by three parameters:

- *Security Parameters Index* (SPI): The SPI assigns a bit string to this SA that has local significance only. The SPI is carried in AH and ESP headers to enable the receiving system to select the SA under which a received packet will be processed.
- *IP destination address*: Currently, only unicast addresses are allowed; this is the address of the destination endpoint of the SA, which may be an end-user system or a network system such as a firewall or router.
- *Security protocol identifier*: This indicates whether the association is an AH or ESP security association.

Hence, in any IP packet, the security association is uniquely identified by the destination address in the IPv4 or IPv6 header and the SPI in the enclosed extension header (AH or ESP).

An IPSec implementation includes a security association database that defines the parameters associated with each SA. A security association is defined by the following parameters:

- *Sequence number counter*: A 32-bit value used to generate the sequence number field in AH or ESP headers
- *Sequence counter overflow*: A flag indicating whether overflow of the sequence number counter should generate an auditable event and prevent further transmission of packets on this SA
- *Anti-replay window*: Used to determine whether an inbound AH or ESP packet is a replay, by defining a sliding window within which the sequence number must fall
- *AH information*: Authentication algorithm, keys, key lifetimes, and related parameters being used with AH
- *ESP information*: Encryption and authentication algorithm, keys, initialization values, key lifetimes, and related parameters being used with ESP
- *Lifetime of this security association*: A time interval or byte count after which an SA must be replaced with a new SA (and new SPI) or terminated, plus an indication of which of these actions should occur
- *IPSec protocol mode*: Tunnel, transport, or wildcard (required for all implementations); these modes are discussed later
- *Path MTU*: Any observed path maximum transmission unit (maximum size of a packet that can be transmitted without fragmentation) and aging variables (required for all implementations)

The key management mechanism that is used to distribute keys is coupled to the authentication and privacy mechanisms only by way of the security parameters index. Hence, authentication and privacy have been specified independent of any specific key management mechanism.

SA Selectors

IPSec provides the user with considerable flexibility in the way in which IPSec services are applied to IP traffic. IPSec provides a high degree of granularity in discriminating between traffic that is afforded IPSec protection and traffic that is allowed to bypass IPSec, in the former case relating IP traffic to specific SAs.

The means by which IP traffic is related to specific SAs (or no SA in the case of traffic allowed to bypass IPSec) is the nominal *Security Policy Database* (SPD). In its simplest form, an SPD contains entries, each of which defines a subset of IP traffic and points to an SA for that traffic. In more complex environments, there may be multiple entries that potentially relate to a single SA or multiple SAs associated with a single SPD entry.

Each SPD entry is defined by a set of IP and upper-layer protocol field values, called *selectors*. In effect, these selectors are used to filter outgoing traffic in order to map it into a particular SA. Outbound processing obeys the following general sequence for each IP packet:

- Compare the values of the appropriate fields in the packet (the selector fields) against the SPD to find a matching SPD entry, which will point to zero or more SAs.
- Determine the SA (if any) for this packet and its associated SPI.
- Do the required IPsec processing (that is, AH or ESP processing).

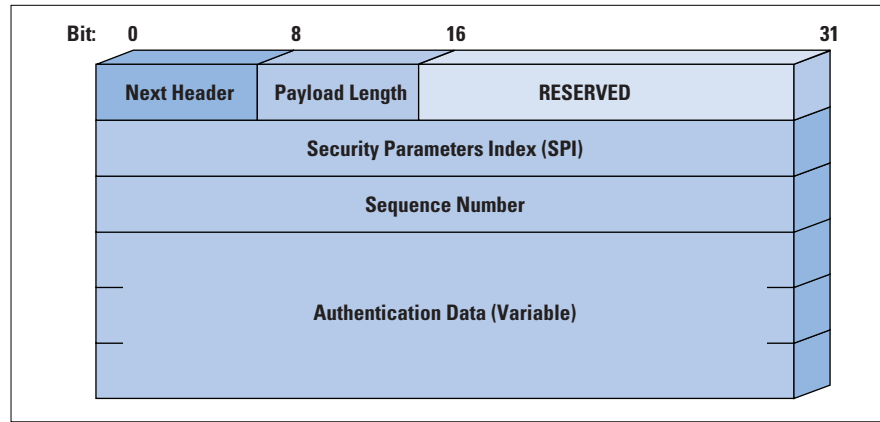
The following selectors determine an SPD entry:

- *Destination IP address*: This may be a single IP address, an enumerated list or range of addresses, or a wildcard (mask) address. The latter two are required to support more than one destination system sharing the same SA (for instance, behind a firewall).
- *Source IP address*: This may be a single IP address, an enumerated list or range of addresses, or a wildcard (mask) address. The latter two are required to support more than one source system sharing the same SA (for instance, behind a firewall).
- *UserID*: UserID is used to identify a policy tied to a valid user or system name.
- *Data sensitivity level*: The data sensitivity level is used for systems providing information flow security (for instance, “Secret” or “Unclassified”).
- *Transport Layer protocol*: This value is obtained from the IPv4 protocol or IPv6 *Next Header* field. This may be an individual protocol number, a list of protocol numbers, or a range of protocol numbers.
- *IPsec protocol* (AH or ESP or AH/ESP): If present, this is obtained from the IPv4 Protocol or IPv6 Next Header field.
- *Source and destination ports*: These may be individual TCP or *User Datagram Protocol* (UDP) port values, an enumerated list of ports, or a wildcard port.
- *IPv6 class*: This class is obtained from the IPv6 header. It may be a specific IPv6 Class value or a wildcard value.
- *IPv6 flow label*: This label is obtained from the IPv6 header. It may be a specific IPv6 flow label value or a wildcard value.
- *IPv4 Type of Service* (TOS): The TOS is obtained from the IPv4 header. It may be a specific IPv4 TOS value or a wildcard value.

Authentication Header

The authentication header provides support for data integrity and authentication of IP packets. The data integrity feature ensures that undetected modification to the content of a packet in transit is not possible. The authentication feature enables an end system or network device to authenticate the user or application and filter traffic accordingly; it also prevents the address spoofing attacks observed in today’s Internet. The AH also guards against the replay attack described later.

Figure 2: IPSec Authentication Header



Authentication is based on the use of a *Message Authentication Code* (MAC); hence the two parties must share a secret key. The authentication header consists of the following fields (Figure 2):

- *Next Header* (8 bits): This field identifies the type of header immediately following this header.
- *Payload Length* (8 bits): This field gives the length of the authentication header in 32-bit words, minus 2. For example, the default length of the authentication data field is 96 bits, or three 32-bit words. With a three-word fixed header, there are a total of six words in the header, and the Payload Length field has a value of 4.
- *Reserved* (16 bits): This field is reserved for future use.
- *Security Parameters Index* (32 bits): This field identifies a security association.
- *Sequence Number* (32 bits): This field contains a monotonically increasing counter value.
- *Authentication Data* (variable): This variable-length field (must be an integral number of 32-bit words) contains the *Integrity Check Value* (ICV), or MAC, for this packet.

Anti-Replay Service

A replay attack is one in which an attacker obtains a copy of an authenticated packet and later transmits it to the intended destination. The receipt of duplicate, authenticated IP packets may disrupt service in some way or may have some other undesired consequence. The *Sequence Number* field is designed to thwart such attacks.

When a new SA is established, the *sender* initializes a sequence number counter to 0. Each time that a packet is sent on this SA, the sender increments the counter and places the value in the Sequence Number field. Thus, the first value to be used is 1. If anti-replay is enabled (the default), the sender must not allow the sequence number to cycle past $2^{32} - 1$ back to zero. Otherwise, there would be multiple valid packets with the same sequence number. If the limit of $2^{32} - 1$ is reached, the sender should terminate this SA, and negotiate a new SA with a new key.

Because IP is a connectionless, unreliable service, the protocol does not guarantee that packets will be delivered in order and does not guarantee that all packets will be delivered. Therefore, the IPSec authentication document dictates that the *receiver* should implement a window of size W , with a default of $W = 64$. The right edge of the window represents the highest sequence number, N , so far received for a valid packet. For any packet with a sequence number in the range from $N - W + 1$ to N that has been correctly received (that is, properly authenticated), the corresponding slot in the window is marked. Inbound processing proceeds as follows when a packet is received:

- If the received packet falls within the window and is new, the MAC is checked. If the packet is authenticated, the corresponding slot in the window is marked.
- If the received packet is to the right of the window and is new, the MAC is checked. If the packet is authenticated, the window is advanced so that this sequence number is the right edge of the window, and the corresponding slot in the window is marked.
- If the received packet is to the left of the window, or if authentication fails, the packet is discarded; this is an auditable event.

Message Authentication Code

The message authentication algorithm is used to calculate a message authentication code, using an algorithm known as *HMAC*. HMAC takes as input a portion of the message and a secret key and produces a MAC as output. This MAC value is stored in the Authentication Data field of the AH header. The calculation takes place over the entire enclosed TCP segment plus the authentication header. When this IP packet is received at the destination, the same calculation is performed using the same key. If the calculated MAC equals the value of the received MAC, then the packet is assumed to be authentic. The authentication data field is calculated over:

- IP header fields that either do not change in transit (immutable) or that are predictable in value upon arrival at the endpoint for the AH SA. Fields that may change in transit and whose value on arrival are unpredictable are set to zero for purposes of calculation at both source and destination.
- The AH header other than the Authentication Data field. The Authentication Data field is set to zero for purposes of calculation at both source and destination.
- The entire upper-level protocol data, which is assumed to be immutable in transit (for instance, a TCP segment or an inner IP packet in tunnel mode).

For IPv4, examples of immutable fields are *Internet Header Length* and *Source Address*. An example of a mutable but predictable field is the *Destination Address* (with loose or strict source routing). Examples of mutable fields that are zeroed prior to ICV calculation are the *Time to Live* (TTL) and *Header Checksum* fields.

Note that both source and destination address fields are protected, so that address spoofing is prevented. For IPv6, examples in the base header are *Version* (immutable), *Destination Address* (mutable but predictable), and *Flow Label* (mutable and zeroed for calculation).

Encapsulating Security Payload

The encapsulating security payload provides confidentiality services, including confidentiality of message contents and limited traffic flow confidentiality. As an optional feature, ESP can also provide the same authentication services as AH.

Figure 3: IPSec ESP Format

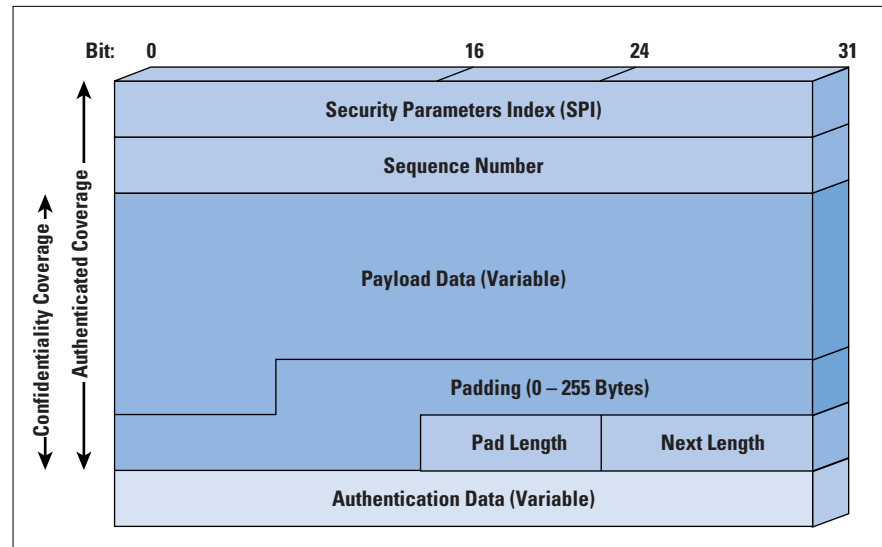


Figure 3 shows the format of an ESP packet. It contains the following fields:

- *Security Parameters Index* (32 bits): Identifies a security association
- *Sequence Number* (32 bits): A monotonically increasing counter value
- *Payload Data* (variable): A transport-level segment (transport mode) or IP packet (tunnel mode) that is protected by encryption
- *Padding* (0–255 bytes): Extra bytes that may be required if the encryption algorithm requires the plaintext to be a multiple of some number of octets
- *Pad Length* (8 bits): Indicates the number of pad bytes immediately preceding this field
- *Next Header* (8 bits): Identifies the type of data contained in the payload data field by identifying the first header in that payload (for example, an extension header in IPv6, or an upper-layer protocol such as TCP)
- *Authentication Data* (variable): A variable-length field (must be an integral number of 32-bit words) that contains the integrity check value computed over the ESP packet minus the Authentication Data field

Encryption and Authentication Algorithms

The Payload Data, Padding, Pad Length, and Next Header fields are encrypted by the ESP service. If the algorithm used to encrypt the payload requires cryptographic synchronization data, such as an *Initialization Vector* (IV), then this data may be carried explicitly at the beginning of the Payload Data field. If included, an IV is usually not encrypted, although it is often referred to as being part of the ciphertext. The current specification dictates that a compliant implementation must support the *Data Encryption Standard* (DES). A number of other algorithms have been assigned identifiers and could, therefore, be used for encryption; these include:

- Three-key triple DES
- RC5
- International Data Encryption Algorithm (IDEA)
- Three-key triple IDEA
- CAST
- Blowfish

It is now well known that DES is inadequate for secure encryption, so it is likely that many future implementations will use triple DES and eventually the *Advanced Encryption Standard* (AES). As with AH, ESP supports the use of a MAC, using HMAC.

Padding

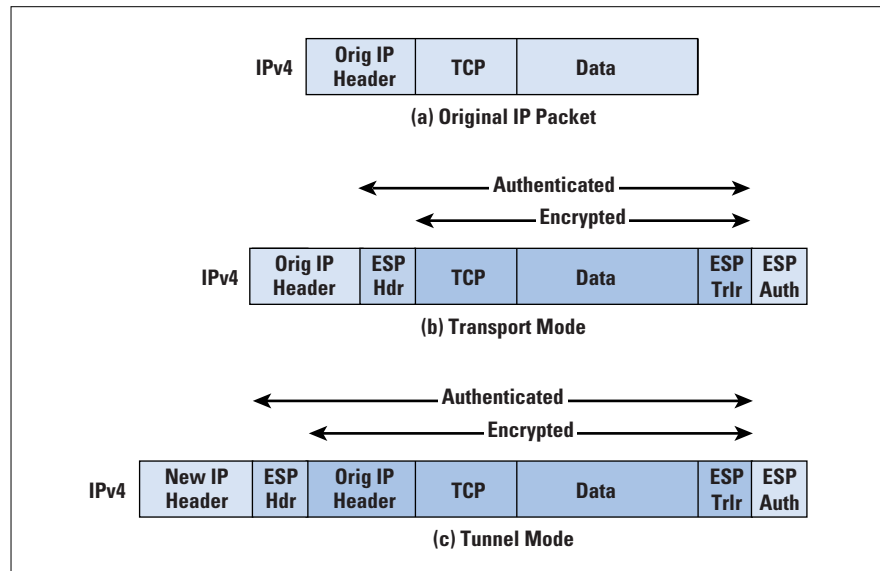
- The Padding field serves several purposes: If an encryption algorithm requires the plaintext to be a multiple of some number of bytes (for instance, the multiple of a single block for a block cipher), the Padding field is used to expand the plaintext (consisting of the Payload Data, Padding, Pad Length, and Next Header fields) to the required length.
- The ESP format requires that the Pad Length and Next Header fields be right aligned within a 32-bit word. Equivalently, the ciphertext must be an integer multiple of 32 bits. The Padding field is used to assure this alignment.
- Additional padding may be added to provide partial traffic flow confidentiality by concealing the actual length of the payload.

Figure 4 indicates the scope of ESP encryption and authentication in both transport and tunnel modes.

Transport and Tunnel Modes

Both AH and ESP support two modes of use: *transport* and *tunnel* mode.

Figure 4: Scope of ESP
Encryption and
Authentication



Transport Mode

Transport mode provides protection primarily for upper-layer protocols. That is, transport mode protection extends to the payload of an IP packet. Examples include a TCP or UDP segment, or an *Internet Control Message Protocol* (ICMP) packet, all of which operate directly above IP in a host protocol stack. For this mode using IPv4, the ESP header is inserted into the IP packet immediately prior to the transport-layer header (for instance, TCP, UDP, ICMP) and an ESP trailer (Padding, Pad Length, and Next Header fields) is placed after the IP packet. This setup is shown in Figure 4b. If authentication is selected, the ESP Authentication Data field is added after the ESP trailer. The entire transport-level segment plus the ESP trailer are encrypted. Authentication covers all of the ciphertext plus the ESP header.

Typically, transport mode is used for end-to-end communication between two hosts (for instance, communications between a workstation and a server, or two servers). When a host runs AH or ESP over IPv4, the payload is the data that normally follows the IP header. For IPv6, the payload is the data that normally follows both the IP header and any IPv6 extensions headers that are present, with the possible exception of the destination options header, which may be included in the protection.

ESP in transport mode encrypts and optionally authenticates the IP payload but not the IP header. AH in transport mode authenticates the IP payload and selected portions of the IP header. All IPv4 packets have a *Next Header* field. This field contains a number for the payload protocol, such as 6 for TCP and 17 for UDP. For transport mode, the IP Next Header field is decimal 51 for AH, or 50 for ESP. This tells the receiving machine to interpret the remainder of the packet after the IP header as either AH or ESP. Both the AH and ESP headers also have a Next Header field.

As an example, let's examine a Telnet session within an ESP packet in transport mode. The IP header would contain 51 in the Next Header field. In the ESP header, the Next Header field would be 6 for TCP. Within the TCP header, Telnet would be identified as port 23.

Transport mode operation may be summarized for ESP as follows:

- At the source, the block of data consisting of the ESP trailer plus the entire transport-layer segment is encrypted and the plaintext of this block is replaced with its ciphertext to form the IP packet for transmission. Authentication is added if this option is selected.
- The packet is then routed to the destination. Each intermediate router needs to examine and process the IP header plus any plaintext IP extension headers but will not need to examine the ciphertext.
- The destination node examines and processes the IP header plus any plaintext IP extension headers. Then, on the basis of the SPI in the ESP header, the destination node decrypts the remainder of the packet to recover the plaintext transport-layer segment. This process is similar for AH, however the payload (upper layer protocol) is not encrypted.

Transport mode operation provides confidentiality for any application that uses it, thus avoiding the need to implement confidentiality in every individual application. This mode of operation is also reasonably efficient, adding little to the total length of the IP packet. One drawback to this mode is that it is possible to do traffic analysis on the transmitted packets.

Tunnel Mode

Tunnel mode encapsulates an entire IP packet within an IP packet to ensure that no part of the original packet is changed as it is moved through a network. The entire original, or inner, packet travels through a "tunnel" from one point of an IP network to another; no routers along the way need to examine the inner IP header. For ESP, this is shown in Figure 4c. Because the IP header contains the destination address and possibly source routing directives and hop-by-hop option information, it is not possible simply to transmit the encrypted IP packet prefixed by the ESP header. Intermediate routers would be unable to process such a packet. Therefore, it is necessary to encapsulate the entire block (ESP header plus ciphertext plus Authentication Data, if present) with a new IP header that will contain sufficient information for routing but not for traffic analysis. Tunnel mode is used when one or both ends of an SA is a security gateway, such as a firewall or router that implements IPSec. With tunnel mode, a number of hosts on networks behind firewalls may engage in secure communications without implementing IPSec. The unprotected packets generated by such hosts are tunneled through external networks by tunnel mode SAs set up by the IPSec process in the firewall or secure router at the boundary of the local network.

Whereas the transport mode is suitable for protecting connections between hosts that support the ESP feature, the tunnel mode is useful in a configuration that includes a firewall or other sort of security gateway that protects a trusted network from external networks. In this latter case, encryption occurs only between an external host and the security gateway or between two security gateways. This setup relieves hosts on the internal network of the processing burden of encryption and simplifies the key distribution task by reducing the number of needed keys. Further, it thwarts traffic analysis based on ultimate destination.

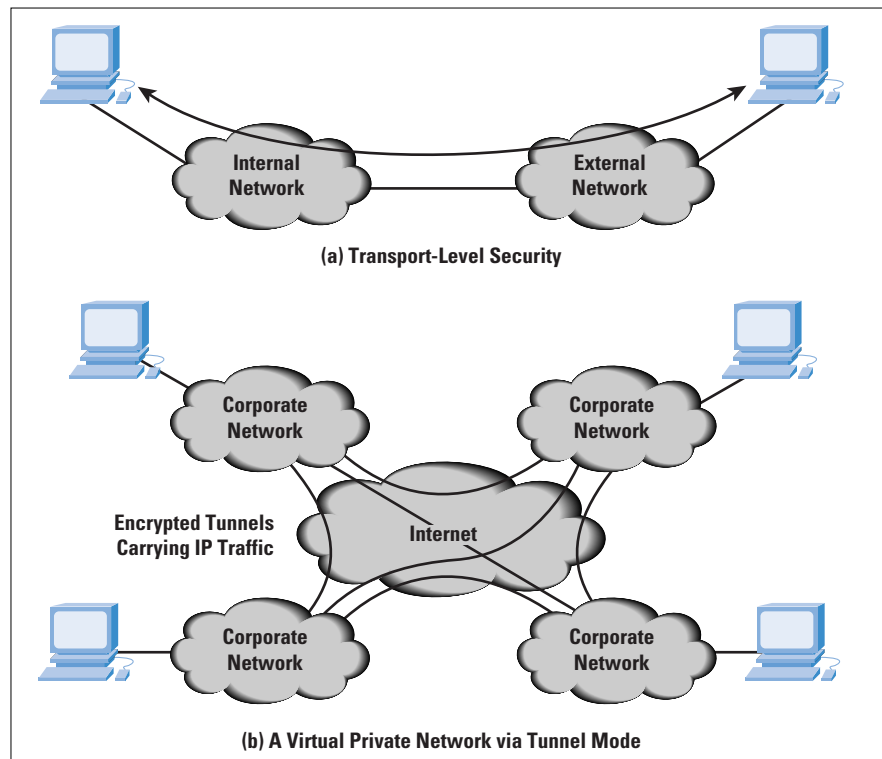
Let's use the diagram in Figure 1 as an example of how tunnel mode IP-Sec operates. The following steps occur for transfer of a transport-layer segment from the user system to one of the servers on one of the protected LANs.

- The user system prepares an inner IP packet with a destination address of the target host on the internal LAN. For a Telnet session, this packet would be a TCP packet with the original SYN flag set with a destination port set to 23. This entire IP packet is prefixed by an ESP header; then the packet and ESP trailer are encrypted and Authentication Data may be added. The Next Header field of the ESP header would be decimal 4 for IP-in-IP, indicating that the entire original IP packet is contained as the "payload." The resulting block is encapsulated with a new IP header (base header plus optional extensions such as routing and hop-by-hop options for IPv6) whose destination address is the firewall; this forms the outer IP packet. The Next Header field for this IP packet is 50 for ESP.
- The outer packet is routed to the destination firewall. Each intermediate router needs to examine and process the outer IP header plus any outer IP extension headers but does not need to examine the ciphertext.
- The destination firewall examines and processes the outer IP header plus any outer IP extension headers. Then, on the basis of the SPI in the ESP header, the gateway decrypts the remainder of the packet to recover the plaintext inner IP packet. This packet is then transmitted in the internal network.
- The inner packet is routed through zero or more routers in the internal network to the destination host. The receiver would have no indication that the packet had been encapsulated and protected by the "tunnel" between the user system and the gateway. It would see the packet as a request to start a Telnet session and would respond back with a TCP SYN/ACK, which would go back to the gateway. The gateway would encapsulate that packet into an IPSec packet and transport it back to the user system through this "tunnel." That return packet would be processed to find the original packet, which would contain the SYN/ACK for the Telnet session.

Common Uses of IPSec in Real Networks

Figure 5 shows two ways in which the IPSec ESP service can be used. In the upper part of the figure, encryption (and optionally authentication) is provided directly between two hosts. Figure 5b shows how tunnel mode operation can be used to set up a *Virtual Private Network* (VPN). In this example, an organization has four private networks interconnected across the Internet. Hosts on the internal networks use the Internet for transport of data but do not interact with other Internet-based hosts. By terminating the tunnels at the security gateway to each internal network, the configuration allows the hosts to avoid implementing the security capability. The former technique is supported by a transport mode SA, while the latter technique uses a tunnel mode SA.

Figure 5: Transport-Mode versus Tunnel-Mode Encryption



Key Management

The key management portion of IPSec involves the determination and distribution of secret keys. The IPSec Architecture document mandates support for two types of key management:

- *Manual:* A system administrator manually configures each system with its own keys and with the keys of other communicating systems. This is practical for small, relatively static environments.
- *Automated:* An automated system enables the on-demand creation of keys for SAs and facilitates the use of keys in a large distributed system with an evolving configuration. An automated system is the most flexible but requires more effort to configure and requires more software, so smaller installations are likely to opt for manual key management.

The default automated key management protocol for IPSec is referred to as *Internet Key Exchange* (IKE). IKE provides a standardized method for dynamically authenticating IPSec peers, negotiating security services, and generating shared keys. IKE has evolved from many different protocols and can be thought of as having two distinct capabilities. One of these capabilities is based on the *Internet Security Association and Key Management Protocol* (ISAKMP). ISAKMP provides a framework for Internet key management and provides the specific protocol support, including formats, for negotiation of security attributes. ISAKMP by itself does not dictate a specific key exchange algorithm; rather, ISAKMP consists of a set of message types that enable the use of a variety of key exchange algorithms. The actual key exchange mechanism in IKE is derived from Oakley and several other key exchange protocols that had been proposed for IPSec. Key exchange is based on the use of the Diffie-Hellman algorithm, but provides added security. In particular, Diffie-Hellman alone does not authenticate the two users that are exchanging keys, making the protocol vulnerable to impersonation. IKE includes mechanisms to authenticate the users.

Public Key Certificates

An important element of IPSec key management is the use of public key certificates. In essence, a public key certificate is provided by a trusted *Certificate Authority* (CA) to authenticate a user's public key. The essential elements include:

- Client software creates a pair of keys, one public and one private. The client prepares an unsigned certificate that includes a user ID and the user's public key. The client then sends the unsigned certificate to a CA in a secure manner.
- A CA creates a signature by calculating the hash code of the unsigned certificate and encrypting the hash code with the CA's private key; the encrypted hash code is the signature. The CA attaches the signature to the unsigned certificate and returns the now signed certificate to the client.
- The client may send its signed certificate to any other user. That user may verify that the certificate is valid by calculating the hash code of the certificate (not including the signature), decrypting the signature using the CA's public key, and comparing the hash code to the decrypted signature.

If all users subscribe to the same CA, then there is a common trust of that CA. All user certificates can be placed in the directory for access by all users. In addition, a user can transmit his or her certificate directly to other users. In either case, once B is in possession of A's certificate, B has confidence that messages it encrypts with A's public key will be secure from eavesdropping and that messages signed with A's private key are unforgeable.

If there is a large community of users, it may not be practical for all users to subscribe to the same CA. Because it is the CA that signs certificates, each participating user must have a copy of the CA's own public key to verify signatures. This public key must be provided to each user in an absolutely secure (with respect to integrity and authenticity) way so that the user has confidence in the associated certificates. Thus, with many users, it may be more practical for there to be many CAs, each of which securely provides its public key to some fraction of the users. In practice, there is not a single CA but rather a hierarchy of CAs. This complicates the problems of key distribution and of trust, but the basic principles are the same.

Whither IP Security

The driving force for the acceptance and deployment of secure IP is the need for business and government users to connect their private WAN/LAN infrastructure to the Internet for (1) access to Internet services and (2) use of the Internet as a component of the WAN transport system. Users need to isolate their networks and at the same time send and receive traffic over the Internet. The authentication and privacy mechanisms of secure IP provide the basis for a security strategy.

Because IP security mechanisms have been defined independent of their use with either the current IP or IPv6, deployment of these mechanisms does not depend on deployment of IPv6. Indeed, it is likely that we will see widespread use of secure IP features long before IPv6 becomes popular.

Recommended Web Sites

- The IPsec Working Group of the IETF. Charter for the group and latest RFCs and Internet Drafts for IPsec:
<http://ietf.org/html.charters/ipsec-charter.html>
- IPsec Resources: List of companies implementing IPsec, implementation survey, and other useful material:
<http://web.mit.edu/tytso/www/ipsec/index.html>

WILLIAM STALLINGS is a consultant, lecturer, and author of over a dozen books on data communications and computer networking. He has a Ph.D. in computer science from M.I.T. His latest book is *Local and Metropolitan Area Networks, Sixth Edition* (Prentice Hall, 2000). His home in cyberspace is WilliamStallings.com and he can be reached at ws@shore.net

Quality of Service—Fact or Fiction?

by Geoff Huston, Telstra

Much has been written about the potential of *Quality of Service* (QoS) and the Internet. However, much of the material is strong on promise, but falls short in critical analysis. In an effort to balance the picture, we present here a brief status report on the QoS effort, exposing some of the weaknesses in the current QoS architectures.

The QoS Service

The default service offering associated with the Internet is a *best-effort* service, where the network treats all traffic in exactly the same way. There is no consistent service outcome from the Internet best-effort service model. When the load level is low, the network delivers a high-quality service. The best-effort Internet does not deny entry to traffic, so as the load levels increase, the network congestion levels increase, and service-quality levels decline uniformly. This decline in service is experienced by all traffic passing through a congestion point, and is not limited to the most recently admitted traffic flows. For many applications, this best-effort response is perfectly acceptable. When network capacity is available, the application can make use of the resource, whereas when the level of contention for network bandwidth is high, each application will experience similar levels of congestion. A best-effort network service is a good match to opportunistic applications that can vary their data transfer rate in response to signaled network load.

The objective of various Internet QoS efforts is to augment this service with a number of selectable service responses. These service responses may be different from the best-effort service by some form of superior service response, such as lower delay, lower jitter, or greater bandwidth. These responses are relative, where the service outcome is claimed to be no worse than best effort at any time, and superior to best-effort under congestion load. Alternatively, QoS service responses may be distinguished by providing a consistent, and therefore predictable, service response that is unaffected by network congestion levels. These are quantitative service responses, where the characteristics of the service can be measured against a constant outcome. A quantitative service may be one that constrains jitter to a maximum level, or one that makes a certain bandwidth available, within parameters of bounded jitter, similar to a conventional leased line. Such constant-rate services may be superior to best-effort services when the network is under load, but they may also offer inferior service when the network is under negligible load. The essential attribute of these services is one of consistency.

Why is there a need for relative or consistent service profiles within the Internet? The underlying reasons for introducing QoS into the Internet appear to be threefold: First is the desire to provide high-quality support for IP voice and video services, second is the desire to manage the ser-

vice response provided to low-speed access devices, such as Internet mobile wireless devices, and third is the desire to provide a differentiated Internet access service, providing a network client with a range of service-quality levels at a range of prices.

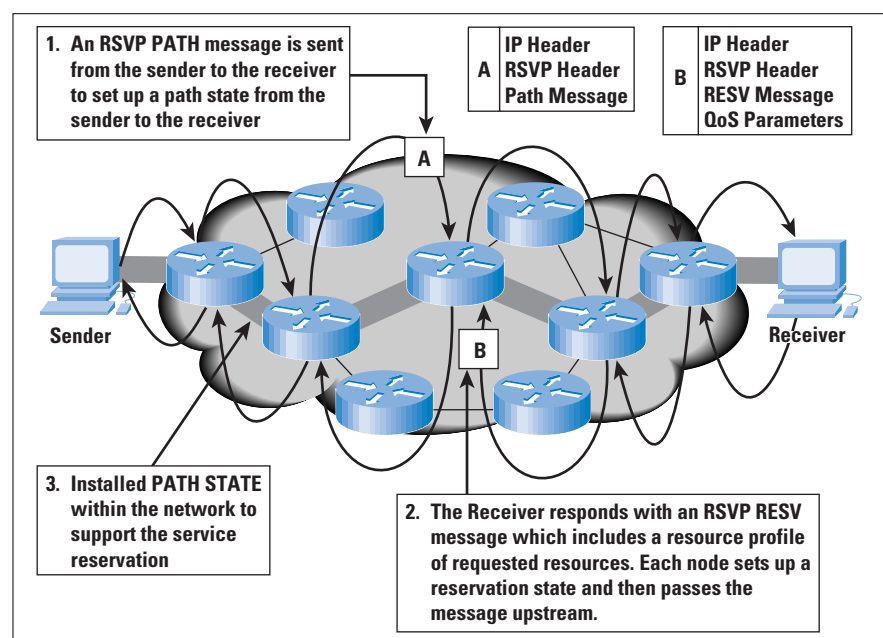
Obviously this is a broad agenda, where there are requirements to extend specific network services to applications, requirements to adapt network services to particular transmission characteristics, and requirements to manage network resources to achieve particular response characteristics for an aggregated collection of traffic.

Approaches to QoS

The relevant efforts within the *Internet Engineering Task Force* (IETF) have been addressing standards for QoS mechanisms within the network.

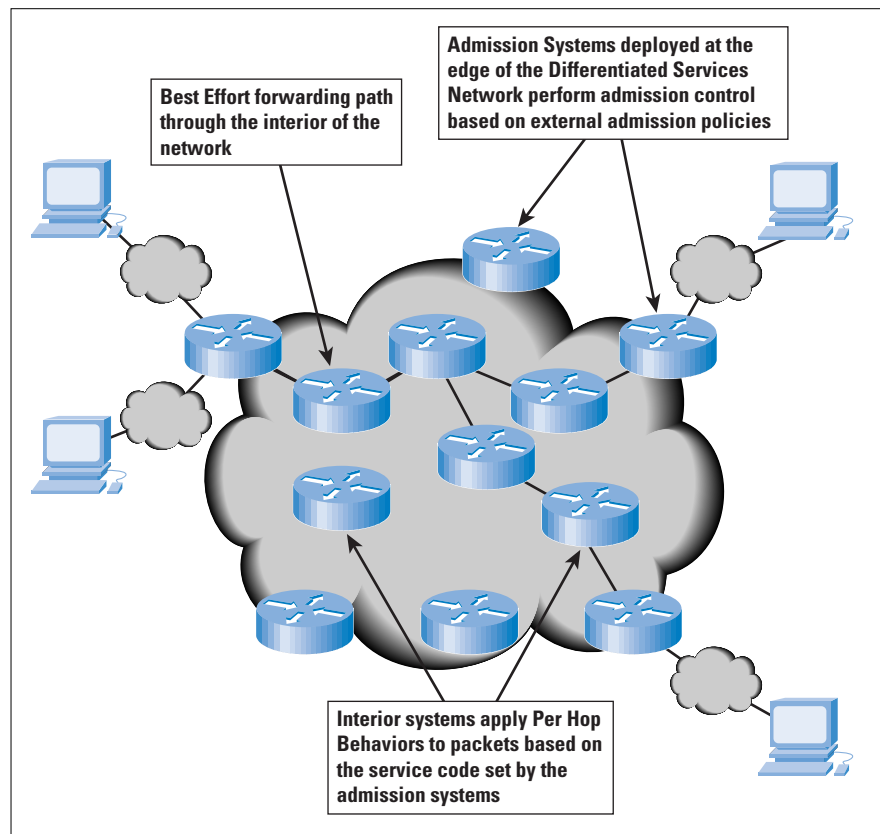
The initial approach to QoS was that of the *Integrated Services* architecture. This approach focuses on the application as the trigger for QoS. Here, the application first signals its service requirements to the network in the form of a reservation, and the network responds to this request. The application proceeds only if the network has indicated that it is able to carry the additional load at the requested service level by committing to the reservation. The reservation remains in force until the application explicitly requests termination of the reservation, or the network signals to the application that it is unable to continue the reservation. The essential feature of this model is the “all-or-nothing” nature of the service model. Either the network commits to the reservation, in which case the application does not have to monitor the level of network response to the service, or the network indicates that it cannot meet the reservation. This approach imposes per-application state within the network, and for large-scale networks, such as the global Internet itself, this approach alone does not appear to be viable (see Figure 1).

Figure 1: The Integrated Services QoS Architecture



The subsequent approach to QoS mechanisms has been to look at the core of the network, and examine those mechanisms that can provide differentiated service outcomes with appropriate scaling properties. This approach, the *Differentiated Services* architecture, includes dropping the concept of a per-application path state across the network using instead the concept of aggregated service mechanisms. Within the aggregated service model, the network provides a smaller number of different service classes and aggregates similar service demands from a set of applications into a single service class. Aggregated services are typically seen as an entry filter, where on entry to the network each packet is classified into a particular service profile. This classification is carried within the IP packet header, using 6 bits from the deprecated IP *Type of Service* (TOS) header to carry the service coding. The network then uses this service code in the packet header to treat this packet identically to all other packets within the same service code. While this approach does possess the ability to scale across the entire Internet, there are numerous unresolved issues relating to the quality signaling between individual applications and the network. The aggregated service model does not allow an individual application to sense if it is receiving the necessary service response from the network (see Figure 2).

Figure 2:
The Differentiated
Services QoS
Architecture



QoS Deployment

Neither approach alone is adequate to meet the QoS requirements. The Integrated Services approach alone imposes an excessive load in the core of large networks through the imposition of a per-application path state. The Differentiated Services approach does provide superior scaling properties through the use of aggregated service elements, but includes no concept of control signaling to inform the traffic conditioning elements of the current state of the network, or the current per-application requirements.

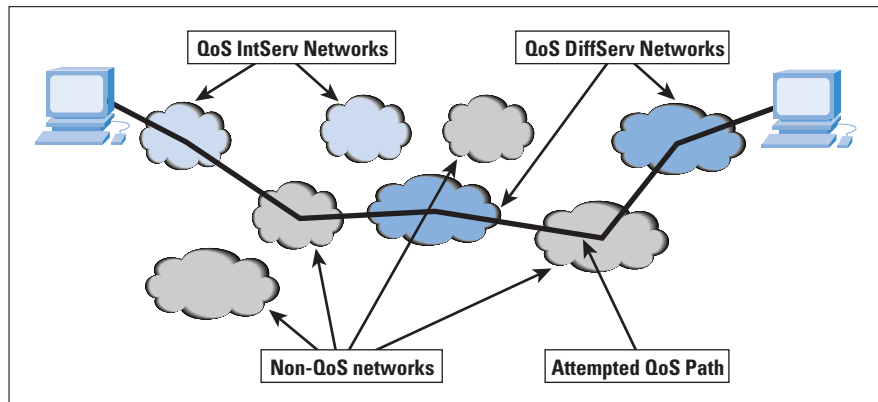
The underlying question then becomes: Is a combination of these two approaches sufficient to allow QoS to be widely deployed on the Internet?

At this stage the response does appear to be a “No.” Perhaps this strong negative response should be further qualified. The existing tools are insufficient to support widespread use of QoS-based services on the multiprovider public Internet. The qualification is that within the enterprise network environment there are much stronger drivers for QoS mechanisms and much greater levels of administrative control over the overall network architecture, while within the multiprovider public Internet, these drivers are not apparent. The enterprise approach may also have some parallels within a single IP carrier’s network, or even across some forms of bilateral agreements between carriers. However, such approaches are not anticipated to be a widespread feature of the public Internet service environment.

Let’s look more closely at the public Internet and QoS to see why there is a mismatch between the two. The major stumbling blocks in attempting to address how QoS could be deployed in the public Internet are both engineering and economic in nature.

From an engineering perspective, we need to remember that in order to actually deliver any reasonable assurance of a quality-differentiated service, the service-quality mechanism chosen must be deployed across all networks along the end-to-end paths of the quality-service traffic. In a heterogeneous multiprovider environment such as the public Internet, this outcome is very unlikely. Within the tens of thousands of component service providers that make up the global Internet, such uniformity of action is highly improbable. The IPv6 transition structure correctly identifies the first step as isolated “islands” of IPv6 functionality, interconnected by some form of IPv6 “bridges.” While the potential scenario of initial QoS deployment may be similar, in terms of isolated islands of deployment of QoS services, there is a much stricter requirement for the “bridges” across the non-QoS-aware parts of the network; namely, that they do not distort the service outcomes. In effect, this scenario requires a QoS response from a non-QoS system (see Figure 3).

Figure 3: Attempted
End-to-End QoS across
the Public Internet



The engineering issues are deeper than simply the considerations of transition within a potential deployment scenario. The issues include:

- The need for QoS-enabled applications that can predict their service requirements in advance, and be able to signal these requirements into the network.
- In the case of the differentiated service approach of admission controls, there is a requirement for the interior of the network to be able to signal current load conditions to the network admission systems. This system also requires that the admission control points be able to use admission-decision support systems in order to include consideration of the service load, the current network load, and the policy parameters of the network that may allow some level of preemption of various admission decisions in order to meet high-priority service requirements.
- The signaling and negotiation aspect of QoS extends into the inter-domain space, where two or more service providers need to negotiate mutually acceptable service profiles, and associated service access. This extends beyond the addition of bilateral agreements and encompasses the requirement to add QoS attributes to interdomain routing protocols. The tools and operating techniques required to support this functionality remain poorly defined.
- Measurement of service performance remains an area in which existing measurement tools are lacking. While it is possible to instrument every active device within a network into a network management system, such an element-by-element view does not readily translate to the end-to-end view of application service performance.

From an economic perspective, we must remember that no current Internet retail tariff includes a concept of end-to-end tariffed transactions. All tariffs are access based, because application transactions are not readily visible to the Internet network. In addition, no technically stable or financially stable structure of interprovider interconnection financial settlements exists today. The financial model of the Internet from an economic viewpoint is very polarized, with only customer and zero-dollar peer arrangements dominating the interprovider space. However, end-to-end QoS transactions demand a different economic model.

The initiator of the end-to-end QoS transaction has the discretion of choosing whether to request an end-to-end service profile. If such a profile is requested, the initiator should pay the initiating provider a retail tariff to cover the entire end-to-end cost of the transaction, and the initiating provider must then indicate a willingness to financially settle with transit peer networks in order for these transit peers to devote additional resources to service the traffic associated with this transaction, and so forth through the entire path of transit providers. The arbitrary nature of the Internet transits, the dynamic nature of routing, and the lack of transaction setups in any scalable form of QoS mechanisms make this entire scenario highly improbable within our current understanding of interprovider policy-management mechanisms.

The relatively loosely coordinated structure of the public Internet will have to change from the state we have today if we want to use QoS-based services. The changes include:

- A common selection of a set of QoS mechanisms to deploy,
- Ubiquitous deployment of these mechanisms across both service provider and client networks,
- The adoption of a uniform set of retail tariffs for QoS services,
- The definition and common acceptance of multi-party QoS-related financial settlements that support fair and equitable cost distribution among multiple providers, and
- The definition of commonly accepted service performance metrics and related measurement methodologies to allow end-to-end and network-by-network service outcomes to be objectively assessed.

This is a significant agenda for the industry at large to undertake, and more so in an environment that features diversity and vigorous competition between various public Internet service providers.

An additional factor is also working against QoS deployment in the public Internet space. The increasing availability of very-high-speed transmission systems is bringing network carriage capacity down to the level of an abundant commodity across large parts of the Internet world. As the unit costs of network capacity decline in the face of increasing levels of availability of transmission systems, the market niche that QoS could occupy in managing a scarce resource is shrinking. The driver for QoS deployment is not that the best-effort service is not good enough. The problem that QoS is attempting to address is one of allocation of network capacity at those points in time when the network is under heavy load, or, in other words, taking on the task of allocating capacity when there is not enough network capacity to meet every demand. When a network is under load, the QoS response is to place additional control functionality in both applications and in the network to manage this allocation function. Obviously such an activity imposes additional costs on the network operators and the network client. Such additional costs have not created any additional network capacity.

The total sum of demand remains in excess of capacity after the deployment of QoS mechanisms. The alternative approach is to incur additional cost by augmenting the capacity of the network. This approach minimizes the impact of load on the network causing disruption to individual transactions. Again this approach imposes additional costs onto the network, but in an environment of abundant transmission capacity, it may often be the more cost-effective approach.

Where does this leave QoS and the public Internet? There is no doubt that QoS is a very stimulating area of research, with much to offer the enterprise network environment, but in asking for QoS to be deployed within the existing incarnation of the public multiprovider Internet, we may be simply asking for too much at this point in time. More effort is required to turn a QoS Internet into a reliable production platform.

Further Reading

- [1] Huston, G., *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, John Wiley & Sons, January 2000.

A detailed examination of Internet Quality of Service technologies and their potential application within the Internet.

- [2] Kilkki, K., *Differentiated Services for the Internet*, ISBN 1578701325, Macmillan Technical Publishing, June 1999.

An in-depth look at the Differentiated Services architecture and its use in enabling networks to handle traffic classes in a specific manner.

- [3] Durham, D., and Yavatar, R., *Inside the Internet's Resource Reservation Protocol: Foundations for Quality of Service*, ISBN 0471322148, John Wiley & Sons, April 1999.

At the core of the Integrated Services architecture is a signaling protocol to undertake service reservations. The Resource ReSerVation Protocol (RSVP) is a signaling protocol that can undertake this role. This book describes both the Integrated Services architecture and RSVP in detail.

- [4] Odlyzko, A., "The Economics of the Internet: Utility, Utilization, Pricing, and Quality of Service," 1998. Available at:

www.research.att.com/~amo

A paper arguing the point of view that overprovisioning data networks is a viable and economically sustainable response to the demands for service quality within data networks, and that such a response is technically and economically superior to implementing QoS responses within the network.

- [5] Braden, R., Clark, D., and Shenker, S., "Integrated Services in the Internet Architecture: An Overview," RFC 1633, June 1994.

This RFC describes the components of the Integrated Services architecture, a proposed extension to the Internet architecture, and protocols to support real-time traffic flows through service-quality commitments.

- [6] Blake, S., Black, D., Carlson, M., Davies, E., Wang, Z., and Weiss, W., "An Architecture for Differentiated Services," RFC 2475, Proposed Standard, December 1998.

The architecture description for the Differentiated Services enhancements to the Internet Protocol. This architecture achieves scalability by aggregating traffic classification state, which is conveyed by means of IP-layer packet marking using the Differentiated Services (DS) field. Packets are classified and marked to receive a particular per-hop forwarding behavior on nodes along their path. Sophisticated classification, marking, policing, and shaping operations need to be implemented only at network boundaries or hosts. Network resources are allocated to traffic streams by service-provisioning policies that govern how traffic is marked and conditioned upon entry to a differentiated services-capable network, and how that traffic is forwarded within that network.

- [7] Gray, T., "Enterprise QoS Survival Guide: 1999 Edition," 1999. Available at:

<http://staff.washington.edu/gray/papers/eqos22.html>

A detailed view of an approach to supporting QoS in an enterprise environment. The paper is an excellent example of the procedural steps involved in network engineering, detailing the intended environment, the available tools and the desired outcomes, and then examining the viability of a number of QoS solutions.

- [8] Huston, G., "Next Steps for the IP QoS Architecture." Available at:

www.ietf.org/internet-drafts/draft-iab-qos-00.txt

While there has been significant progress in the definition of IP QoS architecture, there are a number of aspects of QoS that appear to need further elaboration as they relate to translating a set of tools into a coherent platform for end-to-end service delivery. This document highlights the outstanding issues relating to the deployment and use of QoS mechanisms within the Internet, noting those areas where further standards work may be required. This draft is a work item of the Internet Architecture Board Working Group of the IETF.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. Huston is currently the Chief Technologist in the Internet area for Telstra. He is also an active member of the IETF, and is the chair of the Internet Society Board of Trustees. He is author of *The ISP Survival Guide*, ISBN 0-471-31499-4, *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089 and coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. All three books are published by John Wiley & Sons. E-mail: **gih@telstra.net**

Book Review

Removing the Spam *Removing the Spam: Email Processing and Filtering*, Geoff Mulligan, ISBN 0-201-37957-0, Addison-Wesley, 1999.
<http://cseng.aw.com/bookdetail.qry?ISBN=0-201-37957-0&ptype=0>

Do not be fooled by the title of this book. You might purchase this book, part of the Addison-Wesley Networking Basics Series, thinking you are just getting information dealing with unsolicited commercial e-mail (commonly called, to Hormel's displeasure, "spam"). The title is probably the work of a marketer who thought "spam" in the title would *sell*! The subtitle really describes the meat of the matter. This short, but thorough, book is about e-mail processing and filtering—dealing with spam, yes, but so much more.

A collection of essential information for the Internet e-mail "gatekeeper," *Removing the Spam* is really geared for the gatekeeper using a UNIX-based system, so NT system administrators be forewarned. Being an e-mail gatekeeper on the Internet involves keeping the e-mail flowing, making sure the automated processes in place do the job, supporting e-mail "mailing lists," and providing the services and features your users want or need for e-mail processing.

Commercial products support some of the many requirements, but the best software for most of these functions is freely available on the Internet. Geoff provides answers to the requirements using the most popular and commonly used solutions: *Sendmail* for mail delivery, *procmail* for e-mail filtering, and *majordomo* and *smartlist* for mailing-list management.

The book, however, tries to do a bit too much. Geoff indicates that the intended audience is not only the system administrator, but also the e-mail end users wanting to filter their own personal e-mail as well as those who want to run their own mailing list. Because of this broad audience, there are times when the book delves too long in the basics, giving the impression of topics added to lengthen the book. The overview of IP protocols, the brief history of the Internet, suggestions for users dealing with spammers, and mailing-list etiquette are examples that come to mind. Nevertheless, the other topics covered are "net essentials," and worth skimming over the already known.

The book clearly defines spam and its evils, and presents the tools and techniques available for removing, or at least minimizing, the spam. It is probably too ambitious when covering e-mail forgery and tracing e-mail spam, but leaves no essential unmentioned.

Sendmail coverage is good, dealing with installation as well as configuration, highlighting antispam features, and how to use them. Though not covering as much detail as other books that focus on Sendmail, the important elements of building and modifying are handled, as well as Sendmail's use of data bases, including the infamous "Realtime Black-hole List" (<http://maps.vix.com/rbl/>).

The e-mail gatekeeper, as well as end users of e-mail, can use procmail to preprocess e-mail before final delivery. Procmail is powerful and flexible, and, so, can be difficult to configure properly. Configuration files examples with explanations allow even the procmail-savvy reader to learn and try something new.

The mailing list section again instructs both system administrator and user. Information about subscribing, unsubscribing, and getting information from the mailing list software is useful for the user. The administrator will appreciate the examples of getting, installing, configuring, and running majordomo and smartlist. Geoff gives suggestions about when a manual versus automated solution is best.

About the Author

I knew Geoff back in our Digital Equipment Corporation days when he worked in the Network Systems Lab. My group ran one of the corporate Internet gateways, modeled after the one at NSL. Further, the group I ran also productized and delivered what is arguably the first commercial Internet firewall, based on a design from the team at NSL. All this to say, Geoff certainly has the background to write about these topics. Since those days, Geoff has been busy with other Internet endeavors, such as starting USA.NET and creating the NetAddress product (permanent, follow-you-anywhere e-mail addresses) and helping develop the Sun Microsystems Sunscreen Firewall. He also founded Geocast Network Systems. In various roles, in differing capacities, Geoff has had to wrestle with the matters covered in his book. What he writes is based on experience learned in the danger zone of the Internet gateway.

Organization

The book is divided into four chapters. The first chapter, the introduction really, is strangely entitled "The Dawn of Electronic Mail." This is also the "roughest" chapter. It is difficult to understand why some topics are covered in the order that they are here (and why some are covered at all—the aforementioned "list etiquette" and "Size and Growth of the Internet," for example). It introduces (needlessly, I think) The Internet Protocols, but then reviews the basics of understanding e-mail systems. It introduces spam, along with antispam resources, and the topics in the rest of the book to be covered in detail: e-mail processing, filtering, and e-mail lists.

Chapter 2 is entitled “Sendmail” and covers obtaining, installing, configuring, and running Sendmail on a UNIX machine. It gives the commands to build and install Sendmail and your Sendmail configuration file. This coverage is not detailed enough for *every* situation, but gives the most common configuration information, which should satisfy most readers’ needs. Included are instructions for using Sendmail to help stop (or avert) spam at the mail gateway.

Chapter 3 unravels the mysteries behind procmail configuration for e-mail filtering. This chapter covers getting the software, installing it, and using procmail—the latter for system administrators and users alike. There are example “ready-to-run filters” included. Caveat: Some of the scripts have inherent errors. No doubt these errors are unfortunate publication glitches, but they do detract from the usefulness of this chapter. Geoff has compiled an errata list with corrected scripts. This can be found at: <http://www.hz.com/spam/eratta>

Chapter 4 covers mailing lists, specifically discussing administering them “by hand” (just using Sendmail) or “automatically” (majordomo and smartlist). Again, examples are given with step-by-step commands.

Closing Thoughts

Production errors aside (the serious ones in the procmail chapter and others that are just nits to pick—the “P” in ARPA stands for “Projects,” not “Project”), this book is useful as an introduction as well as a reminder of things forgotten. I can recommend this book to the novice or seasoned e-mail gatekeeper, and I will recommend it to the students in my Sendmail courses.

—Frederick M. Avolio, Avolio Consulting
fred@avolio.com

Would You Like to Review a Book for IPJ?

We receive numerous books on computer networking from all the major publishers. If you’ve got a specific book you are interested in reviewing, please contact us and we will make sure a copy is mailed to you. The book is yours to keep if you send us a review. We accept reviews of new titles, as well as some of the “networking classics.” Contact us at ipj@cisco.com for more information.

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

ICANN Launches Membership Web Site for Individual Internet Users

The *Internet Corporation for Assigned Names and Numbers* (ICANN) recently announced the launch of its At Large Membership Web site. After considerable public input, the ICANN Board has developed this program as a new way for Internet users from all over the globe to participate directly in the ICANN process. Individuals can register to become ICANN members at <http://members.icann.org>

The At Large Membership of ICANN will give individual members of Internet communities worldwide a voice in the selection of Directors to the ICANN Board. By becoming an ICANN member, individuals will have an opportunity to become part of the ICANN “bottom-up” approach to making policy concerning Internet names and addresses. The basic requirements for applying to become an ICANN At Large member are: The completion of an online membership application, a working Internet e-mail address, and a single physical residence verified by a postal mail address. Thanks to a grant from the Markle Foundation, the initial launch of ICANN’s At Large Membership program has been funded without the need for membership dues.

The ICANN Board will consider and adopt further policy about composition and structure of the At Large Membership, and establish rules for the nomination and election of candidates for the At Large Council. It is hoped that the target goal of 5,000 members can be reached in the next few weeks in order to move forward with the At Large Elections later this year.

ICANN is a non-profit, international corporation formed to oversee a select set of Internet technical management functions currently managed by the U.S. Government, or by its contractors and volunteers. Specifically, ICANN is assuming responsibility for coordinating the management of the *Domain Name System* (DNS), the allocation of IP address space, the assignment of protocol parameters, and the management of the root server system.

Online Registration for INET 2000 Now Open

INET 2000, the annual conference of the Internet Society (ISOC) will be held in Yokohama, Japan, July 18–21. You can register for this event by visiting ISOC’s Web site at:

<http://www.isoc.org/inet2000/register.shtml>

Denial of Service Attacks

In early February, several high-profile Internet Web sites were severely disrupted by a number of so-called *Distributed Denial of Service* (DDoS) attacks. We plan to publish an article on this topic in the future. Meanwhile, we recommend you visit the Denial of Service Resource Page at <http://www.denialinfo.com/>

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Engineering
MCI WorldCom, USA

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Edward R. Kozel, Member of The Board of Directors
Cisco Systems, Inc., USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the Chief
Technology Office, Cisco Systems, Inc.
www.cisco.com*

*Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 2000 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-10/5
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

Bulk Rate Mail U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol Journal

June 2000

Volume 3, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

F R O M T H E E D I T O R

In This Issue

From the Editor	1
TCP Performance	2
Internet Mail Standards	25
Book Review	37
Fragments	39

Two protocols used in the Internet are so important that they deserve special attention: the *Internet Protocol* (IP) from which this journal takes its name, and the *Transmission Control Protocol* (TCP). IP is fundamental to Internet addressing and routing, while TCP provides a reliable transport service that is used by most Internet applications, including interactive Telnet, file transfer, electronic mail, and Web page access via HTTP. Because of the critical importance of TCP to the operation of the Internet, it has received much attention in the research community over the years. As a result, numerous improvements to implementations of TCP have been developed and deployed. In this issue, Geoff Huston takes a detailed look at TCP from a performance perspective and describes several enhancements to the original protocol. In a second article, Geoff will look at the challenges facing TCP in a rapidly growing and changing Internet, and describe work to further augment TCP.

Electronic mail is by far the most used of all Internet applications. The fundamental protocols for delivery and retrieval of e-mail have not changed much since the early days of the ARPANET, but as with TCP, many enhancements have been added to accommodate new uses of e-mail. Today, Internet e-mail supports international character sets, includes the ability to send file attachments, and allows roaming e-mail clients to authenticate themselves to servers. All of this has been made possible by continued development in the *Internet Engineering Task Force* (IETF). In our second article, Paul Hoffman of the Internet Mail Consortium gives an overview of Internet mail standards.

This is the second anniversary issue of *The Internet Protocol Journal* (IPJ). By now more than 10,000 people from virtually every country in the world have subscribed to the paper edition of IPJ. In order to serve our readers better, we are developing an online subscription system, which will be deployed in July 2000. With this new system you will be able to modify your mailing address as well as select your preferred delivery method for the journal. You can choose to receive IPJ on paper, or be notified via e-mail when a new issue becomes available on line. More information about this new system can be found on our Web site at www.cisco.com/ipj. We would love to hear your feedback on this system and any other aspect of IPJ. Please send your comments to ipj@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

TCP Performance

by Geoff Huston, Telstra

The *Transmission Control Protocol* (TCP) and the *User Datagram Protocol* (UDP) are both IP transport-layer protocols. UDP is a lightweight protocol that allows applications to make direct use of the unreliable datagram service provided by the underlying IP service. UDP is commonly used to support applications that use simple query/response transactions, or applications that support real-time communications. TCP provides a reliable data-transfer service, and is used for both bulk data transfer and interactive data applications. TCP is the major transport protocol in use in most IP networks, and supports the transfer of over 90 percent of all traffic across the public Internet today. Given this major role for TCP, the performance of this protocol forms a significant part of the total picture of service performance for IP networks. In this article we examine TCP in further detail, looking at what makes a TCP session perform reliably and well. This article draws on material published in the *Internet Performance Survival Guide*^[1].

Overview of TCP

TCP is the embodiment of reliable end-to-end transmission functionality in the overall Internet architecture. All the functionality required to take a simple base of IP datagram delivery and build upon this a control model that implements reliability, sequencing, flow control, and data streaming is embedded within TCP^[2].

TCP provides a communication channel between processes on each host system. The channel is reliable, full-duplex, and streaming. To achieve this functionality, the TCP drivers break up the session data stream into discrete segments, and attach a TCP header to each segment. An IP header is attached to this TCP packet, and the composite packet is then passed to the network for delivery. This TCP header has numerous fields that are used to support the intended TCP functionality. TCP has the following functional characteristics:

- *Unicast protocol*: TCP is based on a unicast network model, and supports data exchange between precisely two parties. It does not support broadcast or multicast network models.
- *Connection state*: Rather than impose a state within the network to support the connection, TCP uses synchronized state between the two endpoints. This synchronized state is set up as part of an initial connection process, so TCP can be regarded as a connection-oriented protocol. Much of the protocol design is intended to ensure that each local state transition is communicated to, and acknowledged by, the remote party.
- *Reliable*: Reliability implies that the stream of octets passed to the TCP driver at one end of the connection will be transmitted across the network so that the stream is presented to the remote process as the same sequence of octets, in the same order as that generated by the sender.

This implies that the protocol detects when segments of the data stream have been discarded by the network, reordered, duplicated, or corrupted. Where necessary, the sender will retransmit damaged segments so as to allow the receiver to reconstruct the original data stream. This implies that a TCP sender must maintain a local copy of all transmitted data until it receives an indication that the receiver has completed an accurate transfer of the data.

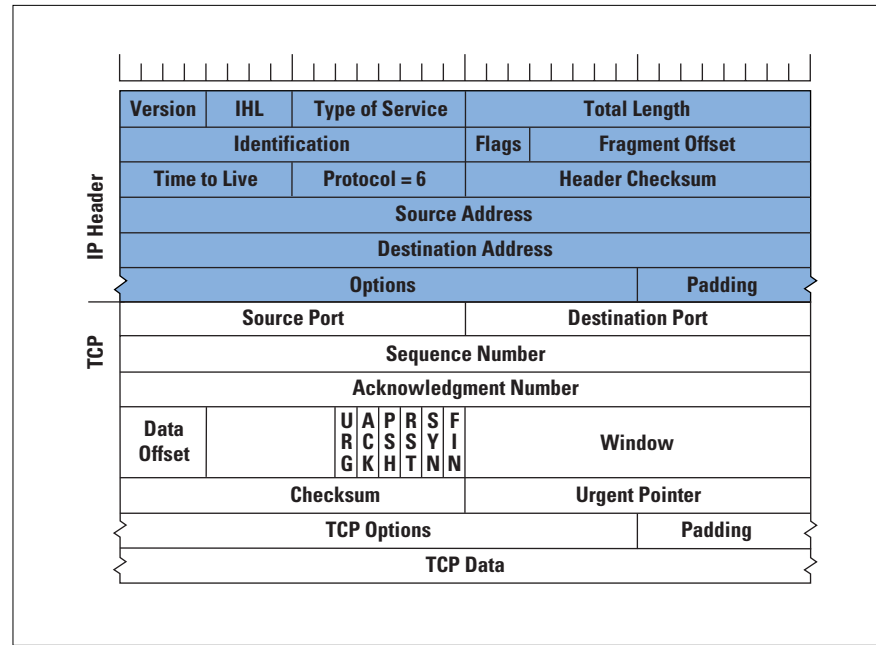
- *Full duplex*: TCP is a full-duplex protocol; it allows both parties to send and receive data within the context of the single TCP connection.
- *Streaming*: Although TCP uses a packet structure for network transmission, TCP is a true streaming protocol, and application-level network operations are not transparent. Some protocols explicitly encapsulate each application transaction; for every *write*, there must be a matching *read*. In this manner, the application-derived segmentation of the data stream into a logical record structure is preserved across the network. TCP does not preserve such an implicit structure imposed on the data stream, so that there is no pairing between *write* and *read* operations within the network protocol. For example, a TCP application may *write* three data blocks in sequence into the network connection, which may be collected by the remote reader in a single *read* operation. The size of the data blocks (segments) used in a TCP session is negotiated at the start of the session. The sender attempts to use the largest segment size it can for the data transfer, within the constraints of the maximum segment size of the receiver, the maximum segment size of the configured sender, and the maximum supportable non-fragmented packet size of the network path (path *Maximum Transmission Unit* [MTU]). The path MTU is refreshed periodically to adjust to any changes that may occur within the network while the TCP connection is active.
- *Rate adaptation*: TCP is also a rate-adaptive protocol, in that the rate of data transfer is intended to adapt to the prevailing load conditions within the network and adapt to the processing capacity of the receiver. There is no predetermined TCP data-transfer rate; if the network and the receiver both have additional available capacity, a TCP sender will attempt to inject more data into the network to take up this available space. Conversely, if there is congestion, a TCP sender will reduce its sending rate to allow the network to recover. This adaptation function attempts to achieve the highest possible data-transfer rate without triggering consistent data loss.

The TCP Protocol Header

The TCP header structure, shown in Figure 1, uses a pair of 16-bit source and destination *Port* addresses. The next field is a 32-bit *sequence number*, which identifies the sequence number of the first data octet in this packet. The sequence number does not start at an initial value of 1 for each new TCP connection; the selection of an initial value is critical, because the initial value is intended to prevent delayed data

from an old connection from being incorrectly interpreted as being valid within a current connection. The sequence number is necessary to ensure that arriving packets can be ordered in the sender's original order. This field is also used within the flow-control structure to allow the association of a data packet with its corresponding acknowledgement, allowing a sender to estimate the current round-trip time across the network.

Figure 1: The TCP/IP Datagram



The *acknowledgment sequence number* is used to inform the remote end of the data that has been successfully received. The acknowledgment sequence number is actually one greater than that of the last octet correctly received at the local end of the connection. The *data offset* field indicates the number of four-octet words within the TCP header. Six single *bit flags* are used to indicate various conditions. URG is used to indicate whether the *urgent pointer* is valid. ACK is used to indicate whether the *acknowledgment* field is valid. PSH is set when the sender wants the remote application to *push* this data to the remote application. RST is used to *reset* the connection. SYN (for *synchronize*) is used within the connection startup phase, and FIN (for *finish*) is used to close the connection in an orderly fashion. The *window* field is a 16-bit count of available buffer space. It is added to the acknowledgment sequence number to indicate the highest sequence number the receiver can accept. The TCP *checksum* is applied to a synthesized header that includes the source and destination addresses from the outer IP datagram. The final field in the TCP header is the *urgent pointer*, which, when added to the sequence number, indicates the sequence number of the final octet of urgent data if the urgent flag is set.

Many options can be carried in a TCP header. Those relevant to TCP performance include:

- *Maximum-receive-segment-size option:* This option is used when the connection is being opened. It is intended to inform the remote end of the maximum segment size, measured in octets, that the sender is willing to receive on the TCP connection. This option is used only in the initial SYN packet (the initial packet exchange that opens a TCP connection). It sets both the maximum receive segment size and the maximum size of the advertised TCP window, passed to the remote end of the connection. In a robust implementation of TCP, this option should be used with path MTU discovery to establish a segment size that can be passed across the connection without fragmentation, an essential attribute of a high-performance data flow.
- *Window-scale option:* This option is intended to address the issue of the maximum window size in the face of paths that exhibit a high-delay bandwidth product. This option allows the window size advertisement to be right-shifted by the amount specified (in binary arithmetic, a right-shift corresponds to a multiplication by 2). Without this option, the maximum window size that can be advertised is 65,535 bytes (the maximum value obtainable in a 16-bit field). The limit of TCP transfer speed is effectively one window size in transit between the sender and the receiver. For high-speed, long-delay networks, this performance limitation is a significant factor, because it limits the transfer rate to at most 65,535 bytes per round-trip interval, regardless of available network capacity. Use of the window-scale option allows the TCP sender to effectively adapt to high-bandwidth, high-delay network paths, by allowing more data to be held in flight. The maximum window size with this option is 2^{30} bytes. This option is negotiated at the start of the TCP connection, and can be sent in a packet only with the SYN flag. Note that while an MTU discovery process allows optimal setting of the maximum-receive-segment-size option, no corresponding bandwidth delay product discovery allows the reliable automated setting of the window-scale option^[3].
- *SACK-permitted option and SACK option:* This option alters the acknowledgment behavior of TCP. SACK is an acronym for *selective acknowledgment*. The SACK-permitted option is offered to the remote end during TCP setup as an option to an opening SYN packet. The SACK option permits selective acknowledgment of permitted data. The default TCP acknowledgment behavior is to acknowledge the highest sequence number of in-order bytes. This default behavior is prone to cause unnecessary retransmission of data, which can exacerbate a congestion condition that may have been the cause of the original packet loss. The SACK option allows the receiver to modify the acknowledgment field to describe noncontinuous blocks of received data, so that the sender can retransmit only what is missing at the receiver's end^[4].

Any robust high-performance implementation of TCP should negotiate these parameters at the start of the TCP session, ensuring the following: that the session is using the largest possible IP packet size that can be carried without fragmentation, that the window sizes used in the transfer are adequate for the bandwidth-delay product of the network path, and that selective acknowledgment can be used for rapid recovery from line-error conditions or from short periods of marginally degraded network performance.

TCP Operation

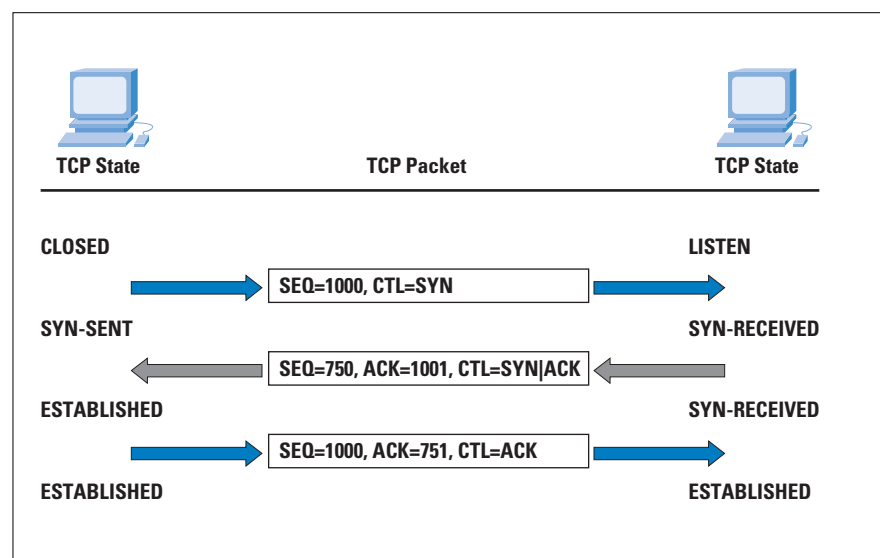
The first phase of a TCP session is establishment of the connection. This requires a *three-way handshake*, ensuring that both sides of the connection have an unambiguous understanding of the sequence number space of the remote side for this session. The operation of the connection is as follows:

- The local system sends the remote end an initial sequence number to the remote port, using a SYN packet.
- The remote system responds with an ACK of the initial sequence number and the initial sequence number of the remote end in a response SYN packet.
- The local end responds with an ACK of this remote sequence number.

The connection is opened.

The operation of this algorithm is shown in Figure 2. The performance implication of this protocol exchange is that it takes one and a half *round-trip times* (RTTs) for the two systems to synchronize state before any data can be sent.

Figure 2:
TCP Connection
Handshake



After the connection has been established, the TCP protocol manages the reliable exchange of data between the two systems. The algorithms that determine the various retransmission timers have been redefined numerous times. TCP is a *sliding-window* protocol, and the general principle of flow control is based on the management of the advertised window size and the management of retransmission timeouts, attempting to optimize protocol performance within the observed delay and loss parameters of the connection. Tuning a TCP protocol stack for optimal performance over a very low-delay, high-bandwidth LAN requires different settings to obtain optimal performance over a dialup Internet connection, which in turn is different for the requirements of a high-speed wide-area network. Although TCP attempts to discover the delay bandwidth product of the connection, and attempts to automatically optimize its flow rates within the estimated parameters of the network path, some estimates will not be accurate, and the corresponding efforts by TCP to optimize behavior may not be completely successful.

Another critical aspect is that TCP is an adaptive flow-control protocol. TCP uses a basic flow-control algorithm of increasing the data-flow rate until the network signals that some form of saturation level has been reached (normally indicated by data loss). When the sender receives an indication of data loss, the TCP flow rate is reduced; when reliable transmission is reestablished, the flow rate slowly increases again.

If no reliable flow is reestablished, the flow rate backs further off to an initial probe of a single packet, and the entire adaptive flow-control process starts again.

This process has numerous results relevant to service quality. First, TCP behaves *adaptively*, rather than *predictively*. The flow-control algorithms are intended to increase the data-flow rate to fill all available network path capacity, but they are also intended to quickly back off if the available capacity changes because of interaction with other traffic, or if a dynamic change occurs in the end-to-end network path. For example, a single TCP flow across an otherwise idle network attempts to fill the network path with data, optimizing the flow rate within the available network capacity. If a second TCP flow opens up across the same path, the two flow-control algorithms will interact so that both flows will stabilize to use approximately half of the available capacity per flow. The objective of the TCP algorithms is to adapt so that the network is fully used whenever one or more data flows are present. In design, tension always exists between the efficiency of network use and the enforcement of predictable session performance. With TCP, you give up predictable throughput but gain a highly utilized, efficient network.

Protocol Performance

In this section we examine the transfer of data using the TCP protocol, focusing on the relationship between the protocol and performance. TCP is generally used within two distinct application areas: short-delay short data packets sent on demand, to support interactive applications such as *Telnet*, or *rlogin*, and large packet data streams supporting reliable volume data transfers, such as mail transfers, Web-page transfers, and *File Transfer Protocol* (FTP). Different protocol mechanisms come into play to support interactive applications, as distinct from short- and long-held volume transactions.

Interactive TCP

Interactive protocols are typically directed at supporting single-character interactions, where each character is carried in a single packet, as is its echo. The protocol interaction to support this is indicated in Figure 3. These 2 bytes of data generate four TCP/IP packets, or 160 bytes of protocol overhead. TCP makes some small improvement in this exchange through the use of *piggybacking*, where an ACK is carried in the same packet as the data, and *delayed acknowledgment*, where an ACK is delayed up to 200 ms before sending, to give the server application the opportunity to generate data that the ACK can piggyback. The resultant protocol exchange is indicated in Figure 4.

Figure 3:
Interactive Exchange

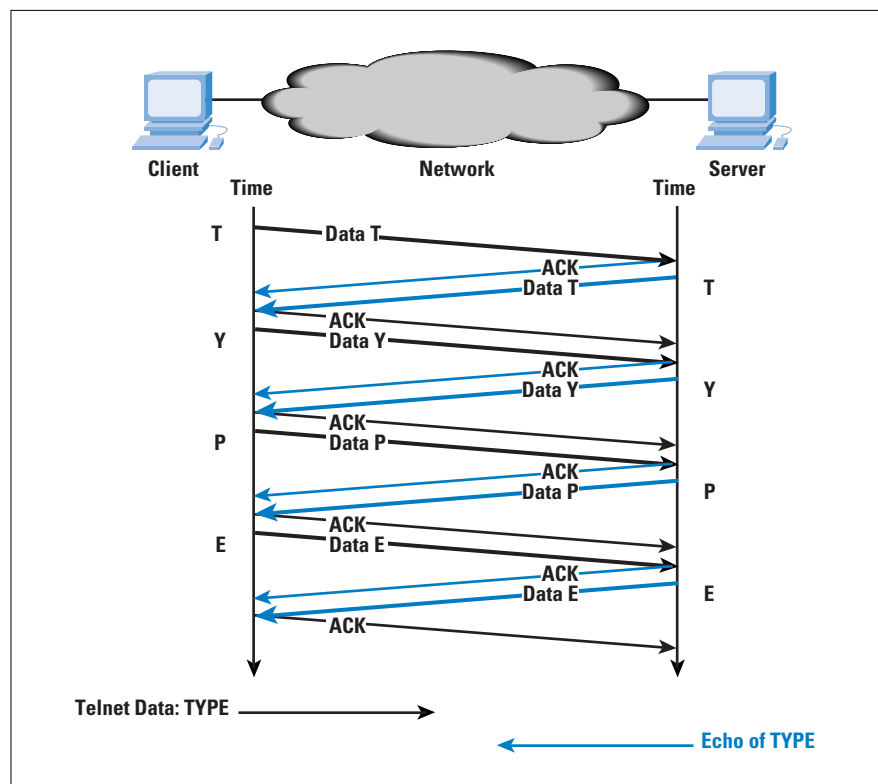
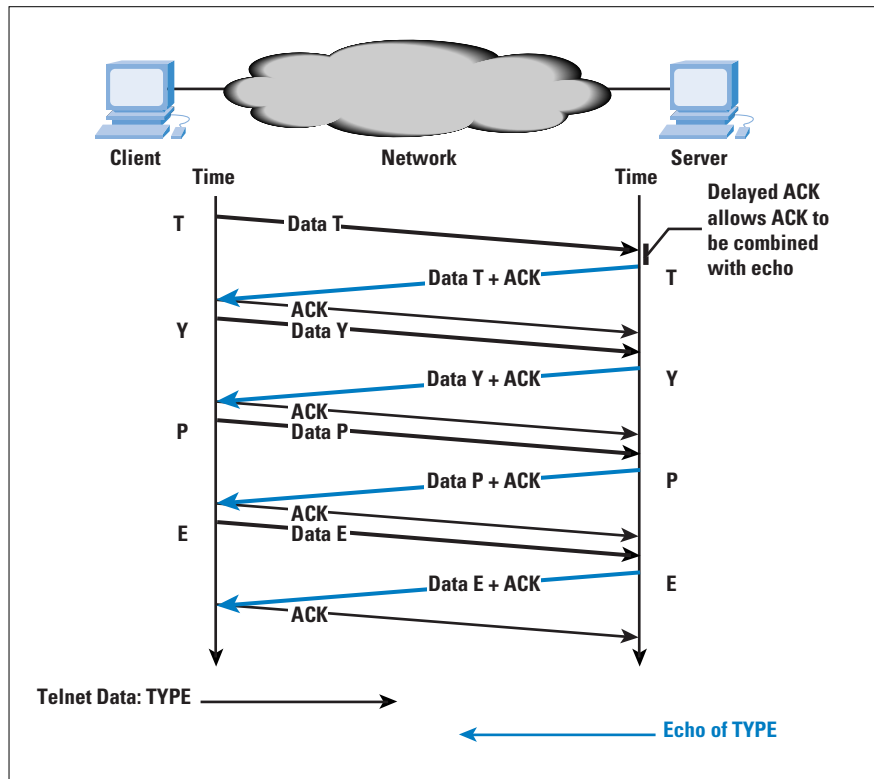


Figure 4:
Interactive Exchange
with Delayed ACK



For short-delay LANs, this protocol exchange offers acceptable performance. This protocol exchange for a single data character and its echo occurs within about 16 ms on an Ethernet LAN, corresponding to an interactive rate of 60 characters per second. When the network delay is increased in a WAN, these small packets can be a source of congestion load. The TCP mechanism to address this small-packet congestion was described by John Nagle in RFC 896^[5]. Commonly referred to as the *Nagle Algorithm*, this mechanism inhibits a sender from transmitting any additional small segments while the TCP connection has outstanding unacknowledged small segments. On a LAN, this modification to the algorithm has a negligible effect; in contrast, on a WAN, it has a dramatic effect in reducing the number of small packets in direct correlation to the network path congestion level (as shown in Figures 5 and 6). The cost is an increase in session jitter by up to a round-trip time interval. Applications that are jitter-sensitive typically disable this control algorithm.

TCP is not a highly efficient protocol for the transmission of interactive traffic. The typical carriage efficiency of the protocol across a LAN is 2 bytes of payload and 120 bytes of protocol overhead. Across a WAN, the Nagle algorithm may improve this carriage efficiency slightly by increasing the number of bytes of payload for each payload transaction, although it will do so at the expense of increased session jitter.

Figure 5: WAN
Interactive Exchange

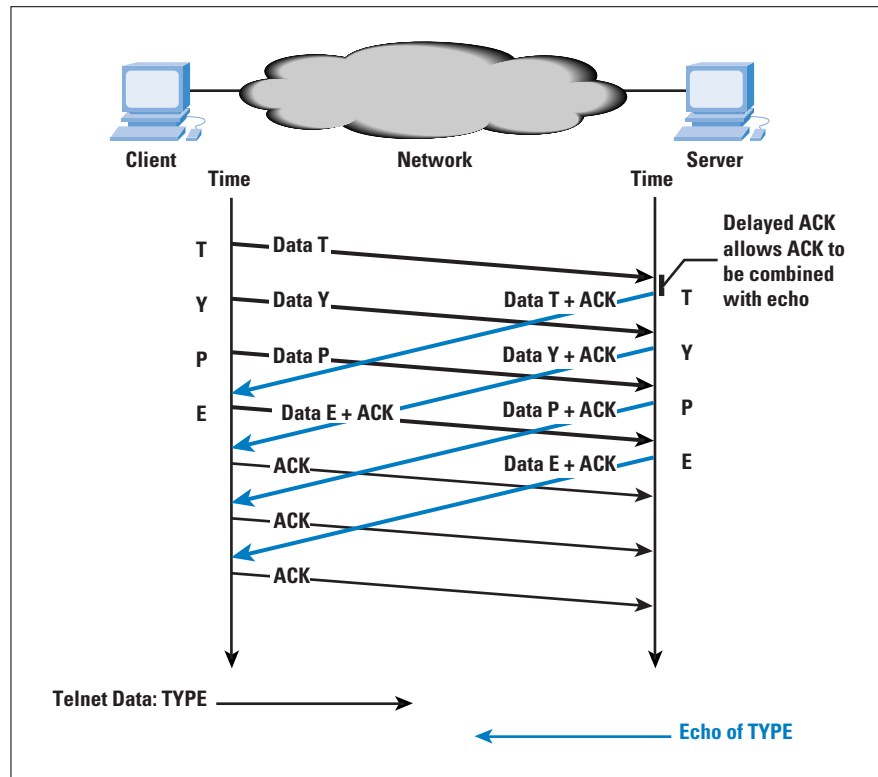
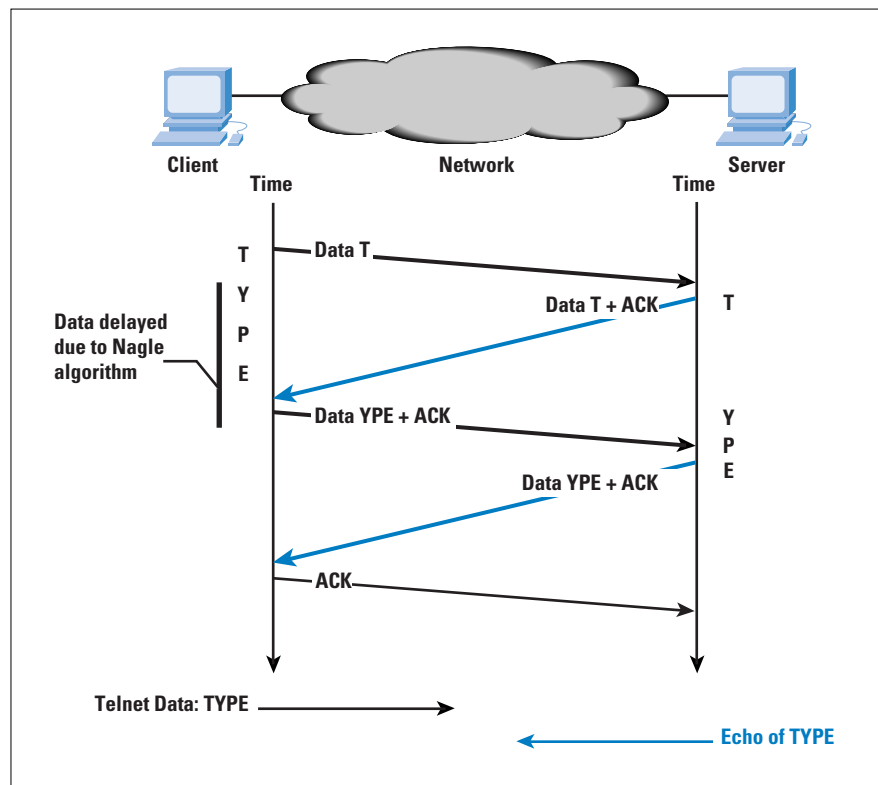


Figure 6: WAN
Interactive Exchange
with Nagle Algorithm



TCP Volume Transfer

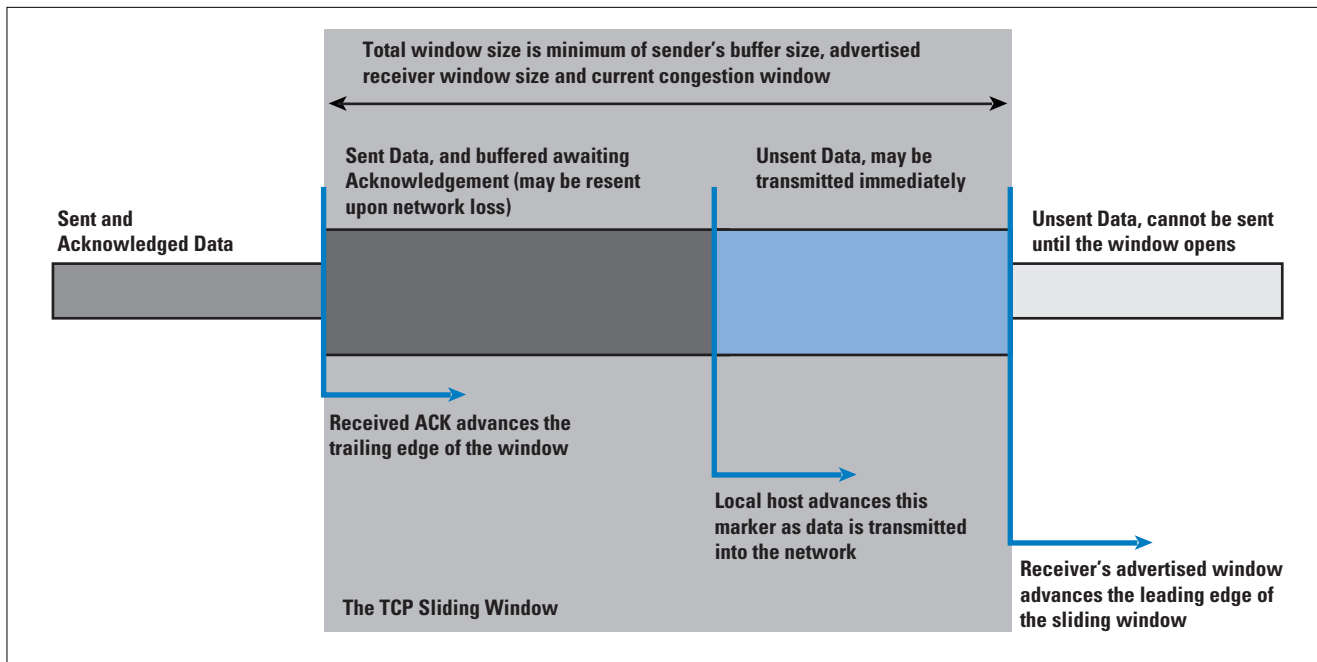
The objective for this application is to maximize the efficiency of the data transfer, implying that TCP should endeavor to locate the point of dynamic equilibrium of maximum network efficiency, where the sending data rate is maximized just prior to the onset of sustained packet loss.

Further increasing the sending rate from such a point will run the risk of generating a congestion condition within the network, with rapidly increasing packet-loss levels. This, in turn, will force the TCP protocol to retransmit the lost data, resulting in reduced data-transfer efficiency. On the other hand, attempting to completely eliminate packet-loss rates implies that the sender must reduce the sending rate of data into the network so as not to create transient congestion conditions along the path to the receiver. Such an action will, in all probability, leave the network with idle capacity, resulting in inefficient use of available network resources.

The notion of a point of equilibrium is an important one. The objective of TCP is to coordinate the actions of the sender, the network, and the receiver so that the network path has sufficient data such that the network is not idle, but it is not so overloaded that a congestion backlog builds up and data loss occurs. Maintaining this point of equilibrium requires the sender and receiver to be synchronized so that the sender passes a packet into the network at precisely the same time as the receiver removes a packet from the network. If the sender attempts to exceed this equilibrium rate, network congestion will occur. If the sender attempts to reduce its rate, the efficiency of the network will drop.

TCP uses a sliding-window protocol to support bulk data transfer (Figure 7). The receiver advertises to the sender the available buffer space at the receiver. The sender can transmit up to this amount of data before having to await a further buffer update from the receiver. The sender should have no more than this amount of data in transit in the network. The sender must also buffer sent data until it has been ACKed by the receiver. The send window is the minimum of the sender's buffer size and the advertised receiver window. Each time an ACK is received, the trailing edge of the send window is advanced. The minimum of the sender's buffer and the advertised receiver's window is used to calculate a new leading edge. If this send window encompasses unsent data, this data can be sent immediately.

Figure 7: TCP Sliding Window



The size of TCP buffers in each host is a critical limitation to performance in WANs. The protocol is capable of transferring one send window of data per round-trip interval. For example, with a send window of 4096 bytes and a transmission path with an RTT of 600 ms, a TCP session is capable of sustaining a maximum transfer rate of 48 Kbps, regardless of the bandwidth of the network path. Maximum efficiency of the transfer is obtained only if the sender is capable of completely filling the network path with data. Because the sender will have an amount of data in forward transit and an equivalent amount of data awaiting reception of an ACK signal, both the sender's buffer and the receiver's advertised window should be no smaller than the *Delay-Bandwidth Product* of the network path. That is:

$$\text{Window size} \geq \text{Bandwidth (bytes/sec)} \times \text{Round-trip time (sec)}$$

The 16-bit field within the TCP header can contain values up to 65,535, imposing an upper limit on the available window size of 65,535 bytes. This imposes an upper limit on TCP performance of some 64 KB per RTT, even when both end systems have arbitrarily large send and receive buffers. This limit can be modified by the use of a window-scale option, described in RFC 1323, effectively increasing the size of the window to a 30-bit field, but transmitting only the most significant 16 bits of the value. This allows the sender and receiver to use buffer sizes that can operate efficiently at speeds that encompass most of the current very-high-speed network transmission technologies across distances of the scale of the terrestrial intercontinental cable systems.

Although the maximum window size and the RTT together determine the maximum achievable data-transfer rate, there is an additional element of flow control required for TCP. If a TCP session commenced by injecting a full window of data into the network, then there is a strong probability that much of the initial burst of data would be lost because of transient congestion, particularly if a large window is being used. Instead, TCP adopts a more conservative approach by starting with a modest amount of data that has a high probability of successful transmission, and then probing the network with increasing amounts of data for as long as the network does not show signs of congestion. When congestion is experienced, the sending rate is dropped and the probing for additional capacity is resumed.

The dynamic operation of the window is a critical component of TCP performance for volume transfer. The mechanics of the protocol involve an additional overriding modifier of the sender's window, the *congestion window*, referred to as *cwnd*. The objective of the window-management algorithm is to start transmitting at a rate that has a very low probability of packet loss, then to increase the rate (by increasing the *cwnd* size) until the sender receives an indication, through the detection of packet loss, that the rate has exceeded the available capacity of the network. The sender then immediately halves its sending rate by reducing the value of *cwnd*, and resumes a gradual increase of the sending rate. The goal is to continually modify the sending rate such that it oscillates around the true value of available network capacity. This oscillation enables a dynamic adjustment that automatically senses any increase or decrease in available capacity through the lifetime of the data flow.

The intended outcome is that of a dynamically adjusting cooperative data flow, where a combination of such flows behaves fairly, in that each flow obtains essentially a fair share of the network, and so that close to maximal use of available network resources is made. This flow-control functionality is achieved through a combination of *cwnd* value management and packet-loss and retransmission algorithms. TCP flow control has three major parts: the flow-control modes of *Slow Start* and *Congestion Avoidance*, and the response to packet loss that determines how TCP switches between these two modes of operation.

TCP Slow Start

The starting value of the *cwnd* window (the *Initial Window*, or IW) is set to that of the *Sender Maximum Segment Size* (SMSS) value. This SMSS value is based on the receiver's maximum segment size, obtained during the SYN handshake, the discovered path MTU (if used), the MTU of the sending interface, or, in the absence of other information, 536 bytes. The sender then enters a flow-control mode termed *Slow Start*.

The sender sends a single data segment, and because the window is now full, it then awaits the corresponding ACK. When the ACK is received, the sender increases its window by increasing the value of *cwnd* by the value of SMSS. This then allows the sender to transmit two segments; at that point, the congestion window is again full, and the sender must await the corresponding ACKs for these segments. This algorithm continues by increasing the value of *cwnd* (and, correspondingly, opening the size of the congestion window) by one SMSS for every ACK received that acknowledges new data.

If the receiver is sending an ACK for every packet, the effect of this algorithm is that the data rate of the sender doubles every round-trip time interval. If the receiver supports delayed ACKs, the rate of increase will be slightly lower, but nevertheless the rate will increase by a minimum of one SMSS each round-trip time. Obviously, this cannot be sustained indefinitely. Either the value of *cwnd* will exceed the advertised receive window or the sender's window, or the capacity of the network will be exceeded, in which case packets will be lost.

There is another limit to the slow-start rate increase, maintained in a variable termed *ssthresh*, or *Slow-Start Threshold*. If the value of *cwnd* increases past the value of *ssthresh*, the TCP flow-control mode is changed from Slow Start to congestion avoidance. Initially the value of *ssthresh* is set to the receiver's maximum window size. However, when congestion is noted, *ssthresh* is set to half the current window size, providing TCP with a memory of the point where the onset of network congestion may be anticipated in future.

One aspect to highlight concerns the interaction of the slow-start algorithm with high-capacity long-delay networks, the so-called *Long Fat Networks* (or LFNs, pronounced "elephants"). The behavior of the slow-start algorithm is to send a single packet, await an ACK, then send two packets, and await the corresponding ACKs, and so on. The TCP activity on LFNs tends to cluster at each epoch of the round-trip time, with a quiet period that follows after the available window of data has been transmitted. The received ACKs arrive back at the sender with an inter-ACK spacing that is equivalent to the data rate of the bottleneck point on the network path. During Slow Start, the sender transmits at a rate equal to twice this bottleneck rate. The rate adaptation function that must occur within the network takes place in the router at the entrance to the bottleneck point. The sender's packets arrive at this router at twice the rate of egress from the router, and the router stores the overflow within its internal buffer. When this buffer overflows, packets will be dropped, and the slow-start phase is over. The important conclusion is that the sender will stop increasing its data rate when there is buffer exhaustion, a condition that may not be the same as reaching the true available data rate. If the router has a buffer capacity considerably less than the delay-bandwidth product of the egress circuit, the two values are certainly not the same.

In this case, the TCP slow-start algorithm will finish with a sending rate that is well below the actual available capacity. The efficient operation of TCP, particularly in LFNs, is critically reliant on adequately large buffers within the network routers.

Another aspect of Slow Start is the choice of a single segment as the initial sending window. Experimentation indicates that an initial value of up to four segments can allow for a more efficient session startup, particularly for those short-duration TCP sessions so prevalent with Web fetches^[6]. Observation of Web traffic indicates an average Web data transfer of 17 segments. A slow start from one segment will take five RTT intervals to transfer this data, while using an initial value of four will reduce the transfer time to three RTT intervals. However, four segments may be too many when using low-speed links with limited buffers, so a more robust approach is to use an initial value of no more than two segments to commence Slow Start^[7].

Packet Loss

Slow Start attempts to start a TCP session at a rate the network can support and then continually increase the rate. How does TCP know when to stop this increase? This slow-start rate increase stops when the congestion window exceeds the receiver's advertised window, when the rate exceeds the remembered value of the onset of congestion as recorded in *ssthresh*, or when the rate is greater than the network can sustain. Addressing the last condition, how does a TCP sender know that it is sending at a rate greater than the network can sustain? The answer is that this is shown by data packets being dropped by the network. In this case, TCP has to undertake many functions:

- The packet loss has to be detected by the sender.
- The missing data has to be retransmitted.
- The sending data rate should be adjusted to reduce the probability of further packet loss.

TCP can detect packet loss in two ways. First, if a single packet is lost within a sequence of packets, the successful delivery packets following the lost packet will cause the receiver to generate a *duplicate* ACK for each successive packet. The reception of these duplicate ACKs is a signal of such packet loss. Second, if a packet is lost at the end of a sequence of sent packets, there are no following packets to generate duplicate ACKs. In this case, there are no corresponding ACKs for this packet, and the sender's retransmit timer will expire and the sender will assume packet loss.

A single duplicate ACK is not a reliable signal of packet loss. When a TCP receiver gets a data packet with an out-of-order TCP sequence value, the receiver must generate an immediate ACK of the highest in-order data byte received. This will be a duplicate of an earlier transmitted ACK. Where a single packet is lost from a sequence of packets, all subsequent packets will generate a duplicate ACK packet.

On the other hand, where a packet is rerouted with an additional incremental delay, the reordering of the packet stream at the receiver's end will generate a small number of duplicate ACKs, followed by an ACK of the entire data sequence, after the errant packet is received. The sender distinguishes between these cases by using three duplicate ACK packets as a signal of packet loss.

The third duplicate ACK triggers the sender to immediately send the segment referenced by the duplicate ACK value (*fast retransmit*) and commence a sequence termed *Fast Recovery*. In fast recovery, the value of *ssthresh* is set to half the current send window size (the send window is the amount of unacknowledged data outstanding). The congestion window, *cwnd*, is set three segments greater than *ssthresh* to allow for three segments already buffered at the receiver. If this allows additional data to be sent, then this is done. Each additional duplicate ACK inflates *cwnd* by a further segment size, allowing more data to be sent. When an ACK arrives that encompasses new data, the value of *cwnd* is set back to *ssthresh*, and TCP enters congestion-avoidance mode. Fast Recovery is intended to rapidly repair single packet loss, allowing the sender to continue to maintain the ACK-clocked data rate for new data while the packet loss repair is being undertaken. This is because there is still a sequence of ACKs arriving at the sender, so that the network is continuing to pass timing signals to the sender indicating the rate at which packets are arriving at the receiver. Only when the repair has been completed does the sender drop its window to the *ssthresh* value as part of the transition to congestion-avoidance mode^[8].

The other signal of packet loss is a complete cessation of any ACK packets arriving to the sender. The sender cannot wait indefinitely for a delayed ACK, but must make the assumption at some point in time that the next unacknowledged data segment must be retransmitted. This is managed by the sender maintaining a *Retransmission Timer*. The maintenance of this timer has performance and efficiency implications. If the timer triggers too early, the sender will push duplicate data into the network unnecessarily. If the timer triggers too slowly, the sender will remain idle for too long, unnecessarily slowing down the flow of data. The TCP sender uses a timer to measure the elapsed time between sending a data segment and receiving the corresponding acknowledgment. Individual measurements of this time interval will exhibit significant variance, and implementations of TCP use a smoothing function when updating the retransmission timer of the flow with each measurement. The commonly used algorithm was originally described by Van Jacobson^[9], modified so that the retransmission timer is set to the smoothed round-trip-time value, plus four times a smoothed mean deviation factor^[10].

When the retransmission timer expires, the actions are similar to that of duplicate ACK packets, in that the sender must reduce its sending rate in response to congestion. The threshold value, *ssthresh*, is set to half of the current value of outstanding unacknowledged data, as in the duplicate ACK case. However, the sender cannot make any valid assumptions about the current state of the network, given that no useful information has been provided to the sender for more than one RTT interval. In this case, the sender closes the congestion window back to one segment, and restarts the flow in slow start-mode by sending a single segment. The difference from the initial slow start is that, in this case, the *ssthresh* value is set so that the sender will probe the congestion area more slowly using a linear sending rate increase when the congestion window reaches the remembered *ssthresh* value.

Congestion Avoidance

Compared to Slow Start, congestion avoidance is a more tentative probing of the network to discover the point of threshold of packet loss. Where Slow Start uses an exponential increase in the sending rate to find a first-level approximation of the loss threshold, congestion avoidance uses a linear growth function.

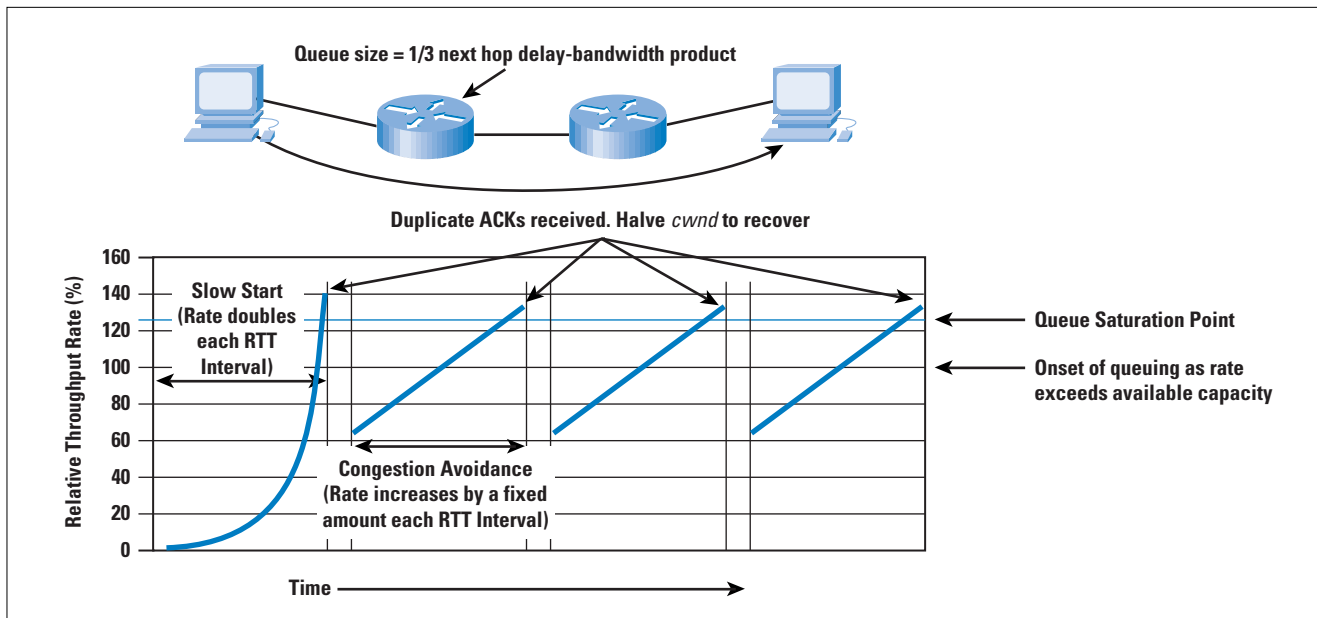
When the value of *cwnd* is greater than *ssthresh*, the sender increments the value of *cwnd* by the value $SMSS \times SMSS/cwnd$, in response to each received nonduplicate ACK^[7], ensuring that the congestion window opens by one segment within each RTT time interval.

The congestion window continues to open in this fashion until packet loss occurs. If the packet loss is isolated to a single packet within a packet sequence, the resultant duplicate ACKs will trigger the sender to halve the sending rate and continue a linear growth of the congestion window from this new point, as described above in fast recovery.

The behavior of *cwnd* in an idealized configuration is shown in Figure 8, along with the corresponding data-flow rates. The overall characteristics of the TCP algorithm are an initial relatively fast scan of the network capacity to establish the approximate bounds of maximal efficiency, followed by a cyclic mode of adaptive behavior that reacts quickly to congestion, and then slowly increases the sending rate across the area of maximal transfer efficiency.

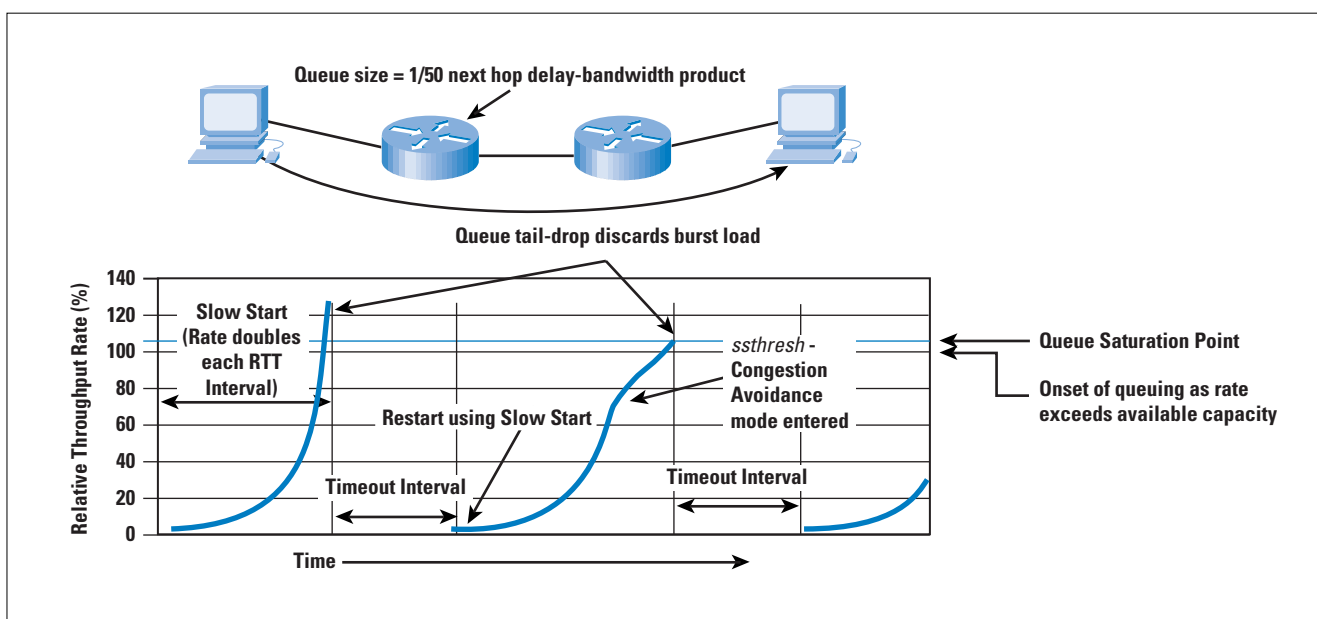
Packet loss, as signaled by the triggering of the retransmission timer, causes the sender to recommence slow-start mode, following a timeout interval. The corresponding data-flow rates are indicated in Figure 9.

Figure 8: Simulation of Single TCP Transfer



The inefficiency of this mode of performance is caused by the complete cessation of any form of flow signaling from the receiver to the sender. In the absence of any information, the sender can only assume that the network is heavily congested, and so must restart its probing of the network capacity with an initial congestion window of a single segment. This leads to the performance observation that any form of packet-drop management that tends to discard the trailing end of a sequence of data packets may cause significant TCP performance degradation, because such drop behavior forces the TCP session to continually time out and restart the flow from a single segment again.

Figure 9: Simulation of TCP Transfer with Tail Drop Queue

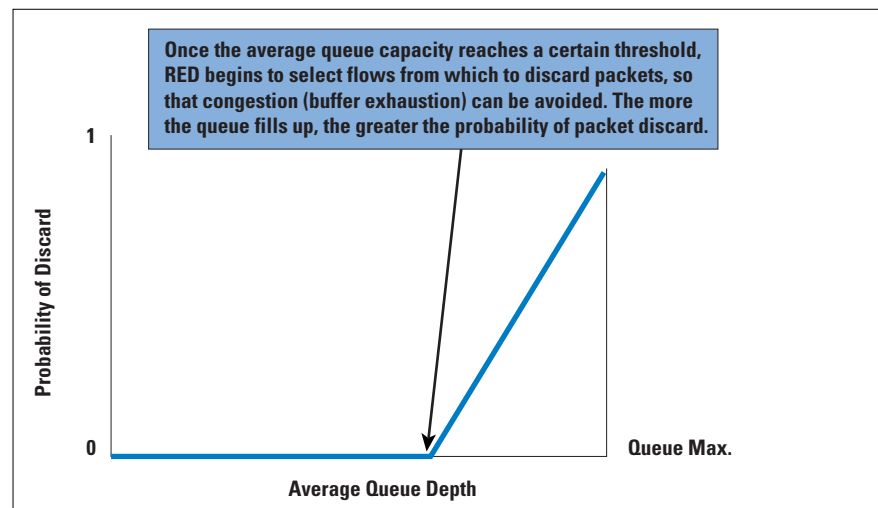


Assisting TCP Performance within the Network—RED and ECN

Although TCP is an end-to-end protocol, it is possible for the network to assist TCP in optimizing performance. One approach is to alter the queue behaviour of the network through the use of *Random Early Detection* (RED). RED permits a network router to discard a packet even when there is additional space in the queue. Although this may sound inefficient, the interaction between this early packet-drop behaviour and TCP is very effective.

RED uses a the weighted average queue length as the probability factor for packet drop. As the average queue length increases, the probability of a packet being dropped, rather than being queued, increases. As the queue length decreases, so does the packet-drop probability. (See Figure 10). Small packet bursts can pass through a RED filter relatively intact, while larger packet bursts will experience increasingly higher packet-discard rates. Sustained load will further increase the packet-discard rates. This implies that the TCP sessions with the largest open windows will have a higher probability of experiencing packet drop, causing a back-off in the window size.

Figure 10: RED Behavior



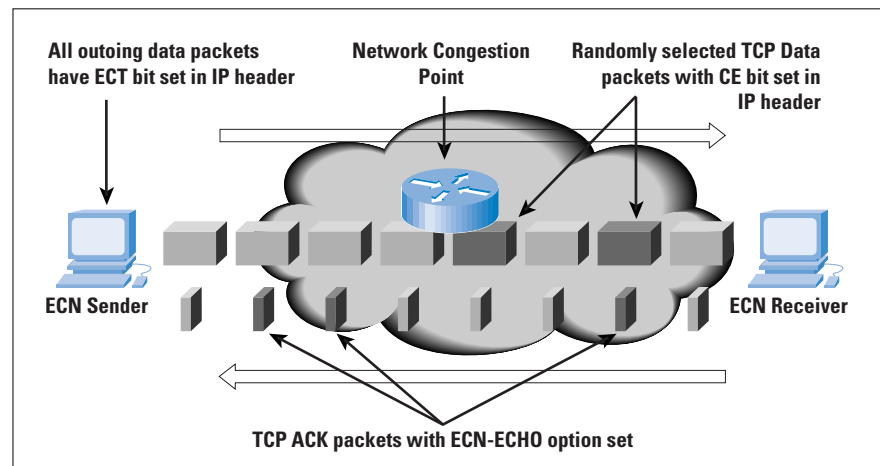
A major goal of RED is to avoid a situation in which all TCP flows experience congestion at the same time, all then back off and resume at the same rate, and tend to synchronize their behaviour^[11,12]. With RED, the larger bursting flows experience a higher probability of packet drop, while flows with smaller burst rates can continue without undue impact. RED is also intended to reduce the incidence of complete loss of ACK signals, leading to timeout and session restart in slow-start mode. The intent is to signal the heaviest bursting TCP sessions the likelihood of pending queue saturation and tail drop before the onset of such a tail-drop congestion condition, allowing the TCP session to undertake a fast retransmit recovery under conditions of congestion avoidance. Another objective of RED is to allow the queue to operate efficiently, with the queue depth ranging across the entire queue size within a timescale of queue depth oscillation the same order as the average RTT of the traffic flows.

Behind RED is the observation that TCP sets very few assumptions about the networks over which it must operate, and that it cannot count on any consistent performance feedback signal being generated by the network. As a minimal approach, TCP uses packet loss as its performance signal, interpreting small-scale packet-loss events as peak load congestion events and extended packet loss events as a sign of more critical congestion load. RED attempts to increase the number of small-scale congestion signals, and in so doing avoid long-period sustained congestion conditions.

It is not necessary for RED to discard the randomly selected packet. The intent of RED is to signal the sender that there is the potential for queue exhaustion, and that the sender should adapt to this condition. An alternative mechanism is for the router experiencing the load to mark packets with an explicit *Congestion Experienced* (CE) bit flag, on the assumption that the sender will see and react to this flag setting in a manner comparable to its response to single packet drop^{[13] [14]}. This mechanism, *Explicit Congestion Notification* (ECN), uses a 2-bit scheme, claiming bits 6 and 7 of the IP Version 4 *Type-of-Service* (ToS) field (or the two *Currently Unused* [CU] bits of the IP *Differentiated Services* field). Bit 6 is set by the sender to indicate that it is an ECN-capable transport system (the ECT bit). Bit 7 is the CE bit, and is set by a router when the average queue length exceeds configured threshold levels.

The ECN algorithm is that an active router will perform RED, as described. After a packet has been selected, the router may mark the CE bit of the packet if the ECT bit is set; otherwise, it will discard the selected packet. (See Figure 11).

Figure 11: Operation of Explicit Congestion Notification



The TCP interaction is slightly more involved. The initial TCP SYN handshake includes the addition of ECN-echo capability and *Congestion Window Reduced* (CWR) capability flags to allow each system to negotiate with its peer as to whether it will properly handle packets with the CE bit set during the data transfer. The sender sets the ECT bit in all packets sent. If the sender receives a TCP packet with the ECN-echo flag set in the TCP header, the sender will adjust its congestion window as if it had undergone fast recovery from a single lost packet.

The next sent packet will set the TCP CWR flag, to indicate to the receiver that it has reacted to the congestion. The additional caveat is that the sender will react in this way at most once every RTT interval. Further, TCP packets with the ECN-echo flag set will have no further effect on the sender within the same RTT interval. The receiver will set the ECN-echo flag in all packets when it receives a packet with the CE bit set. This will continue until it receives a packet with the CWR bit set, indicating that the sender has reacted to the congestion. The ECT flag is set only in packets that contain a data payload. TCP ACK packets that contain no data payload should be sent with the ECT bit clear.

The connection does not have to await the reception of three duplicate ACKs to detect the congestion condition. Instead, the receiver is notified of the incipient congestion condition through the explicit setting of a notification bit, which is in turn echoed back to the sender in the corresponding ACK. Simulations of ECN using a RED marking function indicate slightly superior throughput in comparison to configuring RED as a packet-discard function.

However, widespread deployment of ECN is not considered likely in the near future, at least in the context of Version 4 of IP. At this stage, there has been no explicit standardization of the field within the IPv4 header to carry this information, and the deployment base of IP is now so wide that any modifications to the semantics of fields in the IPv4 header would need to be very carefully considered to ensure that the changed field interpretation did not exercise some malformed behavior in older versions of the TCP stack or in older router software implementations.

ECN provides some level of performance improvement over a packet-drop RED scheme. With large bulk data transfers, the improvement is moderate, based on the difference between the packet retransmission and congestion-window adjustment of RED and the congestion-window adjustment of ECN. The most notable improvements indicated in ECN simulation experiments occur with short TCP transactions (commonly seen in Web transactions), where a RED packet drop of the initial data packet may cause a six-second retransmit delay. Comparatively, the ECN approach allows the transfer to proceed without this lengthy delay.

The major issue with ECN is the need to change the operation of both the routers and the TCP software stacks to accommodate the operation of ECN. While the ECN proposal is carefully constructed to allow an essentially uncoordinated introduction into the Internet without negative side effects, the effectiveness of ECN in improving overall network throughput will be apparent only after this approach has been widely adopted. As the Internet grows, its inertial mass generates a natural resistance to further technological change; therefore, it may be some years before ECN is widely adopted in both host software and Internet routing systems. RED, on the other hand, has had a more rapid introduction to the Internet, because it requires only a local modification to router behavior, and relies on existing TCP behavior to react to the packet drop.

Tuning TCP

How can the host optimize its TCP stack for optimum performance? Many recommendations can be considered. The following suggestions are a combination of those measures that have been well studied and are known to improve TCP performance, and those that appear to be highly productive areas of further research and investigation^[1].

- *Use a good TCP protocol stack:* Many of the performance pathologies that exist in the network today are not necessarily the by-product of oversubscribed networks and consequent congestion. Many of these performance pathologies exist because of poor implementations of TCP flow-control algorithms; inadequate buffers within the receiver; poor (or no) use of path-MTU discovery; no support for fast-retransmit flow recovery, no use of window scaling and SACK, imprecise use of protocol-required timers, and very coarse-grained timers. It is unclear whether network ingress-imposed Quality-of-Service (QoS) structures will adequately compensate for such implementation deficiencies. The conclusion is that attempting to address the symptoms is not the same as curing the disease. A good protocol stack can produce even better results in the right environment.
- *Implement a TCP Selective Acknowledgment (SACK) mechanism:* SACK, combined with a selective repeat-transmission policy, can help overcome the limitation that traditional TCP experiences when a sender can learn only about a single lost packet per RTT.
- *Implement larger buffers with TCP window-scaling options:* The TCP flow algorithm attempts to work at a data rate that is the minimum of the delay-bandwidth product of the end-to-end network path and the available buffer space of the sender. Larger buffers at the sender and the receiver assist the sender in adapting more efficiently to a wider diversity of network paths by permitting a larger volume of traffic to be placed in flight across the end-to-end path.
- *Support TCP ECN negotiation:* ECN enables the host to be explicitly informed of conditions relating to the onset of congestion without having to infer such a condition from the reserve stream of ACK packets from the receiver. The host can react to such a condition promptly and effectively with a data flow-control response without having to invoke packet retransmission.
- *Use a higher initial TCP slow-start rate than the current 1 MSS (Maximum Segment Size) per RTT.* A size that seems feasible is an initial burst of 2 MSS segments. The assumption is that there will be adequate queuing capability to manage this initial packet burst; the provision to back off the send window to 1 MSS segment should remain intact to allow stable operation if the initial choice was too large for the path. A robust initial choice is two segments, although simulations have indicated that four initial segments is also highly effective in many situations.

- *Use a host platform that has sufficient processor and memory capacity to drive the network.* The highest-quality service network and optimally provisioned access circuits cannot compensate for a host system that does not have sufficient capacity to drive the service load. This is a condition that can be observed in large or very popular public Web servers, where the peak application load on the server drives the platform into a state of memory and processor exhaustion, even though the network itself has adequate resources to manage the traffic load.

All these actions have one thing in common: They can be deployed incrementally at the edge of the network and can be deployed individually. This allows end systems to obtain superior performance even in the absence of the network provider tuning the network's service response with various internal QoS mechanisms.

Conclusion

TCP is not a predictive protocol. It is an adaptive protocol that attempts to operate the network at the point of greatest efficiency. Tuning TCP is not a case of making TCP pass more packets into the network. Tuning TCP involves recognizing how TCP senses current network load conditions, working through the inevitable compromise between making TCP highly sensitive to transient network conditions, and making TCP resilient to what can be regarded as noise signals.

If the performance of end-to-end TCP is the perceived problem, the most effective answer is not necessarily to add QoS service differentiation into the network. Often, the greatest performance improvement can be made by upgrading the way that hosts and the network interact through the appropriate configuration of the host TCP stacks.

In the next article on this topic, we will examine how TCP is facing new challenges with increasing use of wireless, short-lived connections, and bandwidth-limited mobile devices, as well as the continuing effort for improved TCP performance. We'll look at a number of proposals to change the standard actions of TCP to meet these various requirements and how they would interact with the existing TCP protocol.

References

- [1] Huston, G., *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, John Wiley & Sons, January 2000.
- [2] Postel, J., "Transmission Control Protocol," RFC 793, September 1981.
- [3] Jacobson, V., Braden, R., and Borman, D., "TCP Extensions for High Performance," RFC 1323, May 1992.
- [4] Mathis, M., Madavi, J., Floyd, S., and Romanow, A., "TCP Selective Acknowledgement Options," RFC 2018, October 1996.

- [5] Nagle, J., "Congestion Control in IP/TCP Internetworks," RFC 896, January 1984.
- [6] Allman, M., Floyd, S., and Partridge, C., "Increasing TCP's Initial Window," RFC 2414, September 1998.
- [7] Allman, M., Paxson, V., and Stevens, W., "TCP Congestion Control," RFC 2581, April 1999.
- [8] Stevens, W. R., *TCP/IP Illustrated, Volume 1*, Addison-Wesley, 1994.
- [9] Jacobson V., "Congestion Avoidance and Control," *ACM Computer Communication Review*, Vol. 18, No. 4, August 1988.
- [10] Jacobson, V., "Berkeley TCP Evolution from 4.3-Tahoe to 4.3, Reno," Proceedings of the 18th Internet Engineering Task Force, University of British Columbia, Vancouver, BC, September 1990.
- [11] Floyd, S., and Jacobson, V., "Random Early Detection Gateways for Congestion Avoidance," *IEEE/ACM Transactions on Networking*, Vol. 1, No. 4, August 1993.
- [12] Braden, R. et al., "Recommendations on Queue Management and Congestion Avoidance in the Internet," RFC 2309, April 1998.
- [13] Floyd, S., "TCP and Explicit Congestion Notification," *ACM Computer Communication Review*, Vol. 24, No. 5, October 1994.
- [14] Ramakrishnan, K., and Floyd, S., "A Proposal to Add Explicit Congestion Notification (ECN) to IP," RFC 2481, January 1999.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. Huston is currently the Chief Technologist in the Internet area for Telstra. He is also an active member of the IETF, and is the chair of the Internet Society Board of Trustees. He is author of *The ISP Survival Guide*, ISBN 0-471-31499-4, *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, and coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. All three books are published by John Wiley & Sons. E-mail: gih@telstra.net

Overview of Internet Mail Standards

by Paul Hoffman, Internet Mail Consortium

People who are new to the Internet often think it is equivalent to “the Web” since that’s what they have heard about most in the media. After a few weeks of using their new Internet account, they tend to say the Internet is “e-mail and the Web,” in that order.

Business users have an even higher regard for e-mail. According to the American Management Association, most business people say that e-mail has surpassed the telephone in importance for business communication. While many companies believe that their Web site will be very important in a few years, their e-mail system is already extremely critical to them today.

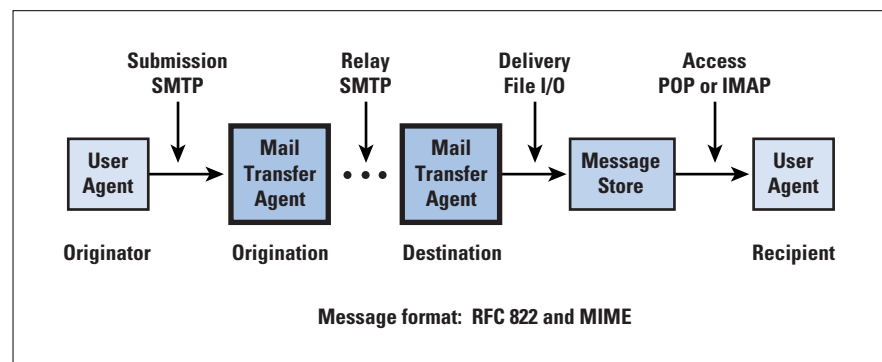
Because mail is one of the oldest services on the Internet, the protocols used to move mail around are more stable and mature than those used for newer services. The flip side of this is that some of the protocols that are used to move the billions of pieces of mail a day are somewhat arcane and even quaint. The *Internet Engineering Task Force* (IETF) motto “if it isn’t broken, don’t fix it!” has prevented people from redesigning Internet mail. Instead, numerous extensions and enhancements have been added to the original set of mail standards, as we shall see later.

Historically, there have been many other mail systems, such as BITNET, Fidonet, MAPI, cc:Mail, and so on. Of course, users of these systems still exist, and there is quite an active market for systems that act as gateways between Internet mail and other systems. However, this article only covers the tried-and-true Internet mail system.

The Internet Mail Model

Many Internet protocols are simple client/server systems with a single message payload format. Mostly due to history, Internet mail doesn’t have this luxury. Figure 1 shows the main protocols and formats used to move Internet mail.

Figure 1: Internet Mail Architecture



A typical mail transaction goes from left to right in the figure. A *Mail User Agent* (MUA), which is most often run by a human but could be controlled by a program, submits a message using the *Simple Mail Transfer Protocol* (SMTP) to the initiating host. That host looks up the IP address associated with the destination host computer and sends the message to the destination host using SMTP. The destination host receives the message and writes it into the local message store (which almost always is a file or database on a hard drive).

The recipient MUA checks the mail store periodically and, if there is mail, retrieves it. Again, the recipient is often a human but might be an automated program such as an order entry system that is controlled by e-mail. The protocols that check for and retrieve mail are usually the *Post Office Protocol* (POP) or *Internet Message Access Protocol* (IMAP), but it could also be any number of proprietary systems.

E-mail messages have a format that is quite easy to understand, so much so that many other protocols have adopted very similar formats. The message consists of ASCII text *headers* followed by one ASCII (or possibly binary) message *body*. The header format is defined in RFC 822^[1], thus the headers are called “RFC 822 headers” or just “822 headers.” A simple message body is a single string of text; a complex body uses the *Multipurpose Internet Mail Extensions* (MIME) message format.

Moving Between Hosts: SMTP

Early host-to-host mail delivery was done using file transfer protocols. Since such methods offer little flexibility and require knowledge of user names and file structures of the remote system, a more general purpose delivery mechanism evolved. The resulting protocol, SMTP was defined in 1982 and has proven to work effectively in the face of orders of magnitude increase in the size of the network.

Even though SMTP moves mail between two host computers, it is a client/server protocol. The host that initiates the contact always acts as a *client*, and the host that was contacted is the *server*. (There are a few rarely-used exceptions to this rule.) The client has a variety of text-based commands that it can give, and the server replies with short responses. The server in the relationship never gives commands on its own, so it is up to the client to ask enough questions, and to carefully watch the server’s responses, to know how best to interact with the server.

When a host wants to send mail somewhere on the Internet, it determines where the mail should go and initiates contact with the target server. Thus, the sender is always the SMTP client, and the hosts that are listening for SMTP traffic are always servers. In reality, most SMTP server software can act both as clients and servers; MUAs almost always only participate in SMTP as clients.

The sending host uses the *Domain Name System* (DNS) to determine the IP address of the target host, contacts that host using TCP port 25, and uses SMTP to deliver the message. Sometimes, as we shall see later, this IP lookup involves a level of indirection,—the target host may use a different host to receive mail on its behalf.

SMTP client commands consist of a keyword, possibly followed by command arguments. The server’s response always starts with a three-digit number that is a status indicator, which is possibly followed by additional information.

Most SMTP interactions follow a typical set of steps, shown in Figure 2. The initiator who has mail to send (the client) is on the left and the host that is receiving the mail (the server) is on the right. The client first opens a TCP connection on port 25 on the the server. Next, the client and server exchange greetings (the HELO command and response). The client then prepares the server to receive the message by telling the server who the message is from and who it is to; the server gives a positive acknowledgment to each of these commands. The client then asks if the server is ready for the body of the message and, when the server says yes, sends the message as a stream of lines that is followed by a single period on a line by itself. After the server says that it has received the message fully, the client says good-bye and closes the connection.

Figure 2: A Typical SMTP Exchange

Client Action	Server Response
<connects to the server using TCP>	
HELO somecompany.com	220 example.com WhizzyMail server version 2.32
MAIL FROM:<chrisj@somecompany.com>	250-example.com says howdy
	250 OK
Text of message...	354 Start mail input; end with <CRLF>.<CRLF>
...more text of message...	
.	
QUIT	250 OK
<disconnects from the server>	221 example.com Service finished

Submission and Relay

After a message is created, the creator uses SMTP to submit the message to one of two places: a local mail-forwarding host (such as the mail server at the sender’s *Internet Service Provider* (ISP) or corporate IS services) or the mail server that the DNS says is definitive for the recipient. The former is typically used by Internet users who do not have persistent network connections; the latter is more common on systems with network connections that are always available.

Messages may be forwarded hop-by-hop from the sending host, via intermediary hosts, to the recipient. This is called “relaying.” In many cases, a message will go through more than two relays, for instance when the recipient’s network is configured to accept all incoming messages on one machine that later relays messages to individual departmental hosts. Note that submission and relay uses the same SMTP commands described above (a recent change to this scheme is described near the end of this article).

The last host in the chain makes the message available to the recipient. This is done by moving the message to the message store, which usually means “write the message out on disk.” There are, of course, many ways to write something on disk; some hosts write out each message as a separate file, some concatenate the message at the end of a file, while others write the message into a database.

Mail Addresses and MX Records

The initiating host’s first job is to determine where a message is supposed to go, that is, how to contact the recipient’s host. SMTP is a hop-by-hop protocol, meaning that a sending host does not know the true destination host for a message: it only knows the designated recipient host. Of course, this might be the recipient’s final host, or it might be a host that will pass the message along further.

The domain name in mail addresses do not necessarily correspond to hosts on the Internet. For example, there is no host whose domain name is **imc.org**. When determining where to send a message, the initiating host first looks in the DNS for a *Mail Exchange* (MX) record that matches the domain name in the recipient’s mail address. If there is no MX record, the initiating host looks for a DNS A record that matches the domain name. If there is no MX record or A record, the message cannot be delivered.

Many people find MX records to be somewhat tricky. Part of the confusion comes from the fact that an SMTP host is supposed to look up MX records before they look for A records; there are very few protocols that don’t rely on A records. Another confusing aspect is that MX records may have wildcards in them. For instance, if a message is being sent to **someone@eng.example.com**, there may be no MX record for **eng.example.com**, but there may be one for ***.example.com**. Wildcard MX records tell the sending host that any message for a domain name that matches the wildcard specification should be sent to the named host.

Modern Mail Extensions

All protocols must evolve, and SMTP has improved over the years. Early mail implementors realized that the initial set of SMTP commands would have to expand. Since the SMTP client gives all commands in an exchange, the client determines which SMTP commands a server will be able to handle. The *SMTP Service Extensions* (ESMTP), defined in RFC 1869^[2], is a small change to SMTP that allows an SMTP server to list the commands it knows at the beginning of an SMTP session.

The bootstrapping process for ESMTP is quite simple. Instead of starting with the “HELO” command, an ESMTP server starts with the “EHLO” command. If the SMTP host indicates that it has no idea what “EHLO” means, the client knows that the server doesn’t understand ESMTP, and therefore doesn’t understand any SMTP extensions. On the other hand, if the server does understand the “EHLO” greeting, the host responds with the entire list of SMTP extensions that the client is allowed to use during the session.

There have been over a dozen extensions to SMTP that are on standards track in the IETF, and many more have been proposed. However, most modern SMTP servers have only implemented a few of these.

Probably the most publicized SMTP extension in the past few years has been the *SMTP Service Extension for Authentication* (AUTH) for authenticating the SMTP client to the server. The AUTH extension, described in RFC 2554^[3], allows roaming users to submit mail from outside their local networks without forcing the servers to accept mail from just anyone. This new method, which is now starting to appear in both mail clients and servers, will reduce the hassle faced by many roaming users as they move from ISP to ISP.

Another significant SMTP extension that has become widely implemented is *Delivery Status Notifications*, or DSNs defined in RFC 1891^[4]. These are similar to return receipts in postal mail, but with some significant differences. DSNs are issued by SMTP servers, not end users. Thus, the meaning of a DSN is interpreted as “the message was received by this SMTP host,” not “the message was received by the intended recipient.”

Retrieving Mail

After the final SMTP server has received a message and written it into the message store, the recipient needs to be able to access the message. In the early days of Internet mail, the message store was nothing more than a text file on disk, and mail was read by reading the text file. In fact, many people still read their mail this way, albeit using somewhat more modern tools.

If the recipient is not directly logged into the host computer that has the message store, reading the disk file can be difficult. To alleviate this problem, the *Post Office Protocol* (POP), described in RFC 1939^[5] introduced a client/server model for an MUA to get mail from the message store and store it on the local computer. The vast majority of mail users today use POP to retrieve their mail.

POP looks like many Internet protocols. The client connects to the server, logs in using a user name and password, checks if it has any messages waiting for it, then asks for the messages one by one. The client has the option of leaving messages that it has read on the server or deleting them after they have been retrieved.

Modern Mail Access with IMAP

Although POP works well for many people, it has its drawbacks. The mail client cannot preview a message to see whether or not it wants to download it. The client has only one mailbox which has no hierarchical structure. In most POP systems, leaving all your mail on the server makes retrieving new mail quite slow. To get around these problems, the mail community developed the *Internet Message Access Protocol* (IMAP), described in RFC 2060^[6].

IMAP is significantly more powerful than POP. IMAP clients give the user much more control over their mail, such as letting them keep some of their mail locally while leaving other mail on the server. IMAP even allows for mailboxes that are shared among users, such as group announcements lists. It also gives mail administrators many more opportunities to support novice users by keeping their mail in a central location. Most modern mail clients support IMAP, and IMAP servers are available from many vendors.

It should be noted that, although IMAP is considered much more useful than POP and is widely available, it has had very little adoption in the ISP market (it has been accepted much more readily in the enterprise mail market). The reasons for this are not clear. Many ISPs say they do not want to incur the costs and responsibilities of storing users' mail, even if this gives them greater ability to administer the mail. It is not clear what, if anything, will shift ISPs away from POP to IMAP.

Access Through Web Browsers

The ubiquity of the Web has introduced a new method for getting mail that has become surprisingly popular: the use of the *HyperText Transfer Protocol* (HTTP). Web access to e-mail lets users read their mail without a POP or IMAP client. Of course, this offers many fewer features than POP or IMAP; for instance, you can't easily store messages after reading them and getting file attachments in your mail takes many more steps. However, the big advantage of this method is that Web browsers are almost everywhere these days, and there are many situations where you don't care about being able to store your mail on your local computer.

Giving users Web browser access to their mail quickly became a commodity market. Now, almost every portal offers such services. In fact, many corporations and ISP also offer this service because it is a fairly easy add-on to POP and IMAP servers. As more and more users want to access their e-mail from small devices such as cellular phones, it is likely that these devices will include Web-like mail interfaces.

Client Extensions

Both POP and IMAP are extensible, and developers have proposed many extensions for both protocols, although most work is being done on IMAP. Because of the slow adoption of IMAP by ISPs (who could make its advantages much more visible), it's not clear when these will appear in clients and servers, even though many of them add interesting functionality that is wanted by both users and administrators.

There are many client extensions that don't rely on either POP or IMAP, however. One of the most popular is *Message Disposition Notifications* (MDNs), which are quite similar to postal return receipts. Unlike DSNs, which say that a particular message got to one of the servers in an SMTP chain, MDNs are truly end-to-end, and are returned by recipients when they open their mail.

Some people find MDNs intrusive (“why should he know when I read this?”), and they aren’t particularly reliable because not all mail clients (most notably Web browser readers) support them. However, they are a good example of what end users are seeing in terms of extensions that add desired functions to the Internet mail system.

The Format of Mail Messages

SMTP, POP, and (to a great extent) IMAP ignore the contents of a message. SMTP uses its own control information to find the recipient of a message; POP and IMAP retrieve messages based on user account names, which may or may not correspond to the address in a message. In users minds, however, the contents of the messages they read are almost always much more important than the way that the message got to them.

Mail messages consist of two parts: the *headers* and the *body*. The headers come first, followed by a blank line, followed by the body, as shown in Figure 3. The basic structure of messages has remained unchanged since it was defined in RFC 822. Originally, the headers were designed to look like inter-office memos, and also to contain control and debugging information; today, some parts of the headers are considered to be as important as the body of the message.

Figure 3: A Typical
E-mail Message

```
Received: from mail.somecompany.com ([198.81.17.2])
  by mail.example.com (8.8.8/8.8.5) with ESMTP id VAA17989
  for <althea@example.com>; Wed, 9 Dec 1998 21:07:44 -0800 (PST)
From: jerry@somecompany.com
Message-ID: <823227a3.366f53ea@somecompany.com>
Date: Wed, 9 Dec 1998 23:54:02 EST
To: Althea Cassidy <althea@example.com>
Mime-Version: 1.0
Subject: I'm outta here
Content-type: text/plain; charset=US-ASCII
Content-transfer-encoding: 7bit

Sorry to make this so brief, but I've got a train to catch.
I'll meet you at the jubilee.
--J
```

Message Headers

Because they were designed to be functional, message headers have a very straight-forward design. Each header has a single token, followed by a colon, followed by the parameters and options of the header. Headers usually consist of a single line, but you can create multi-line headers by starting the continuation lines with blanks.

There are dozens of common headers, and dozens more that are rarely used. Almost all mail users are familiar with “To:”, “From:”, “Subject:”, and “Date:”, and they may have seen additional common headers such as “Cc:” and “Received:”. Depending on the interface of the MUA, users typically see some of these headers after they have retrieved a message with POP or IMAP but before they have “opened” the message to see the message body.

Basic and Advanced Message Bodies

Originally, the body of mail messages consisted of plain ASCII text. This was sufficient for the inventors of e-mail, who spoke mostly English and had access to other information transfer mechanisms such as FTP to move binary data around. Of course, such restrictions would not last.

Probably the biggest advance in Internet mail in the past ten years is in the format of mail messages, not in their transport. In the early 1990s, Internet mail went from being text-only to allowing the transfer of non-text messages and parts of messages. MIME, described in RFCs 2045–2047^[7, 8, 9], revolutionized the usefulness of Internet mail by allowing senders to include files with messages, to use styled text, to give their messages useful structure, and to provide the first interoperable support for international e-mail.

Unfortunately, the term “attachments” became associated with MIME even though it is much more powerful than just allowing files to be attached to a message. The majority of MIME-enabled messages today don’t contain any attachments: instead, they use MIME’s capability of labeling the type of a single message body part. MIME labeling can tell the receiving client the format of the message (for instance, an HTML message) and, if it is a text message, the type of characters in the message.

Another great feature of MIME is that it allows messages to have structure. For instance, Figure 4 shows a message with two representations of the same information: text and HTML. A mail client that cannot display HTML can skip that part of the message and just display the plain text. This allows message content to gradually migrate towards new technology. In the near future, it is likely that similar logic will be used for messages that contain XML, HTML, and plain text.

Figure 4: A Multipart MIME message

```
From: jerry@somecompany.com
Message-ID: <828d83ffzwd.47r7c2dxsa@somecompany.com>
Date: Wed, 10 Dec 1998 03:24:00 EST
To: Althea Cassidy <althea@example.com>
Subject: What you should know
MIME-Version: 1.0
Content-Type: multipart/alternative;
    boundary=ad8ekd2ddr9332dc3df332

--ad8ekd2ddr9332dc3df332
Content-Type: text/plain; charset=us-ascii

Important stuff.
Blah blah blah.

--ad8ekd2ddr9332dc3df332
Content-Type: text/html

<html><head><title>Important stuff</title></head><body>
<h1>Important stuff.</h1>
<p><b>Blah blah blah.</b>
</body></html>

--ad8ekd2ddr9332dc3df332--
```

Because of the capability to structure messages, MIME can be used for multimedia and unified messaging. A single mail message can contain one or more movies, sound files, text files in a variety of formats, binary files such as word processing documents, calendar events, fax images, and so on. The MIME structure tells the recipient software which parts of the message contain particular types of data, as well the relationship between the parts (such as “this part contains three different alternative sound formats”).

With the explosion of the popularity of the Web, users have come to expect that the content they read will look like Web pages. Most users don’t understand that a “Web page” that “contains” graphics in fact isn’t a single entity but is really a page of HTML that has links to other pages that contain individual images. They expect to be able to receive mail messages that look just like the things they see on the Web. The MHTML protocol (described in RFC 2557^[10]) describes how to structure MIME messages that contain both HTML parts and images so that they appear together in mail clients exactly like they appear in Web browsers.

MIME enables a plethora of other uses for e-mail. For example, secure e-mail using S/MIME and PGP uses MIME to structure the messages so that the cryptographic control information is separate from the message itself. For instance, in a digitally-signed message, the signature information (which is unreadable to the human recipient) is in a different part of the structure than the human-readable content. You can even have layers of encryption and signatures, all structured through MIME.

Internationalization of E-mail

You can use character sets other than ASCII in both the headers and body of Internet e-mail messages. Using different character sets in text bodies requires the use of the “charset” parameter in the “Content-type:” header, as described in RFC 2046^[8]. You can also use character sets in message headers with the methods described in RFC 2047^[9].

The Future of E-mail

E-mail is incredibly popular with Internet users, but it is far from finished. The next billion new e-mail users will most likely be much less technically savvy than today’s Internet users, and they will come to the Internet with very different expectations. In order to give these users a more pleasant experience, the Internet mail industry will have to add many new features and make mail clients easier.

The number of ISPs is also increasing, although not as fast as the number of Internet users. Since e-mail is such an integral part of the service that an ISP offers, mail server software will also have to become easier to administer. Internet mail server vendors are working on such enhancements as a way of gaining a competitive advantage.

The most major change that users will see in the next few years are more highly enabled MUAs. These clients will be all-in-one messaging centers that will handle faxes, voice messages, paging, calendar and event management, and probably some sort of instant messaging. In this way, traditional mail will be only one part of what the user sees when they go to their messaging client.

The importance of Internet fax should not be underestimated. The recent standards for Internet fax, defined in RFCs 2301–2306^[11, 12, 13, 14, 15, 16] specify how faxes go through Internet mail. Although there have been a raft of proprietary real-time fax proposals, fax vendors have rallied around faxes in e-mail as an easy way to transition from fax over phone lines. Comparing the high cost of sending international faxes to the near-zero cost of sending e-mail, many companies are quickly moving towards the new standards.

Other mail-enabled services are becoming standardized as well. For example, calendaring over Internet mail is nearing completion. This will allow users to coordinate schedules for meetings, even with people who are not online. E-mail fall-back for phone conversations that were not completed is also being researched.

The e-mail world five and ten years from now will not necessarily look completely different from the way it looks today. Certainly, there will be many more enriched text and multimedia messages being composed by end users. Mailing lists will grow and the mail on them will be more like Web pages than today's text messages. Many people predict that the face of e-mail will change radically if e-mail becomes the "universal inbox" for voicemail, faxes, and other types of communication. Many companies are discovering that regular newsletters sent through e-mail are more effective than expecting users to come to a web site regularly, and it is likely that there will be an increase in the number of publications that are delivered as e-mail.^[18]

There is still plenty of room for additions to Internet mail that resemble today's non-Internet services. For instance, users are already clamoring for features such as true message tracking, which is currently available from many package delivery services. Better security is clearly desired, although there seems to be major impediments caused by the need for trusted certificates before we can see wide deployment of secure mail. More problematic features such as message rescinding also have been proposed.

Forces outside the Internet mail world will also change how Internet mail works. For instance, the rapid increase in wireless users will change the way that large messages are handled by message stores. As more users start reading their mail from more than one system, IMAP may become more popular. At the same time, users will expect to be able to move their configuration information with them from machine to machine, probably using protocols such as the *Application Configuration Access Protocol* (ACAP) defined in RFC 2244^[17].

There are plenty of opportunities in the Internet mail market. The only significant dark cloud is the possibility that increasing unsolicited e-mail—so called “spam”—might scare away users. To date, the technical solutions for battling spam have been limited, and they probably won’t scale well if the amount of spam increases by an order of magnitude. On the bright side, it appears that most legitimate marketers have been scared away from spam and are focusing on opt-in e-mail marketing. This could be a boon for ISPs who specialize in bringing interested e-mail users and potential advertisers together.^[19]

In such an environment, mail with rich media and lots of convenience could become the place where many users want to spend much of their time. To get there, we need to build on today’s well-established mail protocols and to be creative in the kinds of features we add to both the transport and display of e-mail. Fortunately, we don’t need to do much with SMTP, IMAP, and MIME in order to bring these new capabilities to the burgeoning numbers of new users waiting to get on the Internet.

References

- [1] Crocker, D., “Standard for the format of ARPA Internet text messages,” RFC 822, August 1982.
- [2] Klensin, J., Freed, N., Rose, M., Stefferud, E., Crocker, D., “SMTP Service Extensions,” RFC 1865, November 1995.
- [3] Myers, J., “SMTP Service Extension for Authentication,” RFC 2554, March 1999.
- [4] Moore, K., “SMTP Service Extension for Delivery Status Notifications,” RFC 1891, January 1996.
- [5] Myers, J. and Rose, M., “Post Office Protocol—Version 3,” RFC 1939, May 1996.
- [6] Crispin, M., “Internet Message Access Protocol—Version 4rev1,” RFC 2060, December 1996.
- [7] Freed, N. and Borenstein, N., “Multipurpose Internet Mail Extensions (MIME) Part One: Format of Internet Message Bodies,” RFC 2045, November 1996.
- [8] Freed, N. and Borenstein, N., “Multipurpose Internet Mail Extensions (MIME) Part Two: Media,” RFC 2046, November 1996.
- [9] Moore, K., “MIME (Multipurpose Internet Mail Extensions) Part Three: Message Header Extensions for Non-ASCII Text,” RFC 2047, November 1996.
- [10] Palme, J., Hopmann, A., Shelness, N., “MIME Encapsulation of Aggregate Documents, such as HTML (MHTML),” RFC 2557, March 1999.
- [11] McIntyre, L., Zilles, S., Buckley, R., Venable, D., Parsons, G., Rafferty, J., “File Format for Internet Fax,” RFC 2301, March 1998.

- [12] Parsons, G., Rafferty, J., Zilles, S., “Tag Image File Format (TIFF)—image/tiff MIME Sub-type Registration,” RFC 2302, March 1998.
- [13] Allocchio, C., “Minimal PSTN address format in Internet Mail,” RFC 2303, March 1998.
- [14] Allocchio, C., “Minimal FAX address format in Internet Mail,” RFC 2304, March 1998.
- [15] Toyoda, K., Ohno, H., Murai, J., Wing, D., “A Simple Mode of Facsimile Using Internet Mail,” RFC 2305, March 1998.
- [16] Parsons, G., Rafferty, J., “Tag Image File Format (TIFF)—F Profile for Facsimile,” RFC 2306, March 1998.
- [17] Newman, C., Myers J. G., “ACAP—Application Configuration Access Protocol,” RFC 2244, November 1997.
- [18] *Poor Richard’s E-mail Publishing*, by Chris Pirillo, ISBN 0966103254, Top Floor Publishing, 1999.
- [19] *Internet Messaging: From the Desktop to the Enterprise*, by Marshall T. Rose and David Strom, ISBN 0-13-978610-4, Prentice Hall PTR, 1998.
- [20] *Essential Email Standards: RFCs and Protocols Made Practical*, by Pete Loshin, ISBN 0-471-34597-0, Wiley, 1999.
- [21] *Internet Email Protocols: A Developer’s Guide*, by Kevin Johnson, ISBN 0-201-43288-9, Addison-Wesley, 1999.

PAUL HOFFMAN is the director of the Internet Mail Consortium (<http://www.imc.org/>), which is the trade association for Internet mail software vendors and service providers. He is the editor of many recent mail standards, as well as Internet-related books such as *Netscape For Dummies*. He has been active on the Internet for twenty years. E-mail: phoffman@imc.org

Book Review

Introduction to Data Communications and Networking

Introduction to Data Communications and Networking, Behrouz Farouzan, ISBN 0-256-23044-7, WCB/McGraw-Hill, 1998.

As personal computers have proliferated the landscape over the years, they have become the domain of an increasing number of nontechnical end users. Two things assisted in this transformation. The realization of their value as a productivity tool became apparent, as well as their ability to become more user friendly to the masses. Networks, and networking, have followed a similar path. The investment in creating a networked environment in the past may have been a burden—in both time and added complexity—to all but the largest corporations. However, as the world becomes more “wired,” the presence of networks has become commonplace in nearly every work environment, not to mention the movement into private residences. The need to become familiar with concepts and terms as they relate to data communications and networks has become an important part of the technological landscape. *Introduction to Data Communications and Networking* assists the novice in grasping these concepts, as well as serving as a refresher to the more experienced audience.

Organization

The preface explains the ways this book can be useful. The textbook portion is helpful. Multiple choice as well as discussion questions are provided within each chapter, although all the answers are not. In addition, some of the questions asked do not always seem to be posed in the context of the chapter just covered. However, it does turn out to be a rather small inconvenience. The requisite appendices are included as well—such as ASCII and EBCDIC codes, and various representations of numbers. However, two areas that usually receive only fleeting recognition—Fourier analysis and Huffman coding—are covered. Not being an engineer, I’m not sure that I now understand these concepts, but at least now I know why.

Although the areas covered in this book are covered in many introductory network books, this one takes nothing for granted. A good portion of the more experienced readers will know that Layers 2–6 of the OSI model have headers, only Layer 2 will include a trailer. Details such as these are easily forgotten. Introducing concepts in meaningful, practical ways is another positive attribute of this book. One great example is how the author describes the difference between analog and digital. Hands of a traditional, or analog, clock do not jump from minute to minute or hour to hour. The notion of time advancing seems to be a smooth transition, much like an analog signal is a continuous wave form that changes smoothly over time. Digital (as in the case of a digital clock), on the other hand, indicates discrete units of time—usually whole hours and minutes—and can have only limited numbers of defined values. In Chapter 4, analog and digital signals are detailed and explained with clarity and excellent examples are given as well.

In fact, the only subject matter I had difficulty deciphering concerned material presented in Chapter 5. The concepts of polar, unipolar, and bipolar encoding seemed straightforward enough, but digital-to-analog and analog-to-analog encoding will definitely have to be revisited. Amplitude and phase shifting keys may or may not be revisited. In fact, it was at this point that I realized that the material was moving to a different, more difficult, level.

Although the preface states that the first eight chapters are essential for readers being introduced to networking concepts, I found that chapters 5–8 went into a level of depth that would be particularly daunting for an introductory discussion.

Summary

I don't remember exactly how I was introduced to this book—whether I read about it in a journal or it was recommended by a friend—but the book got favorable reviews wherever I inquired about it. It is a practical addition to your bookshelf, regardless of your level of comfort with networks and voice/data communications.

The book is relevant and practical for the professional who has been working in the field for a few years. It is also useful as a textbook for use in the classroom. However, I do not believe that all the information can be adequately covered in a semester, as the author suggests. I believe one of the reasons I enjoyed this book was because of the way it explained ideas and concepts that were never used in any class I had ever taken. I recall promises of receiving a good, comprehensive background in these areas, yet years later I continue to struggle with some of the same concepts I've encountered in classes before. I found myself continually searching for a source that would provide me the information in a comprehensive, understandable fashion. I believe I have finally found it.

—Steve Barsamian, Cisco Systems
sbarsam@cisco.com

Would You Like to Review a Book for IPJ?

We receive numerous books on computer networking from all the major publishers. If you've got a specific book you are interested in reviewing, please contact us and we will make sure a copy is mailed to you. The book is yours to keep if you send us a review. We accept reviews of new titles, as well as some of the “networking classics.” Contact us at ipj@cisco.com for more information.

New Top-Level Domains Are Coming

For several years, there have been proposals to introduce new *generic top-level domains* (gTLDs) into the *Internet Domain Name System* (DNS). Although the introduction of gTLDs raises several issues that are of concern to various members of the Internet community, significant progress has been made recently toward achieving a consensus solution. The *Internet Corporation for Assigned Names and Numbers* (ICANN) Board of Directors is expected to consider adopting a policy to introduce new gTLDs at its meeting in July 2000. The *Names Council* has recommended to the ICANN Board that: "...a limited number of new top-level domains be introduced initially and that the future introduction of additional top-level domains be done only after careful evaluation of the initial introduction."

ICANN Announces CPR Institute as New Dispute Resolution Provider

ICANN recently announced that the *CPR Institute for Dispute Resolution* has been designated an approved provider under their *Uniform Dispute Resolution Policy* (UDRP) for domain name disputes. CPR, an alliance of 500 general counsel of global corporations and partners of major law firms, is the fourth dispute resolution provider to be designated by ICANN to handle domain disputes, joining the *National Arbitration Forum*, the *Disputes.org/Resolution Consortium*, and the *World Intellectual Property Organization*. The UDRP establishes a streamlined, economical process administered by neutral arbitration companies to provide a quick and cheap alternative to litigation. The procedure applies to cases that meet all three of the following criteria: The domain name must be identical or confusingly similar to a name in which the complaining party has trademark rights (either through a registered trademark or a common-law trademark); The domain name holder must have no legitimate right or interest in the name; The domain name must have been registered and used in bad faith.

In its first few months of operation, the UDRP has proven to be a very popular means of quickly resolving trademark/domain name disputes. To date, 691 proceedings have been commenced under the policy involving 1022 domain names. Of those proceedings, 348 have already been resolved. For additional information on UDRP, see <http://www.icann.org/udrp/udrp.htm>

This publication is distributed on an "as-is" basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Engineering
MCI WorldCom, USA

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Edward R. Kozel, Member of The Board of Directors
Cisco Systems, Inc., USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Strategy Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

Copyright © 2000 Cisco Systems Inc.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-10/5
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

Bulk Rate Mail U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol Journal

September 2000

Volume 3, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

F R O M T H E E D I T O R

In This Issue

From the Editor	1
The Future for TCP	2
Securing the Infrastructure.....	28
Book Reviews	45
Call for Papers	49
Fragments	50

In our last issue, Geoff Huston described the basic design and operation of the *Transmission Control Protocol* (TCP). He outlined how numerous enhancements to TCP implementations have been developed over time to improve its performance, particularly in the face of congested networks. The Internet is a rapidly changing environment in which both the applications and the underlying transmission systems are undergoing an evolution, if not a revolution. Some of these changes, such as the introduction of wireless devices, affect the way TCP works, because the protocol makes many implicit assumptions about the network over which it operates. In this issue, Geoff looks at the future for TCP and describes techniques for adopting TCP to today's Internet.

Security continues to be a major concern for everyone involved in the design and operation of networks. Widely publicized "hacker attacks," "denial-of-service attacks," and outright online fraud has brought the topic into sharp focus in the last few years. Because security was not part of the original design of the Internet, numerous solutions at every level of the protocol stack have been proposed and implemented over the last three decades. Today's network manager is, therefore, faced with a *system* of security components that must be carefully configured and monitored in order to provide sufficient security without preventing users from getting their work done. In our second article, Chris Lonvick explores a model for evaluating and securing a network.

The online subscription system for this journal is now up and running at www.cisco.com/ipj. In addition to offering a subscription form, the system allows you to select delivery options, update your mailing and e-mail address, and much more. Please visit our Web site and give it a try. If you encounter any difficulties, please send your comments to ipj@cisco.com.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

The Future for TCP

by Geoff Huston, Telstra

The previous article, “TCP Performance,” examined the operation of the *Transmission Control Protocol* (TCP) protocol^[1]. The article examined the role of TCP in providing a reliable end-to-end data transfer function, and described how TCP incorporates numerous control functions that are intended to make efficient use of the underlying IP network through a host-based congestion control function. Congestion control is an important component of TCP implementations, and today TCP congestion control plays an important role in the overall stability of the Internet.

Today’s Internet spans a very broad base of uses, and ensuring that TCP provides a highly robust, efficient, and reliable service platform for such a diversity of use is a continuing task. The Web has introduced a component of short duration reliable transfers into the public Internet traffic profile. These short sessions are often referred to as “TCP mice” because of the short duration and large number of such TCP sessions. Complementing these short sessions is the increasing size of large transfers as *File Transfer Protocol* (FTP) data sets become larger in response to increasing capacity within the public Internet network^[4]. In addition, there is an increasing diversity of media used within the Internet, both in terms of higher-speed systems and in the use of wireless systems for Internet access. In this article we will extend our examination of TCP by looking at how TCP is being used and adapted to match this changing environment.

A Review of TCP Performance

Within any packet-switched network, when demand exceeds available capacity, the packet switch will use a queue to hold the excess packets. When this queue fills, the packet switch must drop packets. Any reliable data protocol that operates across such a network must recognize this possibility and take corrective action. TCP is no exception to this constraint. TCP uses data sequence numbering to identify packets, and explicit acknowledgements (ACKs) to allow the sender and receiver to be aware of reliable packet transfer. This form of reliable protocol design is termed “end-to-end” control, because interior switches do not attempt to correct packet drops. Instead, this function is performed through the TCP protocol exchange between sender and receiver. TCP uses cumulative ACKs rather than per-packet ACKs, where an ACK referencing a particular point within the data stream implicitly acknowledges all data with a sequence value less than the ACKed sequence.

TCP also uses ACKs to clock the data flow. ACKs arriving back at the sender arrive at intervals approximately equal to the intervals at which the data packets arrived at the sender. If TCP uses these ACKs to trigger sending further data packets into the network, then the packets will be entered into the network at the same rate as they are arriving at their destination. This mode of operation is termed “ACK clocking.”

TCP recovers from packet loss using two mechanisms. The most basic operation is the use of packet timeouts by the sender. If an ACK for a packet fails to arrive within the timeout value, the sender will retransmit the oldest unacknowledged packet. In such a case, TCP assumes that the loss was caused by a network congestion condition, and the sender will enter “Slow Start” mode. This condition causes significant delays within the data transfer, because the sender will be idle during the timeout interval and upon restarting will recommence with a single packet exchange, gradually recovering the data rate that was active prior to the packet loss. Many networks exhibit transient congestion conditions, where a data stream may experience loss of a single packet within a packet train. To address this, TCP introduced the mechanism of “fast recovery.” This mechanism is triggered by a sequence of three duplicate ACKs received by the data sender. These duplicate ACKs are generated by the packets that trail the lost packet, where the sender ACKs each of these packets with the ACK sequence value of the lost packet. In this mode the sender immediately retransmits the lost packet and then halves its sending rate, continuing to send additional data as permitted by the current TCP sending window. In this mode of operation, “congestion-avoidance” TCP increases its sending window at a linear rate of one segment per *Round-Trip Time* (RTT). This mode of operation is referred to as *Additive Increase, Multiplicative Decrease* (AIMD), where the protocol reacts sharply to signs of network congestion, and gradually increases its sending rate in order to equilibrate with concurrent TCP sessions.

TCP Design Assumptions

It is difficult to design any transport protocol without making some number of assumptions about the environment in which the protocol is to be used, and TCP certainly has some inherent assumptions hidden within its design. The most important set of assumptions that lie behind the design of TCP are as follows:

- *A network of wires, not wireless:* As we continually learn, wireless is different. Wireless systems typically have higher *bit error rates* (BERs) than wire-based carriage systems. Mobile wireless systems also include factors of signal fade, base-station handover, and variable levels of load. TCP was designed with wire-based carriage in mind, and the design of the protocol makes numerous assumptions that are typical of such of an environment. TCP makes the assumption that packet loss is the result of network congestion, rather than bit-level corruption. TCP also assumes some level of stability in the RTT, because TCP uses a method of damping down the changes in the RTT estimate.
- *A best-path route-selection protocol:* TCP assumes that there is a single best metric path to any destination because TCP assumes that packet reordering occurs on a relatively minor scale, if at all. This implies that all packets in a connection must follow the same path within the network or, if there is any form of load balancing, the order of packets within each flow is preserved by some network-level mechanism.

- *A network with fixed bandwidth circuits, not varying bandwidth:* TCP assumes that available bandwidth is constant, and will not vary over short time intervals. TCP uses an end-to-end control loop to control the sending rate, and it takes many RTT intervals to adjust to varying network conditions. Rapidly changing bandwidth forces TCP to make very conservative assumptions about available network capacity.
- *A switched network with first-in, first-out (FIFO) buffers:* TCP also makes some assumptions about the architecture of the switching elements within the network. In particular, TCP assumes that the switching elements use simple FIFO queues to resolve contention within the switches. TCP makes some assumption about the size of the buffer as well as its queuing behavior, and TCP works most efficiently when the buffer associated with a network interface is of the same order of size as the delay bandwidth product of the associated link.
- *The duration of TCP sessions:* TCP also makes some assumptions about the nature of the application. In particular, it assumes that the TCP session will last for some number of round-trip times, so that the overhead of the initial protocol handshake is not detrimental to the efficiency of the application. TCP also takes numerous RTT intervals to establish the characteristics of the connection in terms of the true RTT interval of the connection as well as the available capacity. The introduction of short-duration sessions, such as found in transaction applications and short Web transfers, is a new factor that impacts the efficiency of TCP.
- *Large payloads and adequate bandwidth:* TCP assumes that the overhead of a minimum of 40 bytes of protocol per TCP packet (20 bytes of IP header and 20 bytes of TCP header) is an acceptable overhead when compared to the available bandwidth and the average payload size. When applied to low-bandwidth links, this is no longer the case, and the protocol overheads may make the resultant communications system too inefficient to be useful.
- *Interaction with other TCP sessions:* TCP assumes that other TCP sessions will also be active within the network, and that each TCP session should operate cooperatively to share available bandwidth in order to maximize network efficiency. TCP may not interact well with other forms of flow-control protocols, and this could result in unpredictable outcomes in terms of sharing of the network resource between the active flows as well as poor overall network efficiency.

If these assumptions are challenged, the associated cost is that of TCP efficiency. If the objective is to extend TCP to environments where these assumptions are no longer valid, while preserving the integrity of the TCP transfer and maintaining a high level of efficiency, then the TCP operation itself may have to be altered.

There are two basic ways of altering TCP operation: by altering the actions of the end host by making changes to the TCP protocol, or by altering the characteristics of the network, making them more “friendly” to TCP. We will look at the potential for both responses in examining various scenarios for adapting TCP to suit these changing environments.

Some caution should be noted about making changes to the TCP protocol. The major constraint is that any changes that are contemplated to TCP should be backward compatible with existing TCP behavior. This constraint requires a modified TCP protocol to attempt to negotiate the use of a specific protocol extension, and the knowledge that a basic common mode of protocol operation may be required if the negotiation fails. The second constraint is that TCP does assume that it is interacting with other TCP sessions within the network, and the outcome of fair sharing of the network between concurrent sessions depends on some commonality of the protocol used by these sessions. Major changes to the protocol behavior can lead to unpredictable outcomes in terms of sharing of the network resource between “unmodified” and “modified” TCP sessions, and unpredictable outcomes in terms of efficiency of the use of the network. For this reason there is some understandable reluctance to undertake modifications of TCP that radically alter TCP startup behavior or behavior in the face of network congestion.

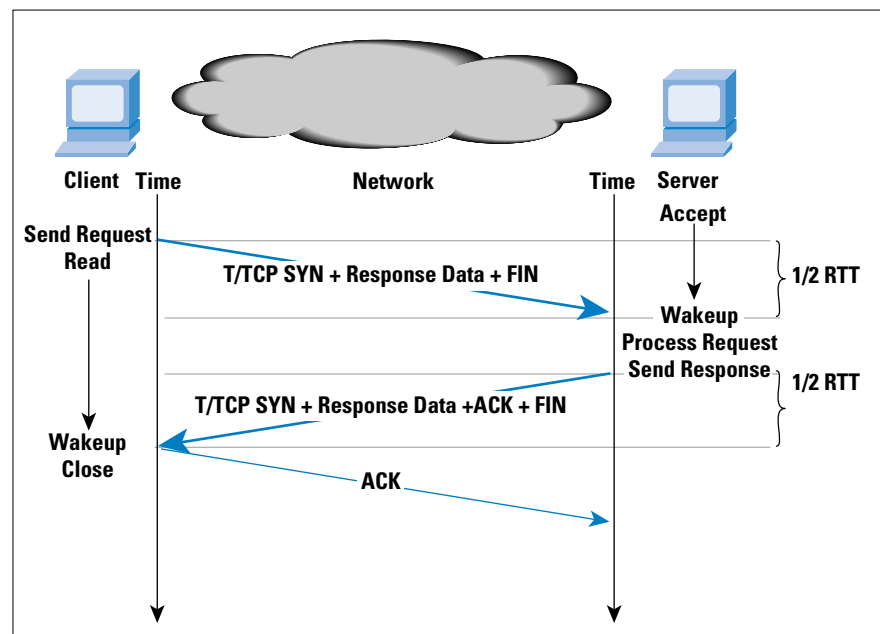
Short-Duration Sessions—TCP for Transactions

For network applications that generate small transactions, the application designer is faced with a dilemma. The application may be able to use the *User Datagram Protocol* (UDP), in which case the sender must send the query and await the response. This operation is highly efficient, because the total elapsed time for the client is a single RTT. However, this speed is gained at the cost of reliability. A missing response is ambiguous, in that it is impossible for the initiator to tell whether the query was lost or the response was lost. If multiple queries are generated, it is not necessarily true that they will arrive at the remote server in the same order as they were generated. Alternatively, the application can use TCP, which will ensure reliability of the transaction. However, TCP uses a three-way handshake to complete the opening of the connection, and uses acknowledged FIN signals for each side to close its end of the connection after it has completed sending data. Under the control of TCP, the sender will retransmit the query until it receives an acknowledgment that the query has arrived at the remote server. Similarly, the remote host will retransmit the response until the server receives an indication that the response has been successfully delivered. The cost of this reliability is application efficiency, because the minimum time to conduct the TCP transaction for the client is two RTT intervals.

TCP for Transactions (commonly referred to as T/TCP^[5]) attempts to improve the performance of small transactions while preserving the reliability of TCP. T/TCP places the query data and the closing FIN in the initial SYN packet. This can be interpreted as attempting to open a session, pass data, and close the sender's side of the session within a single packet. If the server accepts this format, the server responds with a single packet, which contains its SYN response, an ACK of the query data, the server's data in response, and the closing FIN. All that is required to complete the transaction is for the query system to ACK the server's data and FIN (Figure 1). If the server does not accept this format, the client can back off to a conventional TCP handshake followed by a data exchange.

For the client, the time to undertake this T/TCP transaction is one RTT interval, a period equal to the UDP-supported transaction, while still allowing for the two systems to use TCP to negotiate a reliable exchange of data as a backup.

Figure 1: T/TCP Operation



T/TCP requires changes to the protocol stack of both the sender and the receiver in order to operate correctly. The design of the protocol explicitly allows the session initiator to back off to use TCP if the receiver cannot correctly respond to the initial T/TCP packet.

T/TCP is not in common use in the Internet today, because while it improves the efficiency of simple transactions, the limited handshake makes it more vulnerable from a security perspective, and concerns over this vulnerability have been a prohibitive factor in its adoption. This is illustrative of the nature of the trade-offs that occur within protocol design, where optimizing one characteristic of a protocol may be at the expense of other aspects of the protocol.

Long Delay—TCP for Satellite Paths

Satellite-based services pose a set of unique issues to the network designer. Most notably, these issues include delay, bit errors, and bandwidth.

When using a satellite path, there is an inherent delay in the delivery of a packet due to signal propagation times related to the altitude of communications satellites. Geo-stationary orbit spacecraft are located at an altitude of some 36,000 km, and the propagation time for a signal to pass from an earth station directly below the satellite to the satellite and back is 239.6 ms. If the earth station is located at the edge of the satellite view area, this propagation time extends to 279.0 ms. In terms of a round trip that uses the satellite path in both directions, the RTT of a satellite hop is between 480 and 560 ms.

The strength of a radio signal falls in proportion to the square of the distance traveled. For a satellite link, the signal propagation distance is large, so the signal becomes weak before reaching its destination, resulting in a poor signal-to-noise ratio. Typical BERs for a satellite link today are on the order of 1 error per 10 million bits (1×10^{-7}). *Forward error correction* (FEC) coding can be added to satellite services to reduce this error rate, at the cost of some reduction in available bandwidth and an increase in latency due to the coding delay.

There is also a limited amount of bandwidth available to satellite systems. Typical carrier frequencies for commercial satellite services are 6/4 GHz (C-band) and 14/12 GHz (Ku band). Satellite transponder bandwidth is typically 36 MHz^[6].

When used in a data carriage role for IP traffic, satellite channels pose several challenges for TCP.

The delay-bandwidth product of a transmission path defines the amount of data TCP should have within the transmission path at any one time, in order to fully utilize the available channel capacity. The delay used in this equation is the RTT and the bandwidth is the capacity of the bottleneck link in the network path. Because the delay in satellite environments is large, a TCP flow may need to keep a large amount of data within the transmission path. For example, a typical path that includes a satellite hop may have a RTT of some 700 ms. If the bottleneck bandwidth is 2 Mbps, then a sender will need to buffer 180 kB of data to fully utilize the available bandwidth with a single traffic flow. For this to be effective, the sender and receiver will need to agree on the use of TCP Window Scaling to extend the available window size beyond the protocol default limit of 64 kB. A sender using an 8 kB buffer would be able to achieve a maximum transfer rate of 91 kbps, irrespective of the available bandwidth on the satellite path.

Even with advanced FEC techniques, satellite channels exhibit a higher BER than typical terrestrial networks. TCP interprets packet drop as a signal of network congestion, and reduces its window size in an attempt to alleviate the situation. In the absence of certain knowledge about whether a packet was dropped because of congestion or corruption, TCP must assume the drop was caused by congestion in order to avoid congestion collapse^[7, 8]. Therefore, packets dropped because of corruption cause TCP to reduce the size of its sending window, even though these packet drops do not signal congestion in the network. To mitigate this, some care must be taken with the satellite hop *Maximum Transmission Unit* (MTU) size, to reduce the probability of packet corruption. This is an area of compromise, in that the consequence is the potential for a high level of IP packet fragmentation on the satellite feeder router. In addition, the sender needs to use the TCP fast retransmit and fast recovery algorithms^[9] in order to recover from the packet loss in a rapid, but stable fashion. In addition, the sender needs to use larger sending windows to operate the path more efficiently, with a consequent risk of multiple packet drops per RTT window. For this reason the use of *Selective Acknowledgements* (SACKs) is necessary in order to recover from multiple packet drops in a single RTT interval.

The long delay causes TCP to react slowly to the prevailing conditions within the network. The slow start of TCP commences with a single packet exchange, and it takes some number of RTT intervals for the sender's rate to reach the same order of size as the delay bandwidth product of the long delay path. For short-duration TCP transactions, such as much of the current Web traffic, this is a potential source of inefficiency. For example, if a transaction requires the transfer of ten packets, the slow-start algorithm will send a single packet in the first RTT interval, two in the second interval, four in the third, and the remaining three packets in the fourth RTT interval. Irrespective of the available bandwidth of the path, the transaction will take a minimum of four RTT intervals. This theoretical model is further exacerbated by delayed ACKs [RFC 1122], where a receiver will not immediately ACK a packet, but will await the expiration of the 500ms ACK timer, or a second full-sized packet. During slow start, where a sender sends an initial packet, and then awaits an ACK, the receiver will delay the ACK until the expiration of the delayed ACK timer, adding up to 500ms additional delay in the first data exchange. The second part of the delayed ACK algorithm is that it will only ACK every second full-sized data packet, slowing down the window inflation rate of slow start. Also, if congestion occurs on the forward data path, the TCP sender will not be aware of the condition until it receives duplicate ACKs from the receiver. A congestion condition may take many RTT intervals to clear, and in the case of a satellite path, transient congestions may take tens of seconds to be resolved.

The TCP mechanisms that assist in mitigating some of the more serious effects of satellite systems include *Path MTU Discovery*^[10], *Fast Retransmit* and *Fast Recovery*, window scaling options, in order to extend the sender's buffer beyond 65,535 bytes^[11], and the companion mechanisms of *Protection Against Wrapped Sequence Space* (PAWS) and *Round-Trip Time Measurements* (RTTM) and SACKs^[12]. A summary of TCP options is shown in Figure 2.

Figure 2: TCP Options
for Satellite Paths
(after RFC 2488)

Mechanism	Use	Location
Path-MTU Discovery	Recommended	Sender
FEC	Recommended	Link
TCP		
Slow Start	Required	Sender
Congestion Avoidance	Required	Sender
Fast Retransmit	Recommended	Sender
Fast Recovery	Recommended	Sender
Window Scaling	Recommended	Sender and Receiver
PAWS	Recommended	Sender and Receiver
RTTM	Recommended	Sender and Receiver
SACK	Recommended	Sender and Receiver

Further refinements to the TCP stack have been considered in relation to satellite performance^[13].

The options considered include the use of T/TCP as a means of reducing the overhead of the initial TCP three-way handshake. This is effective for short transactions where the data to be transferred can be held in a single packet, or in a small number of packets.

The use of delayed acknowledgements also is an issue for long-delay network paths, particularly if the sender is using slow start with an initial window of a single segment. In this case, the receiver will not immediately acknowledge the initial packet, but will wait up to one-half second for the delayed ACK timer to trigger. Altering the initial window size to two segments allows the receiver to trigger an ACK on reception of the second packet, bypassing the delayed ACK timer. However, even this change to TCP does not completely address the performance issue relating to delayed ACKs on long delay paths for TCP slow start. The delayed ACK algorithm triggers an ACK on every second full-sized packet. Because the sender's congestion window is opened on receipt of ACKs, this causes the slow-start window to open more slowly than if the receiver generated an ACK every packet. One variant of TCP congestion control allows the TCP sender to count the number of bytes acknowledged in an ACK message to control the expansion of the congestion window, making the algorithm less sensitive to delayed ACKs^[9]. Although this approach has some merit for long delay paths, this is a case where the correction is potentially as bad as the original problem. The byte counting mode of congestion control allows a sender to sharply increase its sending rate, causing potential instabilities within the network and impacting concurrent TCP sessions.

One approach to address this is to place a limit on the size of the window expansion, where each increment of the congestion window is limited to the minimum of one or two segment sizes and the size of the data spanned by the ACK. If the limit is set to a single segment size, the window expansion will be in general slightly more conservative to the current TCP ACK-based expansion mechanism. If this upper limit is set to two segments, the congestion window expansion will account for the delayed ACKs, expand at a rate equal to one segment for every successfully transmitted segment during slow start, and expand the window by one segment size each RTT during congestion avoidance. Because a TCP receiver will ACK a large span of data following recovery, this byte counting is bounded to a single segment per ACK in the slow-start phase following a transmission timeout. Another approach that has been explored is for the receiver to disable delayed ACKs until the sender has completed the slow-start phase. Although such an approach shows promising results under simulated conditions, the practical difficulty is that it is difficult for the receiver to remotely determine the current TCP sending state, and the receiver cannot reliably tell if the sender is in slow start, congestion avoidance, or in some form of recovery mode. Explicit signaling of the sender's state as a TCP flag is an option, but the one-half RTT delay in the signaling from the sender to the receiver may prove to be an issue here. This area of congestion control for TCP remains a topic of study.

All of these approaches can mitigate only the worst of the effects of the long delay paths. TCP, as an adaptive reliable protocol that uses end-to-end flow control, can undertake only incremental adjustments in its flow rates in intervals of round-trip times. When the round-trip times extend, then TCP is slower to speed up from an initial start, slower to recover from packet loss, and slower to react to network congestion.

Tuning TCP—ACK Manipulation

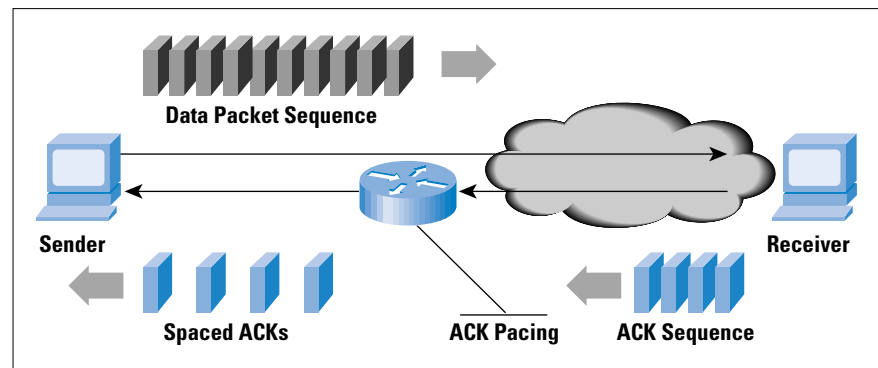
The previous article of TCP Performance discussed numerous network responses to congestion using *Random Early Detection* (RED) for active queue control and *Explicit Congestion Notification* (ECN) as an alternative to RED packet drop. It is feasible for a network control point to impose a finer level of control on a TCP flow by using an approach of direct manipulation of the TCP packets.

The approaches described above to mitigate some of the side effects of satellite paths all share in the side effect of having some latency associated with the congestion response. The sender must await the reception of trailing packets by the receiver, and then await the reception of the matching ACK packets from the data receiver back to the sender to learn of the fate of the original data packet. This may take up to one RTT interval to complete. An alternative approach to congestion management responses is to manipulate the ACK packets to modify the sender's behavior.

The prerequisite to perform this manipulation is that the traffic path be symmetric, so that the congestion point can identify ACK packets traveling in the opposite direction. If this is the case, a couple of control alternatives can mitigate the onset of congestion:

- *ACK Pacing*: Each burst of data packets will generate a corresponding burst of ACK packets. The spacing of these ACK packets determines the burst rate of the next sending packet sequence. For long-delay systems, the size of such bursts becomes a limiting factor. TCP slow start generates packet bursts at twice the bottleneck data rate, so that the bottleneck feeder router may have to absorb one-half of every packet burst within its internal queues. If these queues are not dimensioned to the delay bandwidth product of the next hop, these queues become the limiting factor, rather than the path bandwidth itself. If you can slow down the TCP burst rate, the pressure on the feeder queue is alleviated. One approach to slow down the burst rate is to impose a delay on successive ACKs at a network control point (Figure 3). This measure will reduce the burst rate, but not impact the overall TCP throughput. ACK pacing is most effective on long delay paths, and it is intended to spread out the burst load, reducing the pressure on the bottleneck queue and increasing the actual data throughput.

Figure 3: ACK Pacing



- *Window Manipulation*: Each ACK packet carries a receiver window size. This advertised window determines the maximum burst size available to the sender. Manipulating this window size downward allows a control point to control the maximal TCP sending rate. This manipulation can be done as part of a traffic-shaping control point, enforcing bandwidth limitations on a flow or set of flows.

Both of these mechanisms make some sweeping assumptions about the network control point that must be carefully understood. The major assumption is that these mechanisms assume symmetry of data flows at the network control point, where the data and the associated ACKs flow through this control point (but in opposite directions, of course). Both mechanisms also assume that the control point can cache per-flow state information, so that the current flow RTT and the current transfer rate and receiver window size are available to the service controller.

ACK pacing also implicitly assumes that a single ACK timing response is active at any time along a network path. A sequence of ACK delay actions may cause the sender's timers to trigger, and the sender to close down the transfer and reenter slow-start mode. These environmental conditions are more common at the edge of the network, and such mechanisms are often part of a traffic control system for Web-hosting platforms or similar network service delivery platforms. As a network control tool, ACK manipulation makes too many assumptions, and the per-flow congestion state information represents a significant overhead for large network systems. In general, such manipulations are more appropriate as an edge traffic filter, rather than as an effective congestion management response. For this reason, the more indirect approach of selective data packet discard is more effective as a congestion management measure.

Assisting Short-Duration TCP Sessions—Limited Transmit

One of the challenges to the original set of TCP assumptions is that of short-duration TCP sessions. The Web has introduced a large number of short-duration sessions, and the issue with these sessions is that they use small initial windows. If congestion loss occurs within this early period of TCP slow start, there are not enough packets in the network to generate the three duplicate ACKs required to initiate fast retransmit and fast recovery. Instead the TCP sender must await the expiry of the *retransmission timeout* (RTO), a timer that uses a minimum value of one second. For short-duration TCP sessions that may last six or seven RTT intervals of a small number of milliseconds, the incremental penalty of single packet loss is then extremely severe. A study of this problem indicates that approximately 56 percent of retransmissions are sent following an RTO timeout^[25].

One potential mitigation to this is a mechanism termed “Limited Transmit.” With this mechanism, a duplicate ACK may trigger an immediate transmission of a segment of new data. Two conditions are applied to this; the receiver's advertised window allows the transmission of this segment, and the amount of outstanding data would remain less than the congestion window plus the duplicate ACK threshold used to trigger Fast Retransmit. This second condition implies that the sender can send only two segments beyond the congestion window, and will do so only in response to the receiver lifting a segment off the network. The basic principle of this strategy is to continue the signaling between the sender and receiver in the face of packet loss, increasing the probability that the sender will recover from packet loss using duplicate ACKs and fast recovery, and reducing the probability of the one-second (or longer) RTO timeout as being the recovery trigger. The limited transmit also reduces the potential for the recovery actions to burst into the network at a level that may cause further packet loss.

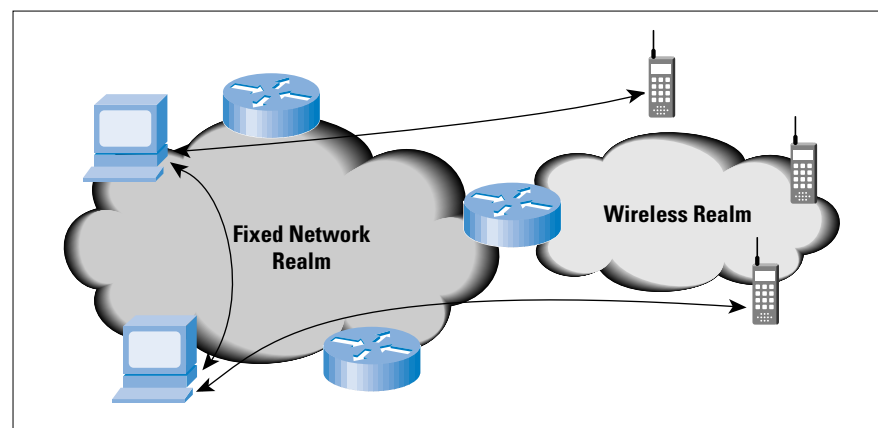
Low Bandwidth and High Error Rates—TCP for Wireless Systems

One of the more challenging environments for the Internet Protocol, and TCP in particular, is that of mobile wireless.

One approach to supporting the wireless environment is that of the so-called “walled garden.” Here the protocols in use within the wireless environment are specifically adapted to the wireless world. The transport protocols can account for the low bandwidth, the longer latency, the BERs, and the variability within all three of these metrics. In this model, Internet applications interact with an application gateway to reach the wireless world, and the application gateway uses a wireless transport protocol and potentially a modified version of the application data to interact with the mobile wireless device. The most common approach is extension of the World Wide Web client into the mobile wireless device, using some form of proxy server at the boundary of the wireless network and the Internet. This is the approach adopted by the *Wireless Access Protocol Forum* (WAP)^[14].

An alternative approach lies in extending not only the World Wide Web to a mobile handset, but also allowing mobile devices to access a complete range of Internet-based services as the functional objective. In this approach, the intent is to allow the mobile wireless device to function as any other Internet-connected device, and there is a consequent requirement for some form of end-to-end direct IP continuity, and an associated requirement for end-to-end TCP functionality, where the TCP path straddles both wired and wireless segments. Ensuring the efficient operation of TCP in this environment is an integral part of the development of such an environment. Given that TCP must now work within a broader environment, it is no longer a case of adjusting TCP to match the requirements of the wireless environment, but one of attempting to provide seamless interworking between the wired and wireless worlds (Figure 4).

Figure 4: Linking the Wired and Wireless Worlds



The wireless environment challenges many of the basic assumptions of TCP noted above. Wireless has significant levels of bit error rates, often with bursting of very high error rates. Wireless links that use forward error correcting codes have higher latency. If the link level protocol includes automatic retransmission of corrupted data, this latency will have high variability. Wireless links may also use adaptive coding techniques that adjust to the prevailing signal to noise ratio of the link, in which case the link will have varying bandwidth. If the wireless device is a hand-held mobile device, it may also be memory constrained. And finally, such an environment is typically used to support short duration TCP sessions.

The major factor for mobile wireless is the BER, where frame loss of up to 1 percent is not uncommon, and errors occur in bursts, rather than as evenly spaced bit errors in the packet stream. In the case of TCP, such error conditions force the TCP sender to initially attempt fast retransmit of the missing segments, and when this does not correct the condition, the sender will have an ACK timeout occur, causing the sender to collapse its sending window and recommence from the point of packet loss in slow-start mode. The heart of this problem is that assumption on the part of TCP that packet loss is a symptom of network congestion rather than packet corruption. It is possible to use a model of TCP AIMD performance to determine the effects of this loss rate on TCP performance. If, for example the link has a 1-percent average packet loss rate, a *Maximum Segment Size* (MSS) size of 1000 bytes, and a 120ms RTT, then the AIMD models predict a best-case performance of 666Kbps throughput, and a more realistic target of 402Kbps throughput^[15]. (See the appendix on page 24 for details of these models.) TCP is very sensitive to packet loss levels, and sustainable performance rapidly drops when packet drop levels exceed 1 percent.

Link-level solutions to the high BER are available to designers, and FEC codes and *automatic retransmission systems* (ARQ) can be used on the wireless link. FEC introduces a relatively constant coding delay and a bandwidth overhead into the path, but cannot correct all forms of bit error corruption. ARQ uses a “stop and resend” control mechanism similar to TCP itself. The consequent behavior is one of individual packets experiencing extended latency as the ARQ mechanisms retransmit link-level fragments to correct the data corruption, because the packet flow may halt for an entire link RTT interval for the link-level error to be signaled and the corrupted level 2 data to be retransmitted. The issue here is that TCP may integrate these extended latencies into its RTT estimate, making TCP assume a far higher latency on the path than is the case, or, more likely, it may trigger a retransmission at the same time as the level 2 ARQ is already retransmitting the same data. An alternative Layer 2 approach to bit-level corruption is to deliver those level 2 frames that were successfully transmitted, while resending any frames that were corrupted in transmission.

The problem for TCP here is that the level 2 drivers are adding packet reordering to the extended latency, and from TCP perspective the delivery of the out-of-order packets will generate duplicate ACKs that may trigger a simultaneous TCP fast retransmit.

Perversely, some approaches have advocated TCP delaying its duplicate ACK response in such situations^[13]. To quote from RFC 2488, “The interaction between link-level retransmission and transport-level retransmission is not well understood.”^[6]

If ARQ is not the best possible answer to addressing packet loss in mobile wireless systems, then what can be done at the TCP level to address this? TCP can take numerous basic steps to alleviate the worst aspects of packet corruption on TCP performance. These include the use of Fast Retransmit and Fast Recovery to allow a single packet loss to be repaired moderately quickly. This mechanism triggers only after three duplicate ACKs, so the associated action is to ensure that the TCP sender and receiver can advertise buffers of greater than four times the MSS. SACKs allow a sender to repair multiple segment losses per window within a single RTT, and where large windows are operated over long delay paths, SACK is undoubtedly useful.

However, useful as these mechanisms may be, they are probably inadequate to allow TCP to function efficiently over all forms of wireless systems. Particularly in the case of mobile wireless systems, packet corruption is sufficiently common that, for TCP to work efficiently, some form of explicit addressing of network packet corruption appears to be necessary.

One approach is to decouple TCP congestion control mechanisms from data recovery actions. The intent is to allow new data to be sent during recovery to sustain TCP ACK clocking. This approach is termed *Forward Acknowledgements with Rate Halving* (FACK)^[13], where one packet is sent for every two ACKs received while TCP is recovering from lost packets. This algorithm effectively reduces the sending rate by one-half within one RTT interval, but does not freeze the sender to wait the draining on one-half of the congestion window’s amount of data from the network before proceeding to sending further data, nor does it permit the sender to burst retransmissions into the network. This is particularly effective for long-delay networks, where the fast recovery algorithm causes the sender to cease sending for up to one RTT interval, thereby losing the accuracy of the implicit ACK clock for the session. FACK allows the sender to continue to send packets into the network during this period, in an effort to allow the sender to maintain an accurate view of the ACK clock. FACK also provides an ability to set the number of SACK blocks that specify a missing segment before re-sending the segment, allowing the sender greater levels of control over sensitivity to packet reordering. The changes to TCP to support FACK are a change in the sender’s TCP to use the FACK algorithm for recovery, and, for optimal performance, use of SACK options by the receiver.

In looking for alternative responses to packet corruption, it is noted that TCP segments that are corrupted are often detected at the link level, and are discarded by the link-level drivers. This discard cannot be used to generate an error message to the packet sender, given that the IP header of the packet may itself be corrupted, nor can the discard signal be reliably passed to the receiver, for the same reason. However, despite this unreliability of information, this signaling from the link level to the transport level is precisely the objective here, because, at the TCP protocol level, the sender needs to be aware that the packet loss was not due to network congestion, and that there is no need to take corrective action in terms of TCP congestion behavior.

One approach to provide this signaling from the data link level to the transport level calls for the link-level device to forward a “corruption experienced” *Internet Control Message Protocol* (ICMP) packet when discarding a corrupted packet^[13]. This approach has the ICMP packet being sent in the forward direction to the receiver, who then has the task of converting this message and the associated lost packet information into a signal to the sender that the duplicate ACKs are the result of corruption, not network congestion. This signal from the receiver to the sender can be embedded in a TCP header option. The sending TCP session will maintain a corruption experienced state for two RTT intervals, retransmitting the lost packets without halving the congestion window size.

As we have noticed, corruption may have occurred in the packet header, and the sender’s address may not be reliable. This approach addresses this by having the router keep a cache of recent packet destinations, and when the IP header information is unreliable because of a failed IP header checksum, the router will forward the ICMP message to all destinations in the cache. The potential weakness in this approach is that if network congestion occurs at the same time as packet corruption, the sender will not react to the congestion, and will continue to send into the congestion for a further two RTT intervals. This approach is not without some deployment concerns. It calls for modification to the wireless routers and to the receiver’s link-level drivers to generate the ICMP corruption experienced messages, modification to the receiver’s IP stack in order to take signals from the IP ICMP processor and from the link-level driver and convert them to TCP corruption loss signals within the TCP header of the duplicate ACKs, and modifications to the TCP processor at the sender to undertake corruption-experienced packet loss recovery. Even with these caveats in mind, this approach of explicit corruption signaling is a very promising approach to addressing performance issues with TCP over wireless.

Of course high levels of bit errors is not the only problem facing TCP over wireless systems. Mobile wireless systems are typically small handsets or personal digital assistants, and the application transactions are often modified to reduce the amount of data transferred, given that a limited amount of data can be displayed on the device.

In this case, the ratio between payload and IP and TCP headers starts to become an issue, and some consideration of header compression is necessary. Header compression techniques typically take the form of stripping out those fields of the header that do not vary on a packet-by-packet basis, or that vary by amounts that can be derived from other parts of the header, and then transmitting the delta values of those fields that are varying^[16, 17].

Although such header compression schemes can be highly efficient in operation, the limitation of such schemes is that the receiver needs to have successfully received and decompressed the previous packet before the receiver can decompress the next packet in the TCP stream. In the face of high levels of bit error corruption, such systems do introduce additional latencies into the data transfer, and multiple packet drops are difficult to detect and signal via SACK in this case.

A more subtle aspect of mobile wireless is that of temporary link outages. For example, a mobile user may enter an area of no signal coverage for a period of time, and attempt to resume the data stream when signal is obtained again. In the same way that there is no accepted way of a link-level driver informing TCP of packet loss due to corruption, there is no way a link-level driver can inform TCP of a link-level outage. In the face of such link-level outages, TCP will assume network-level congestion, and in the absence of duplicate ACKs, TCP retransmission timers will trigger. TCP will then attempt to restart the session in slow-start mode, commencing with the first dropped packet. Each attempt to send the packet will result in TCP extending its retransmission timer using an exponential backoff on each attempt, so that successive probes are less and less frequent. Because the link level cannot inform the sender on the resumption of the link, TCP may wait some considerable time before responding to link restoration. The intention is for the link level to be able to inform the TCP for resumption of the connection following a link outage. One approach is for the link level to retain a packet from each TCP stream that attempted to use the link. When the link becomes operational again, the link-level driver immediately transmits these packets on the link. The result is that the receiver will then generate a response that will then trigger the sender into transmission within a RTT interval. Only a single packet per active TCP stream is necessary to trigger this response, so that the link level does not need to hold an extensive buffer of undeliverable packets during a link outage. Of course if the routing level repaired the link outage in the meantime, the delivery of an out-of-order TCP packet would normally be discarded by the sender.

The bottom line here is the question: Is TCP suitable for the mobile wireless environment? The answer appears to be that TCP can be made to work as efficiently as any other transport protocol for the mobile wireless environment.

However, this does imply that some changes in the operation of TCP need to be undertaken, specifically relating to the signaling of link-level states into the TCP session and use of advanced congestion control and corruption signaling within the TCP session. Although it is difficult to conceive of a change to every deployed TCP stack within the deployed Internet to achieve this added functionality, there does exist a middle ground between the “walled garden” approach and open IP. In this middle ground, the wireless systems would have access to “middle-ware,” such as Web proxies and mail agents. These proxies would use a set of TCP options when communicating with mobile wireless clients that would make the application operate as efficiently as possible, while still permitting the mobile device transparent access to the Internet for other transactions.

Unbundling TCP—Stream Control Transmission Protocol

There are occasions where the application finds the control functions of TCP too limiting. In the case of handling *Public Switched Telephone Network* (PSTN) signaling across an Internet network, the application requirements are somewhat different from those of TCP delivered service. PSTN signaling reliable delivery is important, but the individual transactions within the application are included within each packet, so the concept of preservation of strict order of delivery is unnecessary. Relaxation of this requirement of strict order of packet delivery allows the transport protocol to function more efficiently, because there is no head-of-line blocking at the receiver when awaiting retransmission of lost packets. TCP also assumes the transfer of a stream of data, so that applications that wish to add some form of record delineation to the data stream have to add their own structure to the data stream. In addition, the limited scope of TCP sockets complicates the support of a high-availability application that may use multihomed hosts, and TCP itself is vulnerable to many attacks, such as SYN attacks. The intention of the *Stream Control Transmission Protocol* (SCTP) is to address these application requirements^[16].

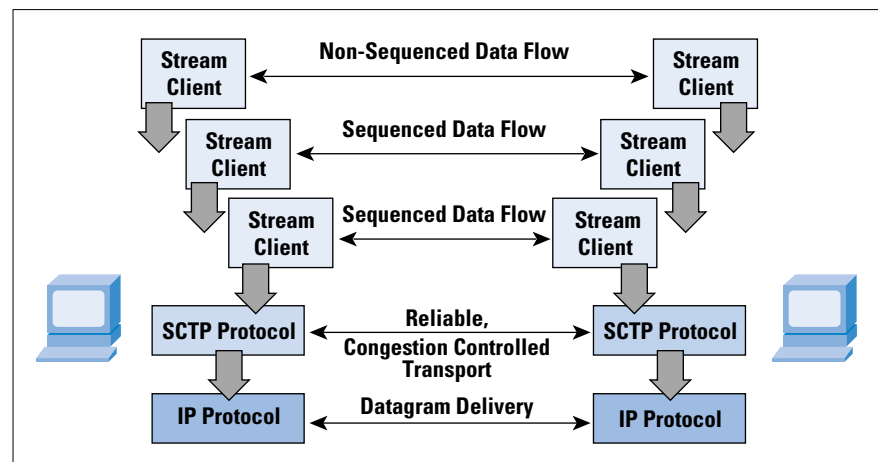
The first major difference between SCTP and TCP occurs during initialization, where the SCTP endpoints exchange a list of SCTP endpoint addresses (IP addresses and port numbers) that are to be associated with the SCTP session. Any pair of these source and destination addresses can be used within the SCTP session.

The startup of SCTP is also altered into a four-way handshake, where the initiator sends a tag value to the other end, which then responds with a copy of this tag and a tag of its own. At this stage the recipient does not allocate any resources for the connection, making the initialization sequence more robust in the face of TCP SYN-styled attacks. The initiator can then respond to this with an echo of the recipient’s tag (COOKIE-ECHO), and can also attach data to the response, allowing data to be transferred as early as possible in the handshake process.

After the recipient ACKs this message, the SCTP session is now established. The closing of an SCTP session is also different from TCP. In TCP, one side can close its sending function via a FIN TCP packet, and continue to receive packets, operating in a “half-open” state. In SCTP, a close from one side will cause the other end to drain its send queues and also shut down.

SCTP also functions in a form of transport-level multiplexing, where numerous logical streams can be supported across a single transport-level association. Although message order within an individual stream is preserved by SCTP, retransmission within one stream does not impact the operation of any other stream that is supported across the same SCTP transport association. Each stream has an explicit identification and a per-stream sequence identification to support this function. SCTP also provides for nonsequenced message delivery, where a message within a stream is marked for immediate delivery, irrespective of the relative order of the message within a stream (Figure 5).

Figure 5: The SCTP Transport Service Model



SCTP explicitly uncouples transport-level reliability and congestion control from per-stream sequenced delivery through the use of a separate transport-level interaction. The transport-level data and ACKs and the corresponding transport-level congestion window controls operate using a transport-level sequence space. This sequence space counts transport-level messages, not byte offsets within the message, so that no explicit window scaling option is necessary for SCTP. The congestion control functions reference those of TCP with fast retransmit and fast recovery, with an explicit specification of the SACK protocol and specification of the maintenance of the transmission timers and congestion control. SCTP also requires the use of MTU path discovery, so that larger transactions will use SCTP-level segmentation, avoiding the IP retransmission problem with lost fragments of a fragmented IP packet. SCTP does use a modified retransmission mechanism to that of TCP. Like TCP, SCTP associates a retransmission timer with each message, and if the timer expires the message is retransmitted and SCTP collapses the congestion window to a single message size. The SCTP receiver will generate SACK reports for a minimum of every second received packet.

If a message is within a SACK gap, then after three further such SACK messages, the sender will immediately send the missing messages, and half its congestion window, analogous to the fast retransmit and fast recovery of TCP.

The use of multiple endpoint addresses assumes that each of the endpoint addresses is associated with the same end host, but with a potentially different network path between the two endpoints. SCTP refreshes path availability to each of the endpoint addresses with a periodic keepalive, so that in the event of primary path failure, SCTP can continue by using one of the secondary endpoint addresses.

One could describe SCTP as being overly inclusive in terms of its architecture, and there is certainly a lot of capability in the protocol that is not contained within TCP. The essential feature of the protocol is to use a single transport congestion state between two systems to allow a variety of applications to attach as stream clients. In itself, this is analogous to TCP multiplexing. It also implicitly assumes that every stream is provided the same service level by the network, an assumption shared by almost all transport multiplexing systems. The essential alteration with SCTP is the use of many transport modes: reliable sequenced message streams, reliable sequenced streams with interrupt message capability, and reliable nonsequenced streams. It remains to be seen whether the utility provided by this protocol will become widely deployed within the Internet environment, or whether it will act as a catalyst for further evolution of transport service protocols.

Sharing TCP information—Endpoint Congestion Management

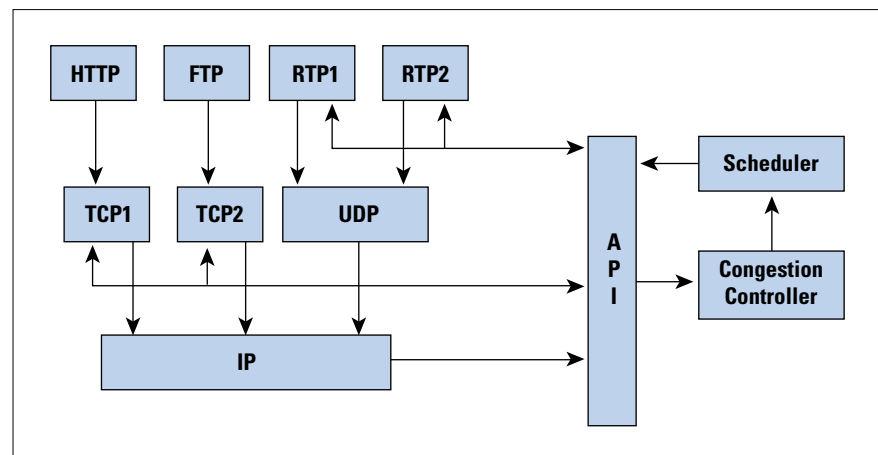
The notion of sharing a single TCP congestion state across multiple reliable streams is one that may also be applied to a mix of reliable and nonreliable data streams that operate concurrently between a pair of endpoints. It is this form of the multiplexing service model that is explored by the congestion manager model. The Congestion Manager is an end-system module that allows a collection of concurrent streams from the host to a single destination to share a common congestion control function, and permits various forms of reliable and nonreliable streams to use the network in a way that cooperates with concurrent congestion controlled flows^[19].

One of the major motivations for the congestion manager is the observation that the most critical part of network performance management is that of managing the interaction between congestion-controlled TCP streams and nonresponsive UDP data streams. In the extreme cases of this interaction, either traffic class can effectively deny service to the other by placing sufficient pressure on the network queuing resources that starve the other traffic class of any usable throughput. The observation made in the motivation for the congestion manager is that applications such as the Web typically open up a set of parallel connections to provide service, sending a mix of reliable flow-controlled data

along one connection and unreliable real-time streaming content along another. If the set of flows used a common congestion-control function at the sending host, the collection of flows would utilize the network resources in a manner analogous to a single TCP connection.

The manner of providing this common congestion control function is an advisory function to applications, as shown in Figure 6. One mechanism is that of a *callback*, where an application inserts a request to send a single message segment with the congestion manager. The Congestion Manager responds with invoking a callback to the requestor when the application may pass the data segment to the protocol driver. The other supported mechanism is that of *synchronous transmission*, where the Congestion Manager has a callback function that updates the application with a maximal available bit rate, the smoothed round-trip time estimate, and the smoothed linear deviation in the round-trip time estimate. In this mode the application can request further notification only when the network state changes by some threshold amount.

Figure 6:
The CM Model,
(after "The Congestion
Manager"^[19])



For the Congestion Manager to maintain a current picture of the congestion state of the path to the destination, each active stream needs to update the congestion manager as to the response from the remote host. It does this by informing the congestion manager of the number of bytes received, the number of bytes lost, and the RTT measurement, as measured at the application level. The application is also expected to provide an indication of the nature of the loss, as a timeout expiry, a transient network condition, or based on the reception of an ECN signal.

There has been little practical experience as yet with this model of shared congestion control within the Internet environment. There also remains a number of issues about how network performance information is passed back from the receiver to the sender in the absence of an active concurrent TCP session. The concurrent operation of a TCP session with a UDP streaming session to the same destination allows Congestion Manager to use the TCP congestion state to determine the sending capability of the streaming flow.

If the TCP session is idle, or if there is no TCP session, then the UDP streaming application will require some form of receiver feedback. The feedback will need to report on the span of data covered by the report, and the data loss rates and jitter levels, allowing the sender to assess the current quality and capacity of the network path.

This approach, and that of SCTP, are both illustrative of the approach of unbundling the elements of TCP and allowing applications to use combinations of these elements in ways that differ from the conventional monolithic transport-level protocol stack, with the intention of allowing the TCP congestion control behavior to be applied to a wider family of applications.

Better than TCP?

Recently, numerous “better-than-TCP” protocol stacks have appeared on the market, most commonly in conjunction with Web server systems, where the performance claim is that these protocol stacks can interoperate with standard TCP clients, but offer superior download performance to a standard TCP protocol implementation.

This level of performance is achieved by modifying the standard TCP flow control systems in a number of ways. The modified implementation may use a lower initial RTT estimate to provide a more aggressive startup rate, and a more finely grained RTT timer system to allow the sender to react more quickly to network state changes. Other modifications may include using a larger initial congestion window size or may use an even faster version of slow start, where the sending rate is tripled, or more, every round-trip time interval. The same technique of incremental modification can be applied to the congestion avoidance state, where the linear rate increase of one segment size per round-trip time interval can be increased to some multiple of the segment size, or use a time base other than the round-trip time for linear expansion of the congestion window. The backoff algorithm can also be altered such that the congestion window is reduced by less than half during congestion backoff. Resetting the TCP session to slow-start mode following the ACK timeout can also be avoided in such modified protocol implementations.

These techniques are all intended to force the sender to behave more aggressively in its transmission of packets into the network, thereby increasing the pressure on the network buffers. The network is not the only subject of this increased sending pressure; such modified protocol systems tend to impose a significant performance penalty on other concurrent TCP sessions that share the path with these modified protocol hosts. The aggressive behavior of the modified TCP systems in filling the network queues tends to cause the other concurrent standard TCP sessions to reduce their sending rate. This in turn opens additional space in the network for the modified TCP session to increase its transmission rate.

In an environment where the overall network resource-sharing algorithm is the outcome of dynamic equilibration between cooperative sending systems, such aggressive flow control modification can be considered to be extremely antisocial behavior at the network level. Paradoxically, such systems can also be less efficient than a standard TCP implementation. TCP server systems modified in this way tend to operate with higher levels of packet loss because their efforts to saturate the network with their own data packets make them less sensitive to the signals of network congestion.

Consequently, when delivering large volumes of traffic, or where there are moderately low levels of competitive pressure for network resources, the modified TCP stack may often perform less efficiently than a standard TCP implementation. Accordingly, these modified better-than-TCP implementations remain in the experimental domain. Within the production environment, their potential to impose undue performance penalties on concurrent TCP sessions and their potential to reduce overall network efficiency are reasonable indicators that such modified stacks should be used in private network environments, and with considerable care and discretion, if at all. Their utility in the public Internet is highly dubious.

TCP Evolution

The evolution of TCP is a careful balance between innovation and considered constraint. The evolution of TCP must avoid making radical changes that may stress the deployed network into congestion collapse, and also must avoid a congestion control “arms race” among competing protocols^[20]. The Internet architecture to date has been able to achieve new benchmarks of network efficiency, and translate this carriage efficiency into ground-breaking benchmark prices for IP-based carriage services. Much of the credit for this must go to the operation of TCP, which manages to work at that point of delicate balance between self-optimization and cooperative behavior.

Widespread deployment of transport protocols that take a more aggressive position on self-optimization will ultimately lead to situations of congestion collapse, while widespread deployment of more conservative transport protocols may well lead to lower jitter and lower packet retransmission rates, but at a cost of considerably lower network efficiency.

The challenges faced with the evolution of TCP is to maintain a coherent control architecture that has consistent behavior within the network, consistent interaction with instances of data flows that use the same control architecture, and yet be adequately flexible to adapt to differing network characteristics and differing application profiles. It is highly likely that we will see continued innovation within Internet transport protocols, but the bounds of such effort are already well recognized.

We can now state relatively clearly what levels of innovation are tolerable within an Internet network model that achieves its efficiency not through enforcement of rigidly enforced rules of sharing of the network resource, but through a process of trust between competing user demands, where each demand is attempting to equilibrate its requirements against a finite network capacity. This is the essence of the TCP protocol.

Appendix: TCP Performance Models

This appendix is an extract from “Advice for Internet Subnet Designers,” work in progress^[15].

The performance of the TCP AIMD Congestion Avoidance algorithm has been extensively analyzed. The current best formula for the performance of the specific algorithms used by Reno TCP is given by Padhye et. al.^[21], this formula is:

$$BW = \frac{MSS}{(RTT \times \sqrt{(1.33 \times \rho)}) + (RTO \times \rho \times [1 + 32 \times \rho^2] \times \min(1, 3 \times \sqrt{0.75 \times \rho}))}$$

MSS is the segment size being used by the connection.

RTT is the end-to-end round-trip time of the TCP connection.

RTO is the packet timeout (based on *RTT*).

ρ is the packet loss rate for the path (that is, 0.01 if there is 1-percent packet loss)

This is currently considered to be the best approximate formula for Reno TCP performance. A further simplification to this formula is generally made by assuming that *RTO* is approximately $5 \times RTT$.

TCP is constantly being improved. A simpler formula, which gives an upper bound on the performance of any AIMD algorithm that is likely to be implemented in TCP in the future, was derived by Ott, et.al.^[22, 23].

$$BW = 0.93 \times \frac{MSS}{RTT \sqrt{\rho}}$$

Assumptions of these formulae:

- Both of these formulae assume that the TCP Receiver Window is not limiting the performance of the connection in any way. Because the receiver window is entirely determined by end hosts, we assume that hosts will maximize the announced receiver window in order to maximize their network performance.
- Both of these formulae allow for bandwidth to become infinite if there is no loss. This is because an Internet path will drop packets at bottleneck queues if the load is too high. Thus, a completely lossless TCP/IP network can never occur (unless the network is being underutilized).
- The *RTT* used is the average *RTT* including queuing delays.

- The formulae are calculations for a single TCP connection. If a path carries many TCP connections, each will follow the formulae above independently.
- The formulae assume long-running TCP connections. For connections that are extremely short (<10 packets) and don't lose any packets, performance is driven by the TCP slow-start algorithm. For connections of medium length, where on average only a few segments are lost, single-connection performance will actually be slightly better than given by the formulae above.
- The difference between the simple and complex formulae above is that the complex formula includes the effects of TCP retransmission timeouts. For very low levels of packet loss (significantly less than 1 percent), timeouts are unlikely to occur, and the formulae lead to very similar results. At higher packet losses (1 percent and above), the complex formula gives a more accurate estimate of performance (which will always be significantly lower than the result from the simple formula).

Note that these formulae break down as ρ approaches 100 percent.

Addendum: An Update on Explicit Congestion Notification

The previous article on TCP performance noted that there was no explicit standardization of the IPv4 header field to carry the *Explicit Congestion Notification* (ECN) signals. As an update to the status of ECN, RFC 2481, the document that describes ECN, categorizes this proposal as an “Experimental” RFC document^[27]. The Internet Standards process^[28] describes this category as follows: “The ‘Experimental’ designation typically denotes a specification that is part of some research or development effort. Such a specification is published for the general information of the Internet technical community ...” ECN is the only experimental proposal to use these two bits of the IP header, and the use of the category “Experimental” reflects the current status of the proposal, in that the Internet Engineering Steering Group has, at the time of publication, yet to make a final decision to allocate these two bits of the IP header to ECN.

Some encouragement to use ECN is certainly timely. As RFC 2481 notes: “Given the current effort to implement RED, we believe this is the right time for router vendors to examine how to implement congestion avoidance mechanisms that do not depend on packet drops alone. With the increased deployment of applications and transports sensitive to the delay and loss of a single packet (e.g., realtime traffic, short web transfers), depending on packet loss as a normal congestion notification mechanism appears to be insufficient (or at the very least, non-optimal).”

References and Further Reading

- [1] Huston, G., TCP Performance, *The Internet Protocol Journal*, Vol. 3, No. 2, Cisco Systems, June 2000.
- [2] Huston, G., *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, John Wiley & Sons, January 2000.
- [3] Postel, J., “Transmission Control Protocol,” RFC 793, September 1981.
- [4] Claffy, K., Miller, G., Thompson, K., “The Nature of the Beast: Recent Traffic Measurements from an Internet Backbone,” INET’98 Proceedings, Internet Society, July 1998. Available at:
http://www.isoc.org/inet98/proceedings/6g/6g_3.htm
- [5] Braden, R., “T/TCP—TCP Extensions for Transactions Functional Specification,” RFC 1644, July 1994.
- [6] Allman, M., Glover, D., Sanchez, L., “Enhancing TCP over Satellite Channels Using Standard Mechanisms,” RFC 2488, January 1999.
- [7] Jacobson, V., “Congestion Avoidance and Control,” ACM SIGCOMM, 1988.
- [8] Floyd, S., Fall, K., “Promoting the Use of End-to-End Congestion Control in the Internet,” Submitted to *IEEE Transactions on Networking*.
- [9] Allman, M., Paxson, V., Stevens, W., “TCP Congestion Control,” RFC 2581, April 1999.
- [10] Mogul, J., Deering, S., “Path MTU Discovery,” RFC 1191, November 1990.
- [11] Jacobson, V., Braden, R., Borman, C., “TCP Extensions for High Performance,” RFC 1323, May 1992.
- [12] Mathis, M., Mahdavi, J., Floyd, S., Romanow, A., “TCP Selective Acknowledgement Options,” RFC 2018, October 1996.
- [13] Allman, M., editor, “Ongoing TCP Research Related to Satellites,” RFC 2760, February 2000.
- [14] Wireless Access Protocol Forum, <http://www.wapforum.org>
- [15] Karn, P., Falk, A., Touch, J., Montpetit, M., Mahdavi, J., Montenegro, G., Grossman, D., Fairhurst, G., “Advice for Internet Subnet Designers,” work in progress, July 2000.
- [16] Jacobson, V., “Compressing TCP/IP Headers for Low-Speed Serial Links,” RFC 1144, February 1990.
- [17] Casner, S., Jacobson, V., “Compressing IP/UDP/RTP Headers for Low-Speed Serial Links,” RFC 2508, February 1999.

- [18] Stewart, R., et al., “Stream Control Transmission Protocol,” work in progress, July 2000.
- [19] Balakrishnan, H., Seshan, S., “The Congestion Manager,” July 2000.
- [20] Floyd, S., editor, “Congestion Control Principles,” work in progress, June 2000.
- [21] Padhye, J., Firoiu, V., Towsley, D., Kurose, J., Modeling TCP Throughput: A Simple Model and Its Empirical Validation, UMASS CMPSCI Tech Report TR98-008, Feb. 1998.
- [22] M. Mathis, M., Semke, J., Mahdavi, J., Ott, T., “The Macroscopic Behavior of the TCP Congestion Avoidance Algorithm,” *Computer Communication Review*, Vol. 27, No. 3, July 1997.
- [23] Ott, T., Kemperman, J., Mathis, M., “The Stationary Behavior of Ideal TCP Congestion Avoidance,” available at:
`ftp://ftp.bellcore.com/pub/tjo/TCPwindow.ps`
- [24] Floyd, S., Mahdavi, J., Mathis, M., Podolsky M., “An Extension to the Selective Acknowledgement (SACK) Option for TCP,” RFC 2883, July 2000.
- [25] Allman, M., Balakrishnan, H., Floyd, S., “Enhancing TCP’s Loss Recovery Using Early Duplicate Acknowledgment Response,” work in progress, June 2000.
- [26] Allman, M., “TCP Congestion Control with Appropriate Byte Counting,” work in progress, July 2000.
- [27] Ramakrishnan, K., Floyd, S., “A Proposal to Add Explicit Congestion Notification (ECN) to IP,” RFC 2481, January 1999.
- [28] Bradner, S., “The Internet Standards Process—Revision 3,” RFC 2026, October 1996.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. Huston is currently the Chief Scientist in the Internet area for Telstra. He is also a member of the Internet Architecture Board, and is the Secretary of the Internet Society Board of Trustees. He is author of *The ISP Survival Guide*, ISBN 0-471-31499-4, *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, and coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. All three books are published by John Wiley & Sons. E-mail: **`gih@telstra.net`**

Securing the Infrastructure

by Chris Lonwick, Cisco Systems

People are becoming much more reliant upon the proper operation of their networks. Consequently, the administrators of these networks are being tasked with providing an ever-increasing level of service. At this time of high reliance upon the network, methods and procedures need to be instilled into the network so the operators can maintain control of their network and they can know with some certainty the effect of each potential change. This may become increasingly difficult as network resiliency techniques are being proposed and deployed with the intent of automatically keeping these networks in top operation. Having a predictable network that is secured in a proper manner results in a network that is more suitable for the users and better meets the intended purpose of the network.

Most of the current network security models start with the physical perimeter of the network as its defining boundary. All things within this boundary are supposed to be protected from the perceived inimical forces that are outside of the perimeter. We are, however, finding that the perimeter of the network is no longer solidly defined. There are many exceptions to the “hard-shell perimeter” model—companies merge, remote sites are linked through *Virtual Private Networks* (Site-to-Site VPNs) across untrusted paths, access is granted in-bound for the network users through *Access Virtual Private Networks* (Access VPNs), and there are several other exceptions. For this article, let’s consider a different model. This model has a boundary of the acceptable network users rather than any geographical or logical perimeter. It is important that these users are allowed access to the services provided by the network. It is equally important that the people who are not authorized to use the network must be prevented from consuming its resources and otherwise disrupting its services.

Other models tend to focus on the restrictions of the users to access devices to provide security to the network. This model, however, looks at the effect that the users and each of the devices have upon the state of the network. To conceptualize this model, visualize that the only time this network would be running at a “steady state” is when there is no user traffic, no administrative or management traffic, and no routing update changes. The insertion of any traffic, or the addition or removal of any device or link, would change the state of this network. Changes to the state of this network may come from any number of sources, but they can be seen as coming from four different, quantifiable areas.

- Operators may enable or disable lines and devices.
- A network device publishing a new route or a different metric to a destination may cause the remainder of the network devices to dynamically recompute paths to all other destinations.
- Servers may insert traffic.
- Users may insert traffic.

Of these, the last two should be the least disruptive to the network as long as the traffic amounts are within the predicted and acceptable ranges. Changes that are within the goals of the network—for example to provide a service to the users—are considered good, while changes that cause outages or other disruptions are to be avoided. As such, it is vital that the network administrators understand the potential impact and consequences of each possible change in their network.

In this model, then, the administrators must know and understand the influences that will change the state of the network. The desire to achieve this goal sometimes leads to improper restrictions placed upon the users. Consider one extreme case of this model where each change in the network must be stringently authorized and authenticated. As a narrow example, this would mean that even traffic that is fundamentally taken for granted as a proper process of the network would have to be authenticated and authorized. *Domain Name System* (DNS) transactions would show that this extreme case is impractical. Each DNS query would have to be associated with a user or authenticated process, and that user or process would have to be authorized to make each specific query. A vastly more practical case for real networks would be for the administrators to allow any DNS query from any device without authentication—as it is done in existing dynamic networks today. In the model, the normal DNS queries and responses would be an influence upon the state of the network. For this influence to change the network in a way that meets the goals of the network, the administrators would have to feel comfortable that the servers and the available bandwidth will adequately handle the amount of DNS traffic as well as all other traffic. On the other hand, the administrators do need to establish a strict set of rules for the influences that they consider sensitive or possibly disruptive to their network. Continuing this example, the administrators may want to place restrictions upon the devices and processes that can insert and update the DNS records. It would be rather inappropriate, and potentially devastating, if any unauthorized person or network device were allowed to overwrite any existing records. If anyone were allowed to perform any DNS update that he or she wished, chaos would soon result. There must be a center position for this example that allows the operators to maintain control but still permits the dynamic freedoms expected by the users. Specifically to address this, the DNS Extensions Working Group has proposed several Internet Drafts^[1].

In the broader sense, this places a very heavy responsibility upon the people who are running the network. They must find some acceptable median between the desire to rigidly control all aspects of the network and the freedoms that are expected by the users, while at the same time satisfying the business requirements of their network. However, defining the freedoms and restrictions of the users is only one part of maintaining the network. The administrators and operators must have an understanding of the influences on the network as described in the model. In this, each aspect of the parts of the network must be under-

stood well enough to predict their behavior as they are normally used, and to limit the potential for disruption if they are used beyond their means. The one area that is vital to the proper working of the network is the infrastructure. This article explores some of the thoughts that may go into the process of securing the network infrastructure.

Table 1: Sources of Change to the Network

Sources of Change to the Network	Some Examples of How the Source Influences the Network	Examples of Device Types within the Network (The 4 Groups)	
Operators and their Devices	<ul style="list-style-type: none"> Add/remove new lines and circuits Install/remove network devices 		
	<ul style="list-style-type: none"> Login to the network devices to change their configuration Poll network devices for their status 	<ul style="list-style-type: none"> Operations Consoles Network Management Stations 	Operators
Network Devices	<ul style="list-style-type: none"> Dynamically route or switch traffic Dynamically mark lines and circuits in or out of service and then use them accordingly Authenticate users and permit their accesses accordingly Dynamically assign addresses and register that information for retrieval by others 	<ul style="list-style-type: none"> Routers and Switches Firewalls 	Infrastructure Devices
		<ul style="list-style-type: none"> Authentication Servers DNS/DHCP Servers 	
Servers	<ul style="list-style-type: none"> Servers send content to User's workstations to fulfill their requests Servers broadcast and multicast content to recipients 	<ul style="list-style-type: none"> Servers offering Content and Servers 	Servers
Users and their Devices	<ul style="list-style-type: none"> Client workstations request content from servers and upload content to servers Client workstations utilize services that are offered within the network 	<ul style="list-style-type: none"> Client Workstations 	Users
	<ul style="list-style-type: none"> A user encourages many others to visit a particular web site which causes a stampede A user tells others that a particular service is down or unavailable causing others to not attempt access 		

Description of Problem

In this abstracted network model, four sources of change were noted. As shown in Table 1, these changes, or influences to the network, may come from the operators, the network devices, the servers, and the users of the network. Let's first look at the influences that each of these groups can effect upon the network by first categorizing the network devices. All the devices on the network may be somewhat separated into four groups that correspond to the four sources. These groups of network devices can be seen in the third column of the table.

- *Operators:* For the purpose of this article, let's describe the Operators as all the people who operate the network, including the network engineers, the installers, the people who monitor the net-

work, and all the other people who make it work. The first group then is made of the operators and the devices that these operators use to run the network, such as the network management stations and all other operations consoles. Operators periodically make changes for moves and additions for better network performance, or to overcome disruptions. They will also monitor the network through polling, receiving alerts, and sometimes directly interacting with the network devices. Generally the amount of traffic inserted into the network from their activities is minimal. Because they generally have physical access to all locations, they can insert or remove network devices. Operators can have influence over all aspects of the network at all layers—from the physical layer, all the way up the stack. Operators can influence the network either in band or out of band, and they should be the only people who directly access the network infrastructure devices such as the routers and DNS servers. Usually this access will be from the management platforms, but in many situations, operators require access from devices that would otherwise be classified as a user's workstation.

- *Infrastructure Devices:* The network infrastructure devices themselves have the ability to change the network as well. This is mostly done through the dynamic nature of the network. At some times the physical portions of the network might fail and cause outages. In some cases, such as self-healing ring topologies, physical-layer devices may heal the network. In other cases, such as when a router is taken out of the network for maintenance, the routing updates will heal the network to the best of their abilities. The network infrastructure devices can be somewhat separated into two categories. The first of these would be the infrastructure devices that have no direct interaction with the users of the network. This category would consist of the devices such as the routers, switches, access control devices, and perhaps even the physical-layer devices such as multiplexers and modems. The user machines and content servers normally would not form sessions or require any information from these devices. The second category would be the devices with which customers indirectly interact. These would be devices such as the DNS servers, *Dynamic Host Configuration Protocol* (DHCP) servers, *Network Time Protocol* (NTP) servers, authentication servers, and the like. The users and servers would form sessions with these supporting devices and would require information from them for the basic operation of the network. In some cases, such as with a DNS/DHCP server, the results of the indirect user interaction would even update the servers with information. This latter group may be called “supporting devices.” These two categories can be taken together with all the wires, circuits, and lines to form the infrastructure of the network. Although the users do not actively see their presence, this infrastructure must be available and functioning before any user can actually do anything productive on the network.

- *Servers:* The servers in this group are those that contain content or services with which the users directly interact. These would be databases, Web servers, application servers, and the like. Like the operators group, this group is not considered to be part of the network infrastructure.
- *Users:* The users and their machines constitute the bulk of the network. The changes that the users make upon the network will probably come through transferring content or requesting and utilizing services. They can change the nature of the network by withdrawing from the network, or by causing others to withdraw from the network. In a nonmalicious way, the user base can degrade the state of the network by using it beyond its expected capacity. In certain situations, users with malicious intent may find exploitable network vulnerabilities. In most normal cases, however, the influence from the users upon the network will be through their interactions with the servers.

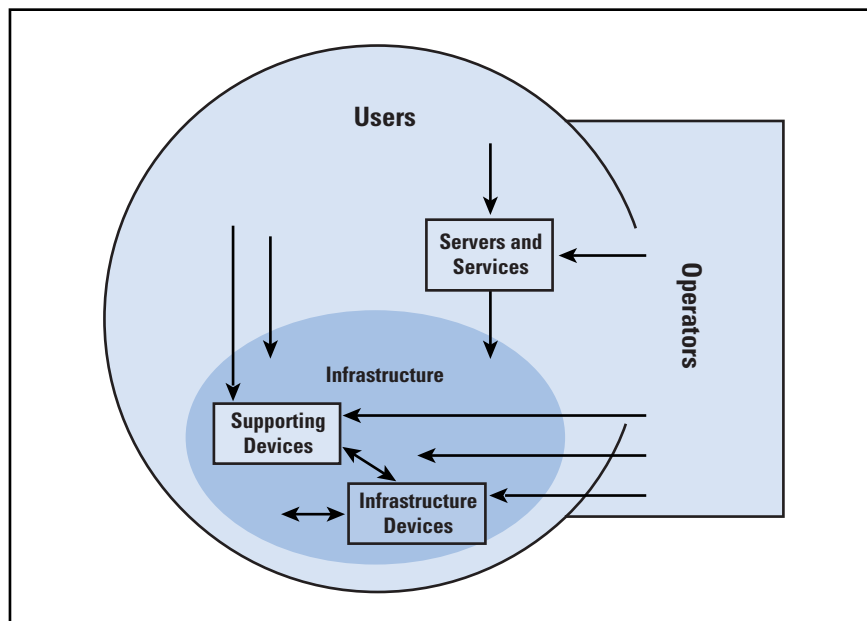
Each type of influence may also be considered to have a different weight. For example, the insertion of a new router into an existing network would be expected to have a larger effect upon the operations of the network than the change to the network caused by a user retrieving some information through a Web browser. To quantify some of the expected network changes, consider that there may be spheres and levels of influence. Any influence that may cause a change over the entire network may be considered to have a global sphere of influence. A router recently inserted into the network would start exchanging routing information with its neighbors. With no restrictions placed upon routing updates, this router could announce a new network, or it could announce the best path to an otherwise difficult-to-reach network. The remainder of the network would be affected, and all other routers would have to recalculate their paths. If the announcements were true, then the network would continue servicing the needs of the users. If the announcements were false, possibly because of an incorrect configuration, then the whole network could suffer. In this case, it is possible to limit the sphere of influence by restricting the acceptance of routing updates. In one method, all the routers could be restricted to disallow the acceptance of an announcement to the “default” network. Additionally, all the routers may be restricted to accept only announcements that are known to be within an acceptable address range. In another method, the routers could be grouped to accept announcements only from a select set of other routers. Additionally, some routing protocols have an option to include an authentication and integrity check through signing the updates. Any of these methods would help to reduce the sphere of influence and thus the potential for changes that could be made by the insertion of a router. There is, however, a cost associated with this; the operators would have to diligently enforce this control.

The level of influence can also be considered a factor in this model. The sphere of influence of a single transmission line can be defined to include any portion of the network that uses that line. If that line develops a fault, it may corrupt or discard packets and the associated network devices may automatically disable that line. If there is a backup line or an alternate path, then this change will be a small problem to the operations staff and its loss may go unnoticed by the users. That would be a low level of influence upon the network. On the other hand, if the line has an intermittent fault that can cause a route flap, or if the line has no backup, then major disruptions may occur. That would be considered a high level of influence.

If the goal of the network is to provide a service to its users, then its operators must try to quantify each of the influences. In a theoretically ideal network, the administrators would appropriately limit the sphere and would try to minimize the level for every influence. As was noted above, however, attempting to do this would require numerous operations tasks. Many of those may be unnecessary for their specific environment. For example, in a small business where there is a high degree of trust that no one has any malicious intent, controls would still be placed upon the influences that would most probably cause network problems through accidents. If the security policy allowed anyone to connect any device to the network, it may still be prudent to disallow the routers from receiving routing updates from any source other than the other routers.

A well-running network is the result of a well-controlled network. These networks must have a separation of authorized administration from other influences, and these other influences must be understood well enough to know how they will change the network. The following diagram shows the network and the groupings of influences upon it, and the table below that describes the elements of this model. This model does not show access paths, but rather the influences that each grouping of devices has upon the infrastructure and upon other devices. As can be seen, the users are pervasive throughout the network (because they are a principal reason for its existence), and they must have the access paths to contact the servers and necessary infrastructure devices. The users will influence the infrastructure as they insert traffic upon the lines, but they should have no direct influence upon the infrastructure devices such as the routers and digital access cross-connects. The operators do have influence upon the infrastructure devices and must have an access path to those devices. It would be most appropriate if the users were not allowed to usurp the access paths of the operators. However, because the two are sometimes nearly indistinguishable, the task of separating the administrative channels from the user channels becomes difficult.

Figure 1: The Network Security Model



The following table describes the elements in this model.

Table 2: Network Model Elements

Element	Description	Example
Operators	The devices and people who operate, manage and support the network	Monitoring and Management Workstations, Syslog servers
Infrastructure	This composite area denotes the entire infrastructure. This is broken out to show the actual infrastructure devices as well as the supporting devices.	<ul style="list-style-type: none"> Infrastructure Devices: Routers and Switches Supporting Devices: DNS and DHCP servers All other infrastructure components: wires, circuits, DSU/CSUs, SONET equipment, repeaters, etc.
Servers and Services	The devices that host content and services for the users	Web servers, file servers
Users	All of the users of the network and their workstations	Alice, Bob, Carol, Dan and their workstations
Arrows	Define which element influences or changes which other element	Users insert traffic into the network and thus influence the Servers and Services. Operators may also influence each of the components of the Infrastructure

Know Your Business

All well-running networks must have a *security policy* defined. This must reflect the goals of the network and must also be acceptable to the users and administrators. There are good examples of policies as well as methods that can be used to generate them. RFC 2196^[2] contains several thoughts about constructing a policy and SANS^[3] offers courses on this. While defining a network security policy, it will be advantageous to list the most likely disruptive influences to the network. This is commonly called *The Threat Model*. All potentially disruptive factors should be considered when forming the threat model, and they must be addressed when writing the security policy. It may, however, be beyond the capabilities of the operations staff to negate all of them. It may also be prohibitively expensive to try. In those cases, the writers of the policy should acknowledge the factors that won't be negated, but they should still find ways to minimize them. For example, in an Enterprise network the operators are somewhat likely to require access to routers and switches from any physical location in the network. In Service Provider networks, there may be less of a chance of that because the operators traditionally reside with the network management devices. In both cases, it would not be considered good for the network if a user could gain control of a router. The security policy for an Enterprise network may explain that network access to routers will be opened and available for any other device within the network. This will allow any operator to access the routers from any location. On the other hand, the security policy for a Service Provider network may state that access to the routers will be opened only for specific address ranges. Implementing this will prevent users, who reside within the address spaces assigned to the users, from accessing the infrastructure devices while allowing the operators, who reside within their own address space, to access the routers. In both cases, however, strong authentication will probably be required to additionally limit access to only authorized people.

Within the business of the network, operators must have the ability to control the infrastructure devices. Traditionally, the ways to interact with a device have been called "interfaces." A terminal with keyboard attached to the console port of a router is an interface just as is a Web browser accessing the router via the *Hypertext Transfer Protocol* (HTTP) through the network. Also, the path between the controlling device and the infrastructure device has traditionally been called a "channel." The wire connecting the terminal to the router is a channel just as is the TCP session that transports the HTTP in the prior example. The channels between the operators and the infrastructure devices must be secured, as well as the channels between the infrastructure devices. The first step in obtaining this goal is to identify all of the interfaces needed by the operations staff to access each of the remote devices. Along with this, they also need to identify each of the interfaces needed for the proper functioning of the network. The following lists are some of the possible network-available interfaces to some of the infrastructure devices in a dynamic network. This is somewhat broken down

into the interfaces needed by the operations staff, the interfaces usually needed by the other infrastructure devices, and some ancillary interfaces.

Table 3: Some Interfaces of Infrastructure Devices

Operations and Interfaces	Infrastructure Interfaces	Ancillary Interfaces
telnet, Kerberized telnet, SSH, rsh, rcmd, rexec, HTTP, FTP, tftp, rcp, scp, SNMP, LDAP, COPS, Finger	Syslog, ICMP, DNS, DHCP, RIP, OSPF, BGP, IS-IS, IGRP, EIGRP, HSRP, NTP, SNMP, Multicast controls	RADIUS, TACACS+, Kerberos Authentication, PAP, CHAP, EAP, chargen, echo, time, discard, Auth (Ident)

Each of these interfaces may be exposed to the nefarious forces that are known to inhabit large networks, and each of these exposures has vulnerabilities that may be exploited. Telnet sessions may be hijacked, DNS queries may be answered by nonauthoritative and possibly maliciously incorrect responses, and sinister people can insert forged routing updates to confound and disrupt the network. The network security policy should expect that these vulnerabilities may be exploited and it should address the mechanisms that may be used to either negate the vulnerabilities or to minimize the exposures. In this model, the process may be used to limit the sphere and level of the influences. The policy may also make some attempt to identify the potential consequences of the disruption caused by the exploitation of these interfaces. It should also describe an escalation procedure for dealing with encountered problems.

Possibly, during the exercise of identifying the open interfaces in an existing network, some of them may be closed or removed if it is determined that they are not needed or if their function can be fulfilled by the use of another interface. As an example, consider a UNIX host that has both the *Secure Shell Protocol* (SSH) and *finger* services running on it. If the policy of the network is to tightly control the information that anyone can obtain from any device, then the operators may want to remove the *finger* service. The operators will be able to obtain similar information by running the *who* command on the UNIX system through an SSH remote execution request. On the other hand, if the operations processes have been built upon the format of the information returned by *finger*, then the operators may want to prevent direct access to *finger* from the network and require that it be run on the device or through the SSH request.

At some point, it would be a good idea to run a scanner against the infrastructure devices. The *Network Mapper* (NMAP)^[4] is a freely available tool that can pick out some of the active interfaces of a device. This, or a similar tool, should be periodically used by the operations staff to ensure that the open ports of an infrastructure device are those that are known to be open. This investigation should not be limited to operations channels, but should also include application channels. For example, the question should be asked if the operations workstations should have open application interfaces—such as *Simple Mail Transfer Protocol* (SMTP) or *Network File System* (NFS). There are exploitable

vulnerabilities associated with some application interfaces that should be addressed in the security policy. In most cases, it would be prudent to remove applications that are not needed from infrastructure devices and supporting servers, as well as from operations devices when they are not needed. In all cases, it is usually considered to be a good practice to review the entries in the *inetd* configuration in UNIX systems.

It should be remembered that there will almost certainly be an access path between the users and the network interfaces of the infrastructure devices. The network security model diagram shows that neither the users nor the servers should have any direct influence to change or control the infrastructure devices. This is somewhat analogous to the policy of giving privileges on a multiuser system. In most well-run multiuser computing systems, the operators give only the most meager of privileges to the users of the system. This prevents most accidental and malicious disruptions. If the users need to run a privileged process or to access the files of other users, processes that utilize *setuid* are used or consensual groups are established. Generally, efforts are made to prevent users from having significant privileges on these machines. The alternative of giving each user high-level privileges usually results in disaster after a short time because the users then have the ability to overwrite or delete files, and may run processes that are generally disruptive to the operating system and to others.

Similarly, giving users high-level access to the routers of a network would have a deleterious effect. In the case of *Quality of Service* (QoS), a user given the privileges to reconfigure routers along a path would be able to provide his/her own designated flows with bandwidth and priority assurances. Subsequent users would also have that capability, and their modifications may leave the first user without his/her expected QoS—and possibly without a session at all. A far better mechanism to fairly deploy QoS is through the use of a brokering service. In a “policy network,” users or authenticated processes may request a level of service for their flows through a *Policy Manager*. This Policy Manager should have the capability to arbitrate requests to provide a semblance of fairness. The Policy Manager would then directly control the appropriate routers within the rules established by the administrators.

Along these lines, conveying security-related policy to infrastructure devices should take a similar path. For example, if the network security policy states that user access to a particularly sensitive network resource must be authenticated and controlled, the operators may elect to place a firewall between the users and that resource. That firewall would be classified as an infrastructure device and users should not directly access or control it. Rather, the users may authenticate themselves to an authentication service, which would notify the firewall that their access to the resource is permitted or denied. The authentication service may also send a set of restrictions for the access method; it may permit HTTP access but deny Telnet and *File Transfer Protocol* (FTP) for one person, but for another it may permit only Telnet.

The reasons for authentication, authorization, and access control must be described in the network security policy. It would be simple to mandate strict controls at many places in the network. However, that may not meet the needs of the business or the tolerance of the users. More to the point in this article is the requirement in the model that the disruptive influences be negated or minimized. Having a firewall or other access control device silently discard disruptive packets may be preferable to having a user or unconstrained process continue to spew garbage around the network.

Decide on the Methods of Securing the Channels and Interfaces

Some of the very first computing devices were designed to be managed locally and not remotely. Consoles consisting of a teletype device and a roll of paper were among the first interfaces to modern computing devices. Various methods were devised to extend these administrative interfaces beyond the confines of the frigid “Computer Room.” The first efforts were to keep these interfaces out of band, a scenario that meant separate wires from the physical port on the machine to a console in the operations room. In many cases, the wires from the remote terminals to the system were still visible because they were laid along the floor and could, therefore, be considered a secure channel. While this maintained a secure administrative channel—or path—that could not be tapped or exploited by others, it didn’t scale as more and more computing and ancillary devices were placed into the computer room, each requiring its own console. When remote terminals became commonplace, administrative functions were allowed over that channel. In almost all cases, the operating systems were mature enough to require some form of authentication before critical management operations were allowed.

The out-of-band channels for secure remote administration of devices may no longer be applicable to large networks. There are costs associated with running separate secure networks for the sole purpose of out-of-band operations, and there is the impracticality of one-at-a-time access through the console port of each device. This applies equally to the practice of placing a modem on the console ports of devices—a deployment that is not considered secure because there are still many automated dialers looking for answering modems. For these reasons, in-band access of operations has become the preferred method for modern networks. Telnet has been the oldest remote channel—and interface—for remote operations. Since then, other remote interfaces have been opened for controlling, commanding, and operating devices.

Many attempts have been made to “secure” Telnet and its use as a command and control channel. These efforts address the vulnerabilities of the protocol, and some address the interface itself. The *Berkeley Software Distribution* (BSD) “r” command set, such as *rlogin*, *rsh*, *rexec*, and others, were meant to be a substitute for the most common uses of Telnet within a trusted environment. It was assumed that the person initiating the command had previously been successfully authenticated.

SSH was meant to be a secure replacement of the Berkeley “r” tools. The SSH console session has been widely deployed to remotely operate devices. This replicated the Telnet interface while replacing the channel. The protocol addressed machine authentication, user authentication, and session confidentiality and integrity. When used as it was intended, it can effectively replace Telnet as a secure interface and channel. The *scp* feature of SSH can also securely replace *rcp*, and it has been used as a replacement for FTP. Likewise, a Kerberized Telnet and Kerberized FTP have been released to do the same.

Several other efforts have also been undertaken to secure some of the other administrative interfaces and channels. For example, the security issues of *Simple Network Management Protocol* (SNMP) are being addressed with the options of SNMPv3^[5]. Also, applications that utilize HTTP can be secured with HTTP over SSL (HTTPS) (*Secure Sockets Layer/Transport Layer Security* [SSL/TLS])^[6]. At this time, it appears that SSL/TLS is emerging as a mechanism that can be utilized to provide some security to many different applications. Beyond the operational interfaces and channels, work has been done to secure some of the infrastructure and ancillary interfaces. Some routing protocols have built-in authentication and integrity through the use of signing the routing updates with a shared key. Each mechanism that has been secured has been the subject of a focused effort to address that specific interface and channel. However, unlike those named above, some channels, such as Syslog and *Trivial File Transfer Protocol* (TFTP), have not been explicitly secured at this time.

IP Security (IPSec)^[7] was developed as a general-purpose mechanism that may be used to provide a secure wrapper around any unicast flow. Its cryptographic mechanisms can provide strong authentication, confidentiality, and integrity. While IPSec can be used to secure any flow, it may require additional infrastructure. A *Public Key Infrastructure* (PKI) must be established within the network. The alternative is to use preshared keys, a solution that is operationally intensive and doesn’t scale well. IPSec also requires consistent time synchronization between the devices, as well as a consistent DNS. If these pieces are in place, the operations staff can utilize IPSec to secure each of the needed operations channels. If the operators and administrators choose this method, then they should ensure that the unsecured channels are unavailable to anyone but themselves. For example, if the Telnet channel is secured with IPSec, then the Telnet port on remote devices should be closed for inbound access.

One method of closing the exposures is through *Access Control Lists*. Routers and switches usually have mechanisms that can be used to allow inbound and outbound sessions from only certain devices. UNIX devices usually have the ability to run TCP wrappers that can provide access-control mechanisms for inbound and outbound sessions. If infrastructure devices can be grouped together, the operators may decide to

place them behind an internal firewall. The decision to do that should be thought through. Generally the internal firewall will limit access of the protected devices to the specified interfaces^[8]. If this is done for a group of network management stations, the net effect may be that any attempts to access those workstations from outside of the firewall would be denied. The only inbound flows may be SNMP responses and traps. This implementation would limit the operations staff to being physically present before they could operate those devices. On the other hand, the firewall would prevent users from mistakenly or intentionally forming sessions with those devices. Because any received packet would have to be assessed by the device, a firewall that would discard packets before they are received by the device would help to prevent denial-of-service attacks. The use of internal firewalls should not be used as an excuse for poor security measures on the protected devices. Regardless of how effective the operators feel their firewall is, the protected devices must be treated as if they were otherwise exposed.

In determining the channels that will be used for the administration of the infrastructure devices, the packages will also be selected. At this time, many devices are sporting Telnet, FTP, and HTTP channels and the operators may utilize workstations that have these packages already loaded onto them. Also, networks comprising Microsoft NT servers may be managed remotely by the NT administrative tools, which commonly run on NetBIOS over TCP/IP (NBT). When given the choice, most often the operations staff will select easy-to-use and commonly available packages to access the interfaces of the infrastructure devices for remote operations and control. In all cases, these will be packages that will be available to the user community of the network as well. The users of the network may also easily download packages of these types if they don't already have them on their machines. For example, the operators may choose to utilize SSH for secured access to some devices. It is a trivial task for the users to also download an SSH client package and to start poking around the network to see what they can find. Even SNMP packages can be easily downloaded to the workstations of the users.

The operators and administrators must avoid the temptation to select a less-well-known package for infrastructure management based upon the thought that the users probably won't know about it. Users may not be initially aware that some packages are being used, but they can also download sniffer packages. Given enough time, even passive sniffing will give them enough clues to determine the channels used for administration. When they know that, they can then probably download the package themselves, and may then attempt to use it to explore the network. It should also be noted that the more heavily used packages have been scrutinized much more than the newer or less used packages. As a very general rule, the older a package gets, the more it becomes trusted because more people have been using it and *probably* attempting to break it.

As described above, and as it is seen in the diagram of the model, some of the channels that are available to the operators are also available to the users. This means that if the operators utilize Telnet to control their routers, it may be possible for a user to also initiate a Telnet session to a router. There must be an extremely strong discriminator to differentiate between the authorized operators and the unauthorized users before access to control the device is granted. Almost exclusively, the discriminator used is some form of authentication. An operator should be able to satisfy an authorization challenge, whereas an unauthorized user should not. A username and password is the most common form of in-band authentication. Specifically within Telnet and FTP, an in-stream challenge is presented to the user attempting a session; the user is asked for a username and then for a password. If these credentials match the values stored on the host, then the session is permitted. In these sessions, the credentials are exposed to casual observation. Anyone with a packet-sniffing device will be able to plainly see the username and password. These credentials must be regarded as secrets that must be protected. If they are compromised or stolen, then the operators have lost their control of their network. Some packages, such as SSH and Kerberos, have addressed these problems and have found ways to prevent secrets from being passed during authentication.

It must also be noted that some infrastructure devices do not offer any in-band channels for control. Many *Channel Service Units/Data Service Units* (CSU/DSUs) are not IP aware and do not offer any in-band channels for control. In cases like those, physical access may be the discriminator that prevents unauthorized users from controlling the device. Typically, a lock on a door or a cabinet would be the “challenge,” and the key would be the authentication credential, which must be treated like a secret. It cannot be emphasized enough that these secrets must be protected. The *CERT Coordination Center* has written a very broad Tech Tip, which explores the topic of password security^[9]. Many companies have found it very beneficial to periodically hold training courses to highlight the importance of this subject both to their operators and to their users.

Ancillary Channels Also Require Security

One of the parallel problems with using authentication credentials is its distribution. Many devices are capable of maintaining a local database of usernames and passwords. However, maintaining identical databases on each device throughout large networks is infeasible. More often, the authentication credentials are stored in a centralized database and an *Authentication, Authorization, and Accounting* (AAA) protocol is used to transfer them as needed. The AAA protocols most often used are *Remote Access Dial-In User Service* (RADIUS), TACACS+, and *Kerberos* authentication. Each of these has different characteristics and security mechanisms. Kerberos authentication was designed to securely transport authentication material. A password is never transferred across the network in this architecture. This protocol has withstood the test of

time, but it has been difficult to establish in networks that aren't committed to maintaining it. This situation seems to be changing because more "productized" versions are becoming available on the market. TACACS+ has a mechanism to hide the exchanges between the TACACS+ client and the server. It is also capable of transferring authorization rules for each user. RADIUS uses a mechanism to hide portions of the exchange between the RADIUS client and server as well.

Beyond this, the channels for telemetry, audit, and accounting may need to be secured. There are no inherent mechanisms to secure syslog at this time, and SNMPv1 may be protected with a Community String, but that solution is considered weak. It is possible to allow read-only access to the SNMP interface, but SNMPv3 has many of the security features that have been requested to secure this protocol. Other channels that are required by the operations staff should also be critically reviewed because many forms of attacks are on open channels.

It would be appropriate for the operations staff to keep up with new exploits and to assume that the users of the network have access to the latest "hacker" tools. It is quite common for people to hear about an exploit or published vulnerability and then "try it out" in the nearest available network. For this reason, it should be in the security policy of the network that "security patches" be given the highest priority and should be loaded on the affected platforms as soon as they are available and have been approved for the environment.

Conclusions

When any security mechanism is applied, the appropriateness and applicability of the solution should be questioned. On the surface, some security solutions may appear to be good; however, their applicability to the situation must be verified. As an example, SSL may be used to secure HTTP traffic, and it is commonly found in many Web browsers. Unfortunately, not many people explore the browser options that are enabled by default. In most browsers, SSLv2 is still available, even though it has published and exploitable vulnerabilities. Additionally, even in SSLv3—which negates the vulnerabilities of SSLv2—low key-length cipher suits are still available and enabled by default. In many cases, a null-cipher crypto algorithm is available. In the internal networks of many companies, SSL may be selected and implemented using a self-signed certificate. Care must be taken to ensure that this certificate is the one distributed to each administrative workstation. SSL sessions may be formed without certificates supplied by either endpoint. An attacker could exploit this through a man-in-the-middle attack. Another example would be the use of SSH. SSHv1 has known vulnerabilities. If the administrators decide to deploy SSH for the control of the remote infrastructure devices, they should first decide if they should be worried about attacks against those known vulnerabilities in their infrastructure. If they are, then they should either deploy SSHv2, which addresses the vulnerabilities of SSHv1, or they should explore the use of Telnet with IPsec.

In many cases, rather than using the “most secure” solution, perhaps a simpler solution would still provide adequate protection. The “most secure” solution—the one that mitigates all perceived threats—is usually too costly to implement. In many cases, network operators and administrators with many years of experience have decided that SSHv1 is adequate for their needs and they can mitigate or minimize the exposure. In other cases, some operators are turning to SSHv2 or IPSec to cover the vulnerabilities that have been found in SSHv1. In some cases, the use of SNMPv1 may also be acceptable as long as its exposures are understood and the operators determine that its use will not pose a problem.

Excessive “security” may also intolerably reduce the usability of the network. It is important to remember that the network is there for the users. Placing security restrictions upon them to keep them out of the infrastructure is like keeping the doors locked to the building boiler room. Untrained people entering that area may hurt themselves or they may cause serious problems to others. If they have malicious intent, they could damage the machinery. Excessive security for that analogy would be similar to locking the boiler room, locking the ingress and egress points to the building, and mandating that armed guards accompany anyone that is permitted to enter the building. In some cases, that may be appropriate for the perceived threat. However, in the case that this applies to an elementary school building, it is inappropriate and would make some parents think of moving their children to other schools.

The model described in this article may be used as a thought process to review an entire network at a high layer to see the relationships between the various devices. It may also be used to design the security policy and the acceptable use policy of the network. Another use for it may be to define the operational procedures for the operators to securely administer the network and to define how the infrastructure devices will communicate. However it is used, some settlements must be made between the desire to provide security and the usefulness of the network. The cost of the security mechanisms cannot be unreasonably high, and the mechanisms cannot change the business model of the company. The enforcement of the policy must be effective, yet above all it must not change the expectations of the users. In all cases, the administrators and operators must find some balance between their need to secure the infrastructure and the need for the users to have the ability to actually use their network.

References

- [1] Internet Engineering Task Force DNS Extensions Working Group, last updated July 2000,
<http://www.ietf.org/html.charters/dnsext-charter.html>
- [2] Fraser, B., "Site Security Handbook," RFC 2196, September 1997.
- [3] System Administration, Networking, and Security Institute,
<http://www.sans.org/>
- [4] Fyodor <fyodor@dhp.com>, "NMAP—The Network Mapper,"
<http://www.insecure.org/nmap/index.html>
- [5] Stallings, William, "Security Comes to SNMP: The New SNMPv3 Proposed Internet Standards," *The Internet Protocol Journal*, Vol. 1, No. 3, December 1998.
- [6] Stallings, William, "SSL: Foundation for Web Security," *The Internet Protocol Journal*, Vol. 1, No. 1, June 1998.
- [7] Stallings, William, "IP Security," *The Internet Protocol Journal*, Vol. 3, No. 1, March 2000.
- [8] Avolio, Fred, "Firewalls and Internet Security," *The Internet Protocol Journal*, Vol. 2, No. 2, June 1999.
- [9] CERT® Coordination Center, Tech Tips, "Protecting Yourself from Password File Attacks," Last revised February 12, 1999.

CHRIS LONVICK holds a Bachelor of Science degree from Christian Brothers College and is in the Consulting Engineering Department of Cisco Systems in Austin, Texas. He is currently the chair of the IETF Syslog Working Group. Chris can be reached at clonvick@cisco.com

Book Reviews

Multiwave Optical Networks *Multiwavelength Optical Networks: A Layered Approach*, by Thomas E. Stern and Krishna Bala, ISBN 020130967X, Addison-Wesley, 1999.

Initial Impressions

This book attempts to fit into two camps; one, an overview of the potential choices that could be offered in wavelength-division multiplexing, or WDM, and the other, an academic text. Because of its scope, the treatment is uneven.

Organization

The first four chapters lay the groundwork. Chapter 1 starts by defining terms and positing why WDM is an enabling technology. The authors believe that the driving application will be LAN interconnection, ostensibly in metro areas. It is worthwhile noting that the authors make no claims about this text relating to an all-optical network. They simply expose the choices available to manipulate the various wavelengths, or lambda. The current methods for performing lambda manipulation are still bound in the electrical domain.

Chapter 2 covers the hierarchy or layering present in a WDM environment and some of the choices for configuration at each point in the hierarchy. The authors spend some time on the concepts of spectrum partitioning and what routing and switching in this domain means. A key point raised relates to the concept of wavelength conversion at network access points. The chapter closes with a brief review of some types of logical overlays that may sit on top of a WDM network. Three types are examined, ATM, *Synchronous Optical Network* (SONET), and IP networks.

The third chapter covers how network interconnection may occur and how the management and control features may be implemented. Four basic topologies are described, each with its salient features highlighted. These topologies include shared channel networks; wavelength routed networks, linear lightwave networks, and hybrid, logically routed networks. It is interesting to note that many commercial implementations, especially from traditional telecom providers, tend to follow the simpler topologies, while we are beginning to see newer telecom providers utilizing the more robust topologies.

Chapter 4 discusses what the authors consider enabling technology. To a large degree, these enabling technologies are the basic components of an optical system, for example, fibers, amplifiers, transmitters, and receivers. Crosstalk is mentioned in particular. The authors then delve into photonic device technologies and wavelength converters, and then they close with some simulation work on end-to-end transmission paths.

Chapters 5, 6, and 7 discuss in depth the ramifications of each of the four techniques. What is fairly intriguing here is that the authors have extensive bibliographies at the end of each chapter, and they include a series of problems that are left as an exercise to the reader.

The eighth chapter touches on the concepts involved with survivability and restoration of service. This chapter should help the practical network engineer in understanding most of the possible failure modes. In the last chapter, the authors look at current trends, and they try to predict business drivers for WDM deployment. Once again, they show their true colors as academics when they close with a statement on the importance of testbeds.

On to the Appendices! I am grateful to the authors for including some basic material on graph theory, scheduling algorithms, Markov chains and queuing, some work on minimal interference routing in the optical domain and, finally, close with a synopsis of the SONET standard.

Good Reference

Overall, there is a fair amount of practical material here, but it is tucked into large amounts of academic detail. I'm not sure this volume would work as a standalone textbook, but it clearly is a good reference for the state of optical networks in the last years of the 20th century.

—Bill Manning,
University of Southern California
Information Sciences Institute
manning@isi.edu

Net Slaves *Net Slaves: True Tales of Working the Web*, Bill Lessard and Steve Baldwin, ISBN 0-07-135243-0, McGraw-Hill, 2000.

How can you not want to read a book that opens with a quote from a Guns&Roses song, “Do you know where you are? You’re in the jungle, baby!”? *Net Slaves* is about the people who maintain the jungle that big game hunters come to exploit. The same jungle marketed as the digital age and the e-generation. This is the land of the “dot-coms” and future big-buck IPOs. Has hubris masked your role in this jungle? *Net Slaves* will set you straight. Exactly who are these net slaves? Well, take the 15 question quiz provided by the authors and determine your Internet exploitation quotient. Don’t be shocked to find yourself among the new media caste; the only question is, what part of the jungle are you assigned to clean after?

The authors spent a year interviewing people who work for Internet-based companies. Based on their findings, they created 11 character composites: Garbagemen; Cops or Streetwalkers; Social Workers; Cab Drivers; Cowboys or Card Sharks; Fry Cooks; Gold Diggers or Gigolos; Priests or Madmen; Robots; Robber Barons; and Mole People.

For each composite the authors cite someone's real-life work experience—of course, in order to protect the innocent (and the guilty), names have been altered.

I was annoyed with David Zorn, Card Shark; his type does nothing but give the industry a bad reputation. The story of Ken Hussein, Robot, both saddened and angered me. I truly hope he and his family are doing better. How can anyone not feel sorry for Kellner after being taken in by Gigolo Mira? Jane, Cab Driver, learned the hard way that you have to roll with the blows to survive in the jungle. Finally, I must confess, I found the most disturbing of all profiles to be of Outis, a Mole Person.

For each profile the authors provide some social-economic statistics. How old is the average Social Worker? How much does it cost to hire a Cowboy? What are the career aspirations of the average Cab Driver? How do you know if a Robot is annoyed with you? You're a Garbage-man; what are your chances of upward mobility? A lot of this is funny, but to leave it at that would be missing the point entirely. Every composite represents scores of real people's lives, and how they live doesn't necessarily match up with the glamour often associated with the high-tech industry.

My favorite profile is of Jason Barstow, a Madman. Barstow arrives on the scene on his Harley, ready to participate in a two-day seminar put on by the Earth Business Network. A former chicken farmer and former guitar player, Barstow now finds himself lecturing to a room full of CEOs. He begins by telling them about the 5 milligrams of LSD he bought the previous night, and proceeds to plant seeds of anxiety—did he spike their morning juice? As Barstow delivers his lecture on the future of e-commerce and builds to the climax, a frustrated Slim Clarkston of NetScathe blurts out, "Mr. Barstow, I want you to tell us the truth about your little prank." With the lecture over, Barstow returns his pass to the security desk. "How did it go?" asks the security guard. "Same bull," Barstow responds, "but they never seem to get tired of it."

Are these stories true? I don't know—it doesn't matter! What are true are the composites. This book is funny. It is also humbling. Most important, it is true. It was fun to read. After each chapter, I found myself wearing an undeniable mischievous grin as I scanned the office looking for the person I just read about; this is all in good fun as long as I remember one important thing: I'm in the book—and you are too. In my experiences, I've found that a certain animosity always exists between people who work call centers, programmers, Web designers, managers, and the like. *Net Slaves* reminds us that we are all in this jungle together.

—Neophytos Iacovou, eBenX Inc
diacovou@ebenx.com

Implementing IPsec *Implementing IPsec: Making Security work on VPNs, Internets, and Extranets*, Elizabeth Kaufman and Andrew Newman, ISBN 0-471-34467-2 Wiley Computers Publishing, 1999.

Organization

The book is organized into four parts. The first three chapters of Part One should be nothing more than review for anyone who has been in networking for even a short time. Chapter 4, “Encrypting within the Law,” analyzes current worldwide regulatory trends for encryption technologies and examines how existing laws will impact your ability to legally purchase and install IPsec products. Included is some good information that may help keep you on the right side of the laws pertaining to encryption. Encryption is an area of potential problems, especially when you are running your network between countries.

Part Two is a primer on the basic technological components of IPsec. Chapter 5, “A Functional Overview of IPv4,” and its basic design characteristics should be old news to anyone who is seriously thinking of running any type of encryption on his/her network. Chapter Six is an overview of cryptographic technologies. Chapter 7 “The Basics of IPsec and Public Key Infrastructures (PKIs) Fundamental to Current IPsec Standards,” has some good information pertaining to IPsec and its different components, but leaves out an explanation of its two basic modes of operation: *transport* and *tunnel*.

Part Three analyzes how and why the IPsec protocols can break existing IP networks, and should provide the reader with some good information. Chapter 8, “What Won’t Work with IPsec,” describes the root cause of IPsec performance problems and protocol conflicts. Chapter 9, “IPsec and PKI Rollout Considerations,” discusses gateway-to-gateway, end host-to-gateway, and end host-to-end host configuration options and explains some of the policy elements of PKI.

Part Four provides some criteria for evaluating vendors and products; this information would be of little interest if you are unfamiliar with writing an RFI. Also included is some reference material, including an appendix, with a complete copy of the IPsec RFC (2401), “Security Architecture for the Internet Protocol.” A glossary, which does *not* offer a description of IPsec, is included as well.

Who Should Read This Book

By trying to appeal to the technical as well as the nontechnical reader, the book has missed both. There are areas that will appeal to the reader with a limited networking background, as well as areas for the more technical. However, if you are the type of reader inclined to read the RFCs, you will find very little reason to read the remainder of the book. Overall the book does not provide enough information for any one group. Inclusion of RFC 2401 seems unnecessary considering how easily RFCs can be obtained from the Internet.

—Al Pruitt, CSG Systems, Inc
al_pruitt@csgsystems.com

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

Scott Bradner Receives Postel Service Award

The Internet Society (ISOC) recently announced that noted Internet standards leader and Internet pioneer Scott O. Bradner has been awarded the prestigious *Jonathan B. Postel Service Award* for 2000. In presenting the award, Geoff Huston, Chair of ISOC, said, “Scott Bradner was introduced to many of us with his accurate and careful measurements of router performance. He has been a long standing participant in the *Internet Engineering Task Force* (IETF), and continues to serve on the *Internet Engineering Steering Group* (IESG) as the Area Director for Transport. He was a ISOC Trustee for six years from 1993 until 1999 and continues to serve as the Society’s Vice-President for Standards. This is an impressive set of contributions and is worthy of recognition in Jon Postel’s name as the 2000 recipient of the Jonathan B. Postel Service Award.”

Don Heath, president and CEO of ISOC, said, “We established the award to honor the late Jon Postel by recognizing his unselfish and substantial contributions to the Internet over a 25 year period.” He added, “Scott Bradner exemplifies the spirit of all that Jon brought to the Internet community and his outstanding contributions have made this year’s choice an easy one. Scott’s careful judgment and good humor has been a major contribution to many of the ISOC’s activities, and we are pleased to be able to recognize his contributions in this unique fashion.”

Bradner has been an active contributor to the IETF for over a decade, and has served as a Working Group Chair, the Area Director for Operations and currently serves as the Area Director for Transport. He also was the Director of the IPv6 area, and oversaw the process of refinement of a number of proposals into the definition of a coherent architecture for IPv6. Bradner has been the prime author of the current Internet Standards Process documents. He has also been an instructor at ISOC’s *Network Training Workshops for Developing Countries* for many years, and has been a catalyst for the development of operationally robust Internet services in many areas of the world.

The Award is named for Dr. Jonathan B. Postel, an Internet pioneer and head of the organization that administered and assigned Internet names, protocol parameters, and Internet Protocol (IP) addresses. He was the primary architect behind what has become the *Internet Corporation for Assigned Names and Numbers* (ICANN), the successor organization to his work. The Award is presented at the Internet Society’s annual INET Conference. It consists of an engraved crystal globe and US \$20,000.00. Scott Bradner becomes the second recipient of the award. The first was presented posthumously to Dr. Postel in 1999.

The Internet Society is a non-profit, non-governmental, open membership organization whose worldwide individual and organization members make up a veritable “who’s who” of the Internet industry. It provides leadership in technical and operational standards, policy issues, and education. ISOC hosts two annual Internet conferences, trains people from all over the world in networking technologies, conducts workshops for educators, and publishes an award-winning magazine, *OnTheInternet*. ISOC provides an international forum to address the most important economic, political, social, ethical and legal initiatives influencing the evolution of the Internet. This includes facilitating discussions on key policy decisions such as taxation, copyright protection, privacy and confidentiality, and initiatives towards self-governance of the Internet. ISOC created the Internet Societal Task Force as an on-going forum for discussion, debate, and development of position papers, white papers, and statements on Internet related societal issues.

ISOC is the organizational home of the IETF, the Internet Architecture Board, the IESG, and the Internet Research Task Force—the standards setting and research arms of the Internet community. These organizations operate in an environment of bottom-up consensus building made possible through the participation of thousands of people from throughout the world. For more information, see <http://www.isoc.org/>

APNIC Policy Meeting

The *Asia Pacific Network Information Centre* (APNIC) will host an Open Policy Meeting October 25–27, 2000 in Brisbane, Australia. The meeting is open to anyone with an interest in Internet addressing issues. For more information see: <http://apnic.org>

APRICOT 2001

The *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) will be held in Kuala Lumpur, Malaysia, February 26 to March 2, 2001. APRICOT is a forum that facilitates knowledge sharing among key Internet builders in the region, with peers and leaders from the Internet community worldwide. Since 1996, APRICOT has established itself as Asia Pacific’s premier regional Internet Summit where related organisations converge and host their annual general meetings and other special events. The week-long summit comprises seminars, workshops, tutorials, conference sessions, Birds of a Feather (BOFs), and other forums, all geared towards spreading and sharing the knowledge required to operate the Internet within the Asia Pacific region. For more information see: <http://www.apricot2001.net>

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Technology
WorldCom, USA

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Edward R. Kozel, Member of The Board of Directors
Cisco Systems, Inc., USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

Copyright © 2000 Cisco Systems Inc.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-10/5
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

Bulk Rate Mail U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol *Journal*

December 2000

Volume 3, Number 4

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

F R O M T H E E D I T O R

In This Issue

From the Editor	1
The Trouble with NAT	2
The Social Life of Routers	14
New Frontiers for Research Networks.....	26
Book Review.....	40
Call for Papers	41
Fragments	42

Numerous technologies have been developed to protect or isolate corporate networks from the Internet at large. These solutions incorporate security, either end-to-end (IP security, or IPSec), or at the Internet/intranet border (firewalls). A third class of systems allows a range of IP addresses to be used internally in a corporate network, while preserving IP address consumption through the use of a *single* public address. This latter class of device is called a *Network Address Translator* (NAT), and while many Internet engineers consider NATs to be “evil,” they are nonetheless very popular. Combining IPSec, NATs, and firewalls can be quite challenging, however. In our first article Lisa Phifer explains the problem and offers some solutions.

Successful network design is the result of many factors. In addition to the basic building blocks of routers, switches and circuits, network planners must carefully consider how these elements are interconnected to form an overall system with as few single points of failure as possible. In our second article, Valdis Krebs looks at how lessons learned from social network analysis can be applied to the design of computer networks.

The current Internet grew out of several government-funded research efforts that began in the late 1960s. Today, basic technology development as well as research into new uses of computer networks continues in many research “testbeds” all over the world. Bob Aiken describes the past, present and future state of network research and research networks.

The online subscription system for this journal will be up and running in January at www.cisco.com/ipj. In addition to offering a subscription form, the system will allow you to select delivery options, update your mailing and e-mail address, and much more. Please visit our Web site and give it a try. If you encounter any difficulties, please send your comments to ipj@cisco.com.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

The Trouble with NAT

by Lisa Phifer, Core Competence

Those who are implementing virtual private networks often ask whether it is possible to safely combine *IP Security* (IPSec) and *Network Address Translation* (NAT). Unfortunately, this is not a question with a simple “yes” or “no” answer. IPSec and NAT can be employed together in some configurations, but not in others. This article explores the issues and limitations associated with combining NAT and “NAT-sensitive” protocols like IPSec. It examines configurations that do not work, and explains why. It illustrates methods for using NAT and IPSec together, and discusses an emerging protocol that may someday prove more IPSec friendly.

This article builds upon “IP Security and NAT: Oil and Water?”^[1] and “Realm-Specific IP for VPNs and Beyond”^[2], works previously published by *ISP-Planet*.

What Is Network Address Translation?

NAT was originally developed as an interim solution to combat IPv4 address depletion by allowing globally registered IP addresses to be re-used or shared by several hosts. The “classic” NAT defined by RFC 1631^[3] maps IP addresses from one realm to another. Although it can be used to translate between any two address realms, NAT is most often used to map IPs from the nonroutable private address spaces defined by RFC 1918^[4], shown below.

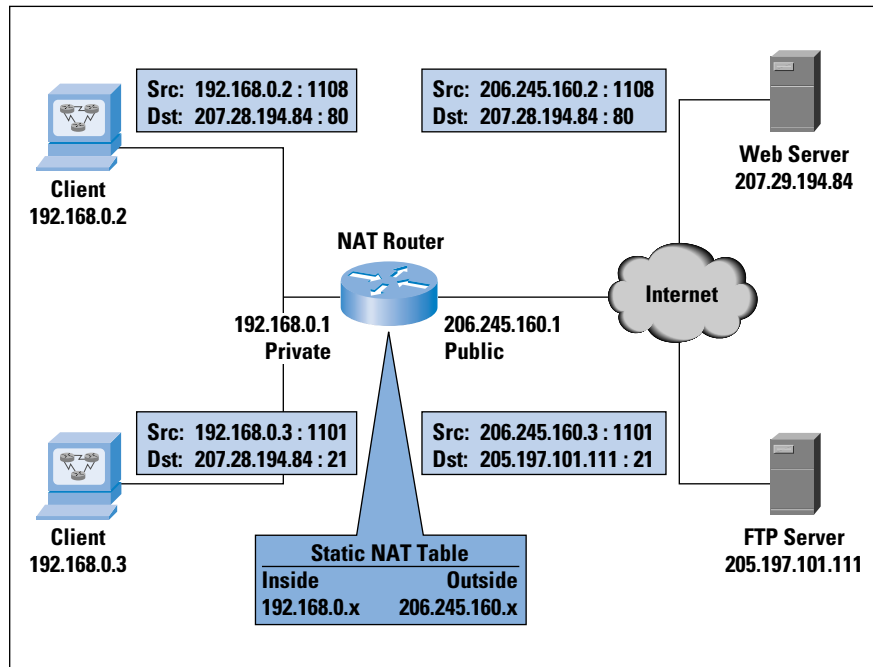
Class	Private Address Range
A	10.0.0.0 ... 10.255.255.255
B	172.16.0.0 ... 172.16.255.255
C	192.168.0.0 ... 192.168.255.255

These addresses were allocated for use by private networks that either do not require external access or require limited access to outside services. Enterprises can freely use these addresses to avoid obtaining registered public addresses. But, because private addresses can be used by many, individually within their own realm, they are nonroutable over a common infrastructure. When communication between a privately addressed host and a public network (like the Internet) is needed, address translation is required. This is where NAT comes in.

NAT routers (or NATifiers) sit on the border between private and public networks, converting private addresses in each IP packet into legally registered public ones. They also provide transparent packet forwarding between addressing realms. The packet sender and receiver (should) remain unaware that NAT is taking place. Today, NAT is commonly supported by WAN access routers and firewalls—devices situated at the network edge.

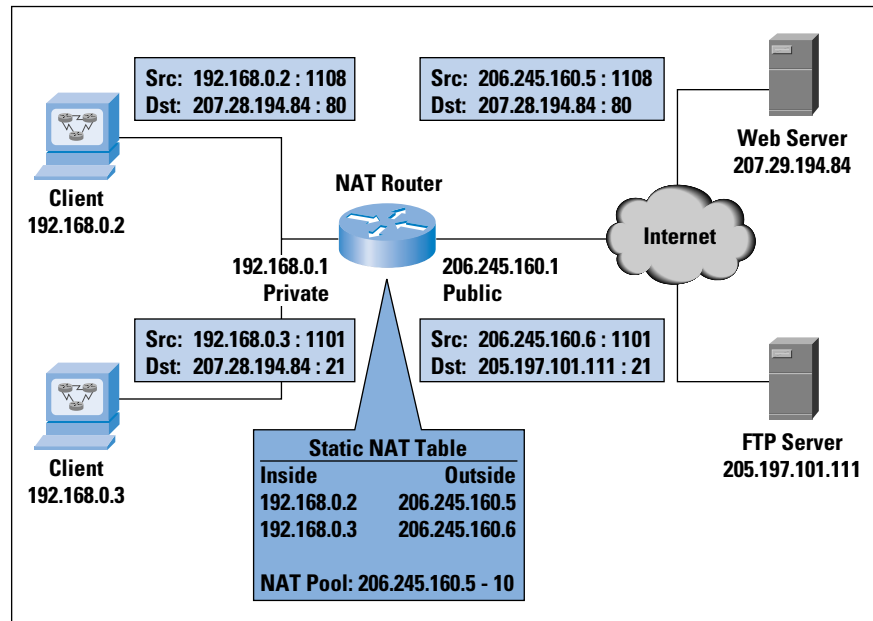
NAT works by creating bindings between addresses. In the simplest case, a one-to-one mapping may be defined between public and private addresses. Known as static NAT, this can be accomplished by a straightforward, stateless implementation that transforms only the network part of the address, leaving the host part intact. The payload of the packet must also be considered during the translation process. The IP checksum must, of course, be recalculated. Because TCP checksums are computed from a pseudo-header containing source and destination IP address (prepended to the TCP payload), NAT must also regenerate the TCP checksum.

Figure 1: Static NAT



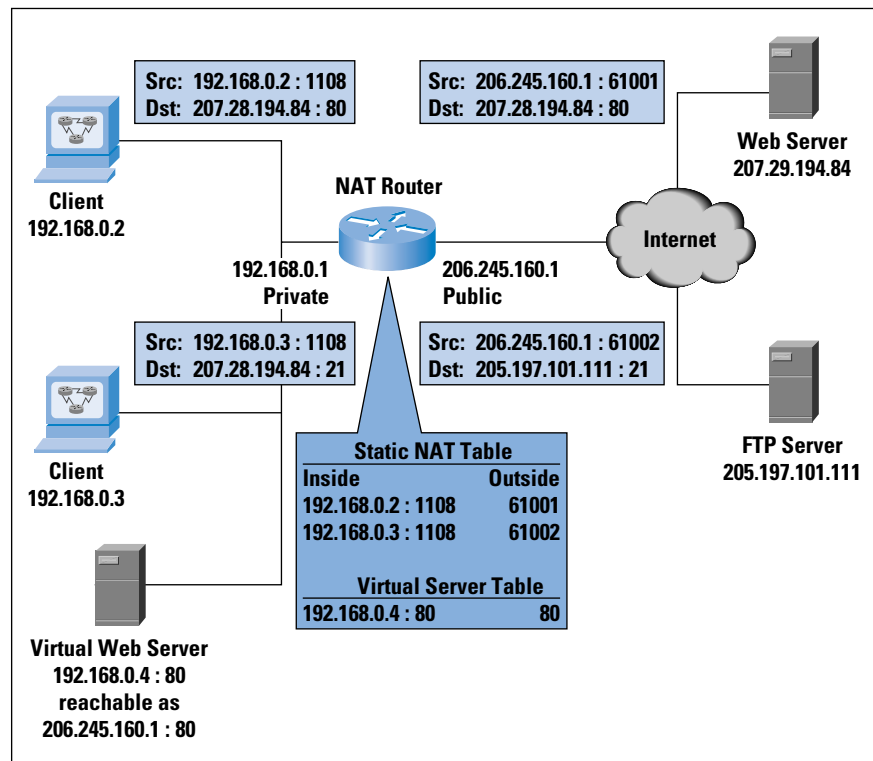
More often, a pool of public IP addresses is shared by an entire private IP subnet (dynamic NAT). Edge devices that run dynamic NAT create bindings “on the fly,” building a NAT Table. Connections initiated by private hosts are assigned a public address from a pool. As long as the private host has an outgoing connection, it can be reached by incoming packets sent to this public address. After the connection is terminated (or a timeout is reached), the binding expires, and the address is returned to the pool for reuse. Dynamic NAT is more complex because state must be maintained, and connections must be rejected when the pool is exhausted. But, unlike static NAT, dynamic NAT enables address reuse, reducing the demand for legally registered public addresses.

Figure 2: Dynamic NAT



A variation of dynamic NAT known as *Network Address Port Translation* (NAPT) may be used to allow many hosts to share a single IP address by multiplexing streams differentiated by TCP/UDP port number. For example, suppose private hosts 192.168.0.2 and 192.168.0.3 both send packets from source port 1108. A NAPT router might translate these to a single public IP address 206.245.160.1 and two different source ports, say 61001 and 61002. Response traffic received for port 61001 is routed back to 192.168.0.2:1108, while port 61002 traffic is routed back to 192.168.0.3:1108.

Figure 3: NAPT



NAPT (masquerading) is commonly implemented on small Office/Home Office (SOHO) routers to enable shared Internet access for an entire LAN through a single public address. Because NAPT maps individual ports, it is not possible to “reverse map” incoming connections for other ports unless another table is configured. A virtual server table can make a server on a privately addressed DMZ reachable from the Internet via the public address of the NAPT router (one server per port). This is really a limited form of static NAT, applied to incoming requests.

In some cases, static NAT, dynamic NAT, NAPT, and even bidirectional NAT or NAPT may be used together. For example, an enterprise may locate public Web servers outside of the firewall, on a DMZ, while placing a mail server and clients on the private inside network, behind a NAT-ing firewall. Furthermore, suppose there are applications within the private network that periodically connect to the Internet for long periods of time. In this case:

- Web servers can be reached from the Internet without NAT, because they live in public address space.
- *Simple Mail Transfer Protocol* (SMTP) sent to the private mail server from the Internet requires incoming translation. Because this server must be continuously accessible through a public address associated with its *Domain Name System* (DNS) entry, the mail server requires static mapping (either a limited-purpose virtual server table or static NAT).
- For most clients, public address sharing is usually practical through dynamically acquired addresses (either dynamic NAT with a correctly sized address pool, or NAPT).
- Applications that hold onto dynamically acquired addresses for long periods could exhaust a dynamic NAT address pool and block access by other clients. To prevent this, long-running applications may use NAPT because it enables higher concurrency (thousands of port mappings per IP address).

Where is NAT used today? Outbound NAT is commonly employed by multihost residential users, teleworkers, and small businesses that share a single public IP for outbound traffic while blocking inbound session requests. In other words, small LANs connected via ISDN, *Digital Subscriber Line* (DSL), or cable modem.

Bidirectional static NAT/NAPT combinations are typically used by enterprises that host services behind a masquerading firewall. NAT can also be employed by enterprises wishing to insulate themselves from *Internet Service Provider* (ISP) address changes, or by those wanting to obscure private network topology for security reasons.

NAT-Sensitive Protocols

Our need to conserve IPv4 addresses has prompted many to overlook the inherent limitations of NAT, recognized in RFC 1631 but deemed acceptable for a short-term solution.

As noted previously, NAT regenerates TCP checksums. This, of course, requires the TCP header containing the checksum to be visible (that is, not encrypted). If only the TCP payload is encrypted and immutable between the application source and destination (for instance, *Secure Shell Protocol* [SSH], *Secure Sockets Layer* [SSL]), then the checksum in the TCP header can be recalculated without a visible TCP payload. But if the TCP header is encrypted (for instance, IPsec transport mode), the TCP checksum field in the TCP header cannot be modified.

Furthermore, many application protocols carry IP addresses in an application-level protocol. In such cases, an *Application-Level Gateway* (ALG) is needed to complete the translation. For example:

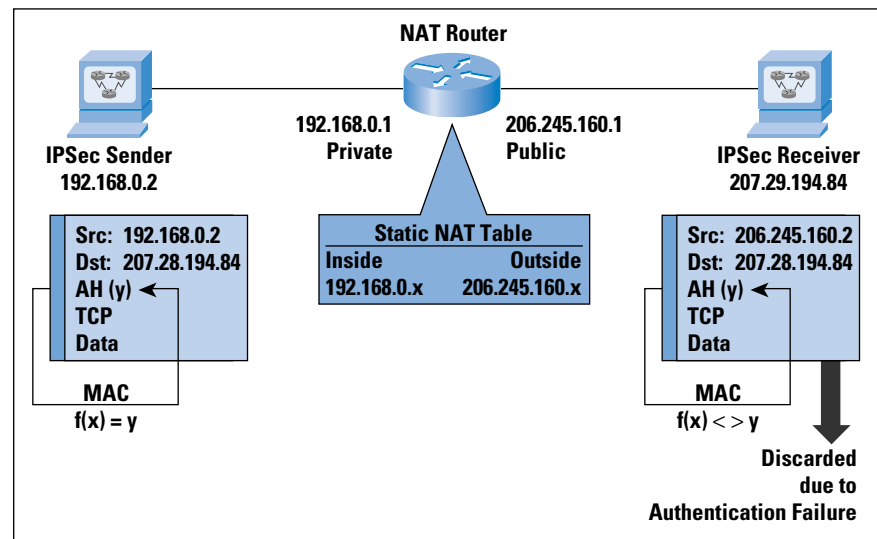
- Many *Internet Control Message Protocol* (ICMP) packets (for instance, “Destination Unreachable”) carry embedded IP packets in ICMP payload. These require both address translation and checksum regeneration.
- A *File Transfer Protocol* (FTP) ALG is needed to rewrite IP addresses carried by FTP PORT and PASV control commands. In the IP header, these addresses are fixed-length words. Unfortunately, in the FTP protocol, these IP addresses are carried as human-readable, variable-length strings; rewriting can change the length of the TCP segment. If the segment is shortened, it can be padded. If the segment is lengthened, SEQ and ACK numbers must be transformed for the duration of the connection.
- Protocols like H.323 use multiple TCP connections or UDP streams to form “session bundles.” If all connections in the bundle originate from the same end system, an ALG may be avoided. But H.323 presents other challenges, including ephemeral ports and embedded, ASN.1-encoded IP addresses in application payload.
- *NetBIOS over TCP/IP* (NBT) can be challenging to translate correctly because packet-header information is placed in NetBIOS payload at inconsistent offsets, and many embedded IP addresses are exchanged during an NBT session. Fortunately, most companies do not let NBT beyond their firewall anyway.
- *Simple Network Management Protocol* (SNMP) packets also carry IP addresses that identify trap source and object instance. Perhaps more important, dynamic NAT makes it impossible to uniquely identify hosts by IP address; public addresses are transient and shared. Remote management of private hosts can thus be impeded by NAT.
- Obviously DNS, responsible for domain name/IP address mapping, is impacted by NAT. From simple query handling to zone transfers, a robust DNS ALG is defined by RFC 2694^[9].

NAT-sensitive protocols such as Kerberos, X-Windows, remote shell, Session Initiation Protocol (SIP), and others are further described in the Internet Draft “*Protocol Complications with the IP Network Address Translation*”^[12]. Another Internet Draft, “*NAT Friendly Application Design Guidelines*”^[13], explains how new application protocols can integrate smoothly with NAT. But there are still cases where ALGs simply cannot “fix” packets modified by NAT.

Impact of NAT on IPSec

The IPSec *Authentication Header* (AH)^[5] is an example. AH runs the entire IP packet, including invariant header fields such as source and destination IP address, through a message digest algorithm to produce a keyed hash. This hash is used by the recipient to authenticate the packet. If any field in the original IP packet is modified, authentication will fail and the recipient will discard the packet. AH is intended to prevent unauthorized modification, source spoofing, and man-in-the-middle attacks. But NAT, by definition, modifies IP packets. Therefore, AH + NAT simply cannot work.

Figure 4: NAT vs. AH
(Transport Mode)

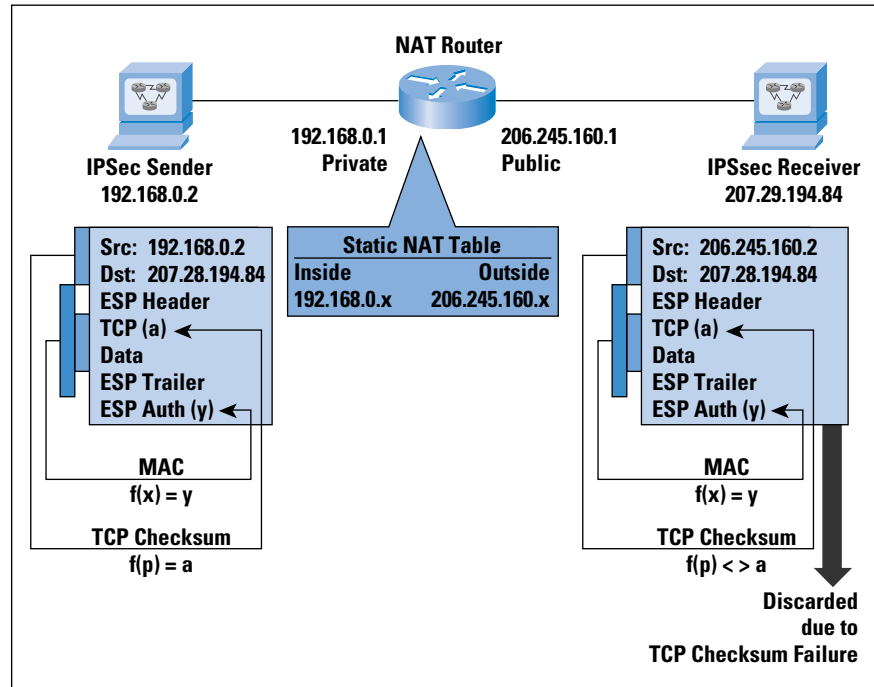


The IPSec *Encapsulating Security Payload* (ESP)^[6] also employs a message digest algorithm for packet authentication. But, unlike AH, the hash created by ESP does not include the outer packet header fields. This solves one problem, but leaves others.

IPSec supports two “modes.” Transport mode provides end-to-end security between hosts, while tunnel mode protects encapsulated IP packets between security gateways—for example, between two firewalls or between a roaming host and a remote access server. When TCP or UDP are involved—as they are in transport mode ESP—there is a catch-22. Because NAT modifies the TCP packet, NAT must also recalculate the checksum used to verify integrity. If NAT updates the TCP checksum, ESP authentication will fail. If NAT does not update the checksum (for example, payload encrypted), TCP verification will fail.

If the transport endpoint is under your control, you might be able to turn off checksum verification. In other words, ESP can pass through NAT in tunnel mode, or in transport mode with TCP checksums disabled or ignored by the receiver.

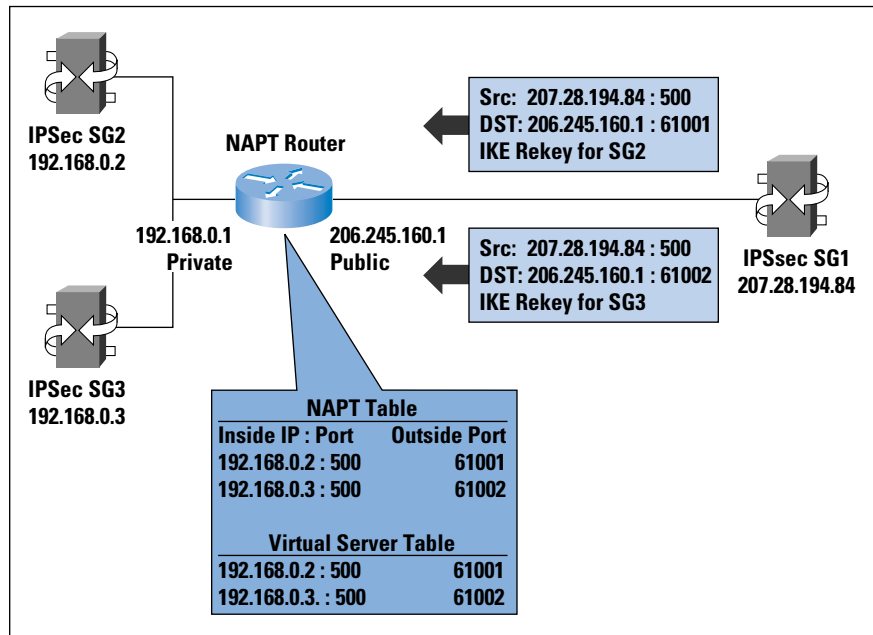
Figure 5: NAT vs. ESP
(Transport Mode)



If we stick to ESP in tunnel mode or turn off checksums, there's still another obstacle: the *Internet Key Exchange* (IKE)^[7]. IPsec-based *Virtual Private Networks* (VPNs) use IKE to automate security association setup and authenticate endpoints. The most basic and common method of authentication in use today is preshared key. Unfortunately, this method depends upon the source IP address of the packet. If NAT is inserted between endpoints, the outer source IP address will be translated into the address of the NAT router, and no longer identify the originating security gateway. To avoid this problem, it is possible to use another IKE “main mode” and “quick mode” identifier (for example, user ID or fully qualified domain name).

A further problem may occur after a *Security Association* (SA) has been up for awhile. When the SA expires, one security gateway will send a rekey request to the other. If the SA was initiated from the well-known IKE port UDP/500, that port is used as the destination for the rekey request. If more than one security gateway lies behind a NAT router, how can the incoming rekey be directed to the right private IP address? Rekeys can be made to work by “floating” the IKE port so that each gateway is addressable through a unique port number, allowing incoming requests to be demultiplexed by the NAT router.

Figure 6: NAT vs.
IKE Rekey



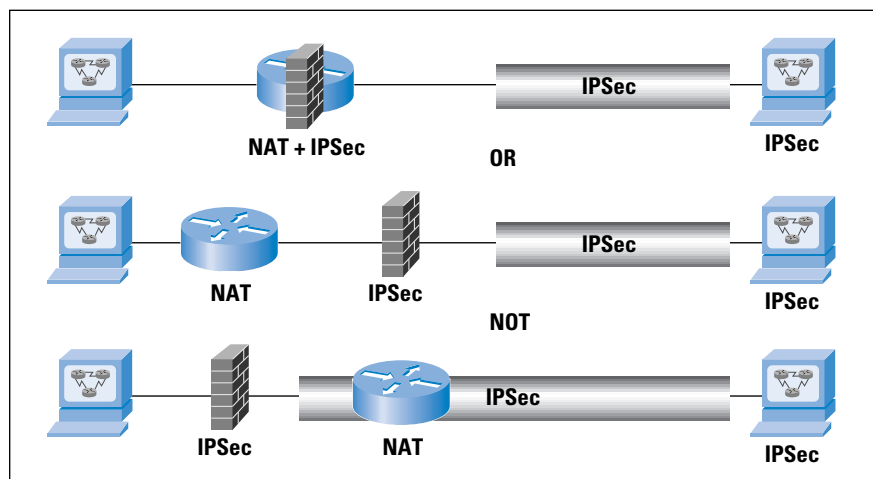
At this point, two things should be clear: (1) it is possible to find a “flavor” of IPSec that will run through NAT, but (2) one must do so with great care and attention to detail. Recent Internet Drafts^{[12] [14]} have recorded these problems for further consideration, and RFC 2709^[10] describes a security model for running tunnel-mode IPSec through NAT.

One Solution: Avoid the Problem

By far the easiest way to combine IPSec and NAT is to completely avoid these problems by locating IPSec endpoints in public address space. That is, NAT before IPSec; don’t perform IPSec before NAT. This can be accomplished in two ways:

- Perform NAT on a device located behind your IPSec security gateway; or
- Use an IPSec device that also performs NAT.

Figure 7: Combining
IPSec and NAT



Many routers, firewalls, security gateways, and Internet appliances implement IPSec and NAT in the same box. These products perform outbound address translation before applying security policies; the order is reversed for inbound packets. A typical “any-to-any” security policy is easily specified with such a product. Granular policies can be a bit more difficult because filters are often based on IP address, and care must be taken to avoid overlapping filters.

If you cannot avoid translating IPSec-protected traffic midstream, limit use of IPSec to tunnel-mode ESP and design security policies with care. If you simply cannot NAT before IPSec or require transport-mode ESP, there may still be hope. The *Internet Engineering Task Force* (IETF) is now defining *Realm-Specific IP* (RSIP), an alternative that may someday prove kinder to IPSec.

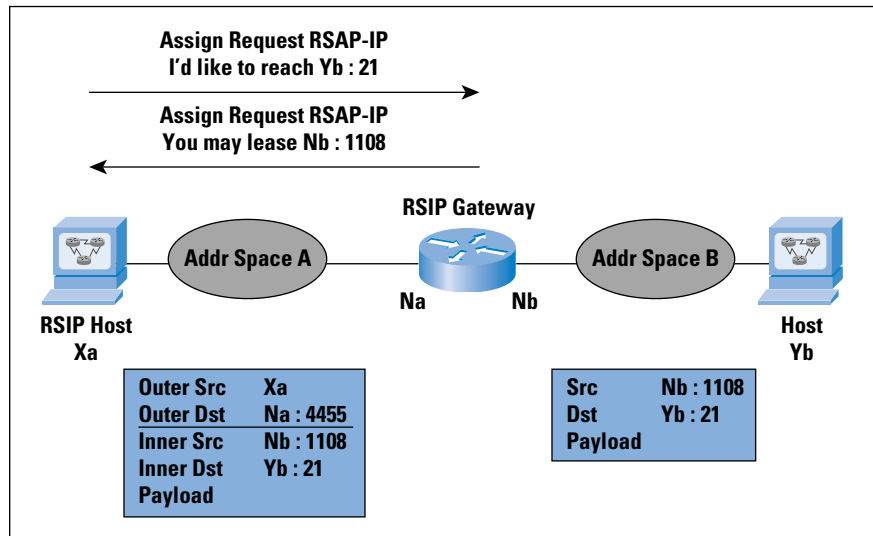
What Is RSIP?

RSIP^[16] leases public IP addresses and ports to RSIP hosts located in private addressing realms. Unlike NAT, RSIP does not operate in stealth mode and does not translate addresses on the fly. Instead, RSIP allows hosts to directly participate concurrently in several addressing realms. Although RSIP does require host awareness, it avoids violating the end-to-end nature of the Internet. With RSIP, IP payload flows from source to destination without modifications that cripple IPSec AH and many other NAT-sensitive protocols.

RSIP gateways are multihomed devices that straddle two or more addressing realms, just as NAT-capable firewalls and routers do today. When an RSIP-savvy host wants to communicate beyond its own private network, it registers with an RSIP gateway. The RSIP gateway allocates a unique public IP address (or a shared public IP address and a unique set of TCP/UDP ports) and binds the private address of the RSIP host to this public address. The RSIP host uses this public source address to send packets to public destinations until its lease expires or is renewed.

But the RSIP host cannot send a publicly addressed packet as-is; it must first get the packet to the RSIP gateway. To do this, the host wraps the original packet inside a privately addressed outer packet. This “encapsulation” can be accomplished using any standard tunneling protocol: IP-in-IP, the *Generic Routing Encapsulation* (GRE), or the *Layer 2 Tunneling Protocol* (L2TP). Upon receipt, the RSIP gateway strips off the outer packet and forwards the original packet across the public network, toward its final destination.

Figure 8: RSIP



For simplicity, we talk about RSIP linking one private network to the public Internet, but RSIP can also be used to relay traffic between several privately addressed networks. An RSIP host can lease several different addresses as needed to reach different destinations networks. We've also focused on outgoing traffic, but an RSIP host can ask the RSIP gateway to "listen" and relay incoming packets addressed to a public IP and port.

Combining RSIP and IPSec

At first glance, RSIP sounds like a promising way for hosts to share public addresses while avoiding the pitfalls associated with applying NAT to IPSec traffic. But it turns out that RSIP extensions are needed to accommodate end-to-end IPSec^[17].

Basic RSIP relies on unique port numbers to demultiplex arriving packets, but IPSec ESP encrypts port numbers. When several RSIP hosts use the same RSIP gateway to relay ESP, another discriminator is needed. Fortunately, every IPSec packet carries a unique *Security Parameters Index* (SPI), assigned during security association setup. Unfortunately, the SPI is guaranteed unique only for the responder. To enable demultiplexing, the tuple (SPI, protocol [AH or ESP], destination IP address) must also be unique at the initiating RSIP gateway.

A similar problem occurs during association setup with the IKE. IKE packets usually carry the well-known source port UDP/500. Using different source ports is the preferred solution, but if several RSIP hosts use the same RSIP gateway to relay IKE from port UDP/500, another discriminator is needed. Again, there is a convenient answer: every IKE packet carries the initiator cookie supplied in the first packet of an IKE session. The RSIP gateway can route IKE responses to the correct RSIP host using the tuple (initiator cookie, destination port [IKE], destination IP address). But rekeys may still be an issue.

To fix these problems, extensions have been proposed to allow RSIP hosts to register with an RSIP gateway for IPSec support, and allow hosts to request and receive unique SPI values along with leased IP addresses and ports.

Possible Applications for RSIP

RSIP specifications^{[16][17][18]} are still at the Internet Draft stage. If and when RSIP matures, there may be a wide variety of applications:

- Residential power users and teleworkers with multihost LANs that share a single, publicly known IP address leased by an RSIP-enabled Internet appliance, DSL router, or cable modem;
- Small-to-midsize enterprise customers with dozens or hundreds of hosts, sharing a small pool of public IPs leased by an RSIP-enabled WAN access router or firewall;
- Multidwelling units (apartments, shared office buildings) with many private LANs, sharing public Internet access through an RSIP-enabled device;
- Hospitality networks (airports, hotels) where roaming hosts briefly lease the public IP(s) shared by the entire network;
- Remote access concentrators that use RSIP to lease private IP(s) to roaming corporate users that access the Internet via dynamically assigned public addresses; and
- Wireless devices (cell phones, personal digital assistants [PDAs]) that lease public IP(s) for “sticky sessions” that persist even when the mobile device moves from one location to another, updating its local access IP.

These scenarios, and the relationship of RSIP to IP multicast and differentiated services, are more fully explored in the RSIP framework^[18].

Conclusion

Although NAT can be combined with IPSec and other NAT-sensitive protocols in certain scenarios, NAT tampers with end-to-end message integrity. RSIP—or whatever RSIP evolves into—may someday prove to be a better address-sharing solution for protocols that are adversely impacted by NAT. If RSIP fails to mature, another solution may be developed to broaden use of NAT with IPSec. Alternatives now under discussion within the IETF include UDP encapsulation and changes to IKE itself^{[14][15]}.

Despite its origin as a short-term solution, NAT is unlikely to disappear in the very near future. Until it does, understanding the relationship between NAT and IPSec and alternatives for safe combined deployment will remain an important aspect of VPN design.

References

- [1] Phifer, L., "IP Security and NAT: Oil and Water?" *ISP-Planet*, June 15, 2000.
http://www.isp-planet.com/technology/nat_ipsec.html
- [2] Phifer, L., "Realm-Specific IP for VPNs and Beyond," *ISP-Planet*, June 23, 2000.
<http://www.isp-planet.com/technology/rsip.html>
- [3] Egevang, K. and Francis, P., "The IP Network Address Translator (NAT)," RFC 1631, May 1994.
- [4] Rekhter, Y., Moskowitz, B., Karrenberg, D., de Groot, G.J., and Lear, E., "Address Allocation for Private Internets," RFC 1918, February 1996.
- [5] Kent, S. and Atkinson, R., "IP Authentication Header," RFC 2402, November 1998.
- [6] Kent, S. and Atkinson, R., "IP Encapsulating Security Payload (ESP)," RFC 2406, November 1998.
- [7] Harkins, D. and Carrel, D., "The Internet Key Exchange (IKE)," RFC 2409, November 1998.
- [8] Srisuresh, P. and Holdrege, M., "IP Network Address Translator (NAT) Terminology and Considerations," RFC 2663, August 1999.
- [9] Srisuresh, P., Tsirtsis, G., Akkiraju, P. and Heffernan, A., "DNS Extensions to Network Address Translators (DNS_ALG)," RFC 2694, September 1999.
- [10] Srisuresh, P., "Security Model with Tunnel-Mode IPSec for NAT Domains," RFC 2709, October 1999.
- [11] Tsirtsis, G. and Srisuresh, P., "Network Address Translation-Protocol Translation (NAT-PT)," RFC 2766, February 2000.
- [12] Srisuresh, P. and Holdrege, M., "Protocol Complications with the IP Network Address Translator," Internet Draft, Work in Progress, July 2000.
- [13] Senie, D., "NAT Friendly Application Design Guidelines," Internet Draft, Work in Progress, July 2000.
- [14] Aboba, B., "NAT and IPSec," Internet Draft, Work in Progress, July 2000.
- [15] Stenberg, M., Paavolainen, S., Ylonen, T., and Kivinen, T., "IPSec NAT-Traversal," Internet Draft, Work in Progress, July 2000.
- [16] Borella, M. and Lo, J., "Realm-Specific IP: Protocol Specification," Internet Draft, Work in Progress, March 2000.
- [17] Montenegro, G. and Borella, M., "RSIP Support for End-to-End IPSec," Internet Draft, Work in Progress, March 2000.
- [18] Borella, M., Lo, J., Grabelsky, D., and Montenegro, G., "Realm-Specific IP: Framework," Internet Draft, Work in Progress, March 2000.

LISA PHIFER is vice president of Core Competence, Inc. (www.corecom.com), a consulting firm specializing in Internet, network management, and security technologies. She earned her Master's Degree in Computer Science from Villanova University. A Bellcore award recipient for her work in ATM network operations, Lisa has been involved in the design and deployment of networking protocols for over 18 years. She represented Bellcore and Unisys in several industry-standards organizations, and has participated in The Internet Security Conference (TISC) since its inception. Lisa consults, teaches, and writes about a variety of technologies, including caching, load balancing, DSL, ISDN, IPSec, PKI, OSS, and VPNs. Her monthly column on virtual private networking is published by *ISP-Planet*. E-mail: lisa@corecom.com

The Social Life of Routers

Applying Knowledge of Human Networks to the Design of Computer Networks

by Valdis Krebs

We often forget that computer networks are put in place to support human networks—person-to-person exchanges of information, knowledge, ideas, opinions, insights, and advice. This article looks at a technology that was developed to map and measure human networks—social network analysis—and applies some of its principles and algorithms to designing computer networks. And as we see more peer-to-peer (P2P) models of computer-based networks, the P2P metrics in human network analysis become even more applicable.

Social network analysts look at complex human systems as an interconnected system of nodes (people and groups) and ties (relationships and flows)—much like an internetwork of routers and links. Human networks are often unplanned, emergent systems. Their growth is sporadic and self-organizing^[1]. Network ties end up being unevenly distributed, with some areas of the network having a high density of links and other areas of the network sparsely connected. These are called “small world networks”^[2]. Computer networks often end up with similar patterns of connections—dense interconnectivity within subnetworks, and sparser connections uniting subnetworks into a larger internetwork.

Social network researchers and consultants focus on *geodesics*—shortest paths in the network. Many of today’s social network algorithms are based on a branch of mathematics called *graph theory*. Social network scientists have concentrated their work, and therefore their algorithms, in the following areas:

- Individual node centrality within a larger network—network dependency and load upon individual routers
- Overall path distribution—good connectivity without excessive routing tables
- Improving communication flow within and between groups—designing better topologies
- Network patterns surrounding ego networks—strategies for analyzing and manipulating individual router connections
- Analyzing information flow behavior of client organization—how computer networks can support human networks

One of the methods used to understand networks and their participants is to evaluate the location of actors in the network. Measuring the network location is finding the *centrality* of a node^[3]. All network measures discussed here are based on geodesics—the shortest path between any two nodes. We will look at a social network, called the *kite network*, that effectively shows the distinction between the three most popular centrality measures—the ABCs—Activity, Betweenness, and Closeness.

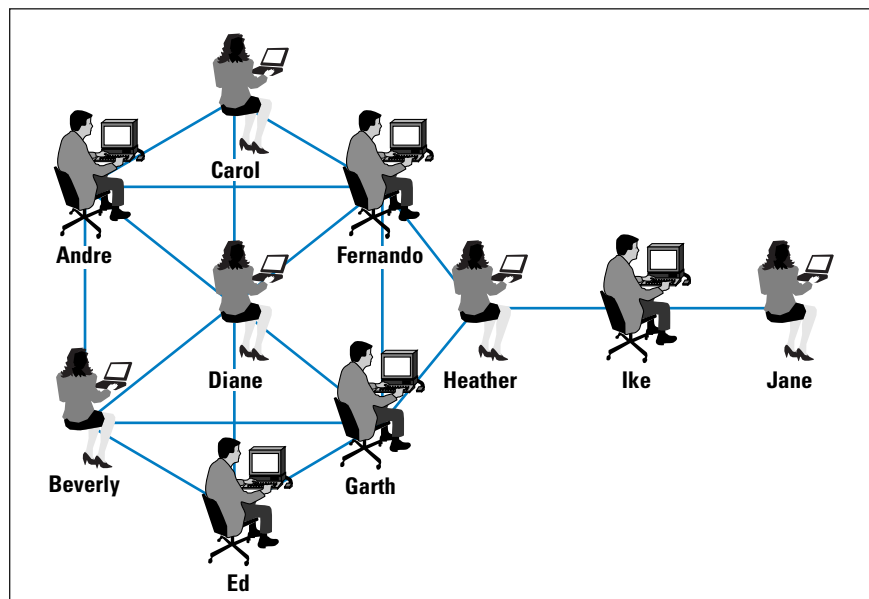
This model^[4] was first developed by David Krackhardt, a leading researcher in social networks.

Activity

Figure 1 shows a simple social network. A link between a pair of nodes depicts a bidirectional information flow or knowledge exchange between two individuals. Social network researchers measure network activity for a node by using the concept of *degrees*—the number of direct connections a node has.

In this human network, Diane has the most direct connections in the network, making hers the most active node in the network with the highest degree count. Common wisdom in personal networks is “the more connections, the better.” This is not always so. What really matters is where those connections *lead to*—and how they connect the otherwise unconnected!^[5] Here Diane has connections only to others in her immediate cluster—her clique. She connects only those who are already connected to each other—does she have too many redundant links?

Figure 1: Human Network



Betweenness

While Diane has many direct ties, Heather has few direct connections—fewer than the average in the network. Yet, in many ways, she has one of the best locations in the network—she is a boundary spanner and plays the role of broker. She is *between* two important constituencies, in a role similar to that of a border router. The good news is that she plays a powerful role in the network, the bad news is that she is a single point of failure. Without her, Ike and Jane would be cut off from information and knowledge in Diane’s cluster.

Closeness

Fernando and Garth have fewer connections than Diane, yet the pattern of their ties allow them to *access* all the nodes in the network more quickly than anyone else. They have the shortest paths to all others—they are *close* to everyone else. Maximizing closeness between *all* routers improves updating and minimizes hop counts. Maximizing the closeness of only one or a few routers leads to counterproductive results, as we will examine below.

Their position demonstrates that when it comes to network connections, quality beats out quantity. Location, location, location—the golden rule of real estate also works in networks. In real estate it is geography—your physical neighborhood. In networks, it is your virtual location determined by your network connections—your network neighborhood.

Network Centralization

Individual network centralities provide insight into the individual's location in the network. The relationship between the centralities of all nodes can reveal much about the overall network structure. A very centralized network is dominated by one or a few very central nodes. If these nodes are removed or damaged, the network quickly fragments into unconnected subnetworks. Highly central nodes can become critical points of failure. A network with a low centralization score is not dominated by one or a few nodes—such a network has no single points of failure. It is resilient in the face of many local failures. Many nodes or links can fail while allowing the remaining nodes to still reach each other over new paths.

Average Path Length in Network

The shorter the path, the fewer hops/steps it takes to go from one node to another. In human networks, short paths imply quicker communication with less distortion. In computer networks, the signal degradation and delay is usually not an issue. Nonetheless, a network with many short paths connecting all nodes will be more efficient in passing data and reconfiguring after a topology change.

Average Path Length is strongly correlated with Closeness throughout the network. As the closeness of all nodes to each other improves (average closeness), the average path length in the network also improves.

Internetwork Topology

In the recent network design book, *Advanced IP Network Design*^[6], the authors define a well-designed topology as the basis of a well-behaved and stable network. They further propose that “three competing goals must be balanced for good network design”:

- Reducing hop count
- Reducing available paths
- Increasing the number of failures the network can withstand

Our social network algorithms can assist in measuring and meeting all three goals.

- Reducing the hop count infers minimizing the average path length throughout the network—maximize the closeness of all nodes to each other.
- Reducing the available paths leads to minimizing the number of geodesics throughout the network.
- Increasing the number of failures a network can withstand focuses on minimizing the centralization of the whole network.

On the following pages we examine various network topologies and evaluate them using social network measures while remembering these three competing goals of network design.

The models we examine do *not* cover hierarchical structures—with Core, Distribution, and Access layers—found in networks of hundreds or thousands of routers. We examine flat, nonhierarchical topologies such as those found in smaller internetworks, area subnetworks, or within core backbones. The topologies we model are the most commonly used—Star, Ring, Full Mesh, and Partial Mesh. We compute the social network measures on each of the topologies and discuss how the various measures help us meet the competing goals discussed above.

Star Topology

The Star topology, shown in Figure 2, has many advantages—but one glaring fault. The advantages include ease of management and configuration for the network administrators. For the Star, the three competing goals delineate as follows:

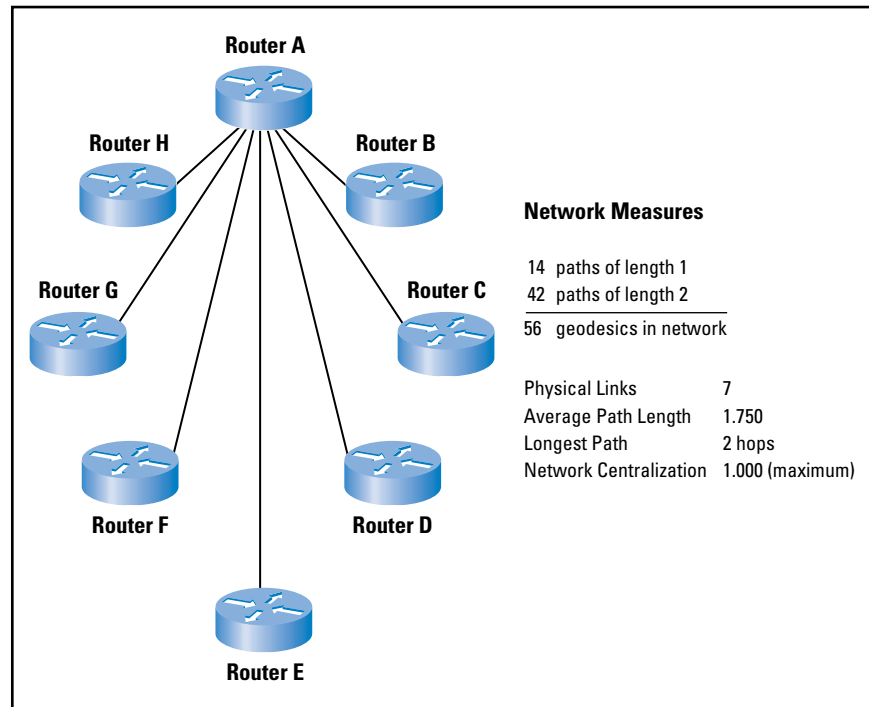
- *Reducing hop count*: The short average path length (1.75) throughout the network meets this goal well. Any router can *reach* any other router in two steps or less.
- *Reducing available paths*: The fact that there are a minimum number of possible available paths (56) to reach all other nodes—will not overload the routing tables, nor cause delays during routing table updates. It takes only seven bidirectional links to create the available paths.

- *Reducing network failures:* The network fails miserably if Router A goes down. Also, any link failure isolates the attached router—there are no multiple paths to reach each router.

Router A is not only a single point of failure—it is also a potential bottleneck—it will likely become overburdened with packet flows and routing updates as more routers are added in the star structure.

Router A receives the top score (1.000) in Activity, Betweenness, and Closeness. As a result, the network is very centralized around Router A from the perspective of all measures.

Figure 2: Routers in Star Topology



Ring Topology

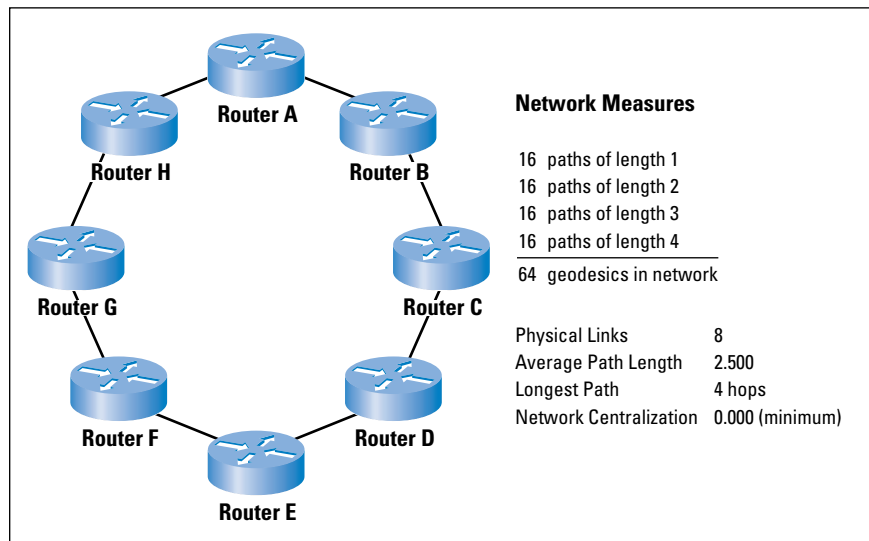
The Ring topology, shown in Figure 3, is an improvement over the Star. It has some of the same advantages, but does not eliminate all of the drawbacks of the Star. The advantages include ease of management and configuration for the network administrators—adding another router is very simple. Unlike the Star topology, the Ring provides some redundancy and, therefore, eliminates the single point of failure—all nodes have an alternate path through which they can be reached. Yet it is still vulnerable to both link and router failures. For the Ring, the three competing goals delineate as follows:

- *Reducing hop count:* The average path length of 2.5 is quite long for a small network of eight nodes. Some routers (that is, A and E) require four steps to reach each other! Many ring physical layers hide this complexity from the IP layers in order to make those hops invisible to routing protocols.

- *Reducing available paths:* This configuration has more geodesics (64) than Star, yet not significantly more to overload the routing tables, nor cause delays during table updates.
- *Reducing network failures:* Even though network centralization is at the minimum (no node is more central than any other), this network reaches failure quickly because of its weak redundancy. The Ring topology can withstand one link failure or one router failure and still keep a contiguous network. Two simultaneous failures can cause unreachable segments because of the lack of redundancy.

Most modern ring technologies such as *Synchronous Optical Network* (SONET) or the Cisco *Dynamic Packet Transport Protocol* (DPT) add a measure of redundancy by running a dual ring that heals itself if a link gets cut. The network “wraps” to avoid the downed line and operates at lower speed. A two-hop path can become a six-hop path if a single link fails. This can cause network congestion if the original dual ring was being used for data in all directions.

Figure 3: Routers in Ring Topology



Full Mesh Topology

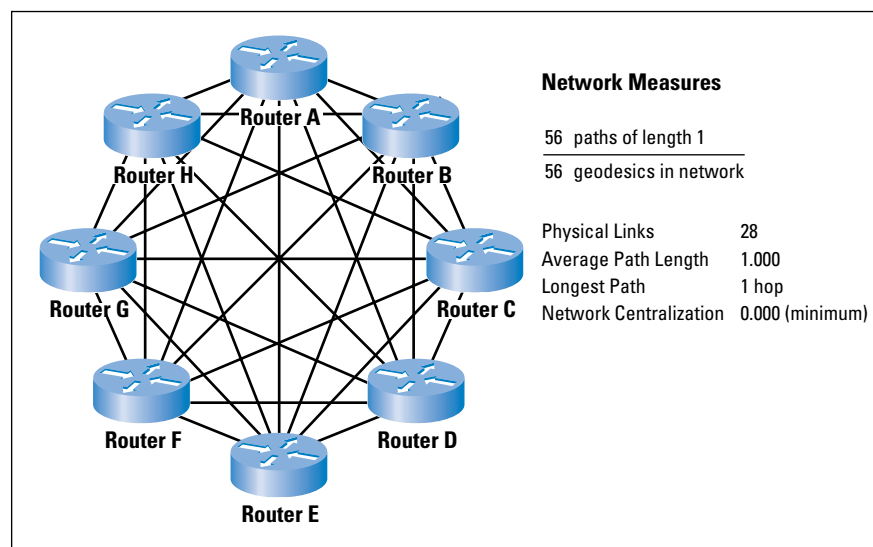
The Full Mesh topology has several big advantages and several faults. The advantages include short path length (one hop) to all other routers and maximum resilience to failure if links or routers start failing. The disadvantages revolve around the complexity created by this topology. For the Full Mesh, the three competing goals delineate as follows:

- *Reducing hop count:* The shortest path length possible is attained for all routes—all nodes can reach each other in one hop.
- *Reducing available paths:* There are a minimum number of possible available paths (56) to reach all other nodes. The routing entries will not overload the routing tables, nor cause delays during routing table updates.

- *Reducing network failures:* The network is not dependent upon any single node (network centralization = 0.000). This configuration represents the most robust topology available—chances are very slim that the number of failures necessary to fragment the network will actually occur within the same time period.

The disadvantages of the Full Mesh topology all focus on one glaring fault—there are too many physical links. If the routers are far apart, the link costs can quickly become prohibitively expensive because adding routers creates a geometrical explosion in links required—soon the routers do not have enough ports to support this topology. Administering the system and keeping an up-to-date topology map becomes more and more complex as routers are added. The network in Figure 4 has 28 two-way links. Double the routers, in a full mesh topology, and the link count increases by a factor greater than 4.

Figure 4: Routers in Full Mesh Topology



Partial Mesh Topology

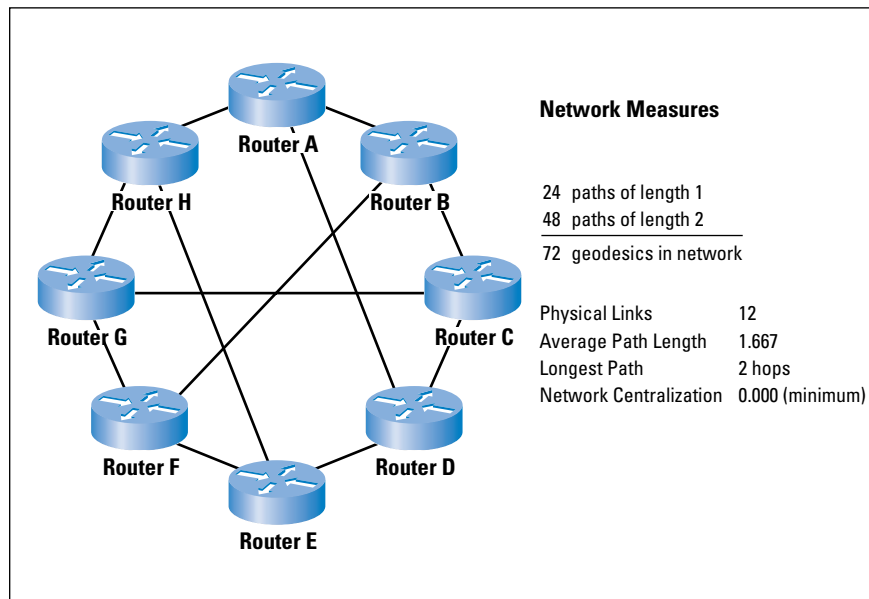
The Partial Mesh topology is quite different. It is the most difficult to build—there is no simple rule to follow (rule for Star: connect everyone to Router A; rule for Full Mesh: connect everyone to everyone). If built incorrectly, the partial mesh layout can have many of the disadvantages of the former topologies without many of the benefits. If built correctly, the opposite is true—more advantages, fewer disadvantages.

Building a successful partial mesh topology is where the interactive use of our social network measures really comes into play. The design below evolved after several iterations. With every iteration the average path length dropped until it appeared to reach a plateau where no further changes lowered the hop count without noticeably increasing the number of physical links. For the Partial Mesh, the three competing goals delineate as follows:

- *Reducing hop count:* The short average path length (1.667) throughout the network meets this goal well. Any router can *reach* any other router in two steps or less. Path length is less than that for the Star and Ring topologies.
- *Reducing available paths:* The number of available paths in the network (72) is the highest among all topologies, though not significantly more than the Ring topology. As the number of nodes in a network increases, this could become a problem—the average path length vs. path count trade-off needs to be closely monitored.
- *Reducing network failures:* Network centralization (0.000) is the same as for the Full Mesh topology—no router, nor link, is more important than any other. As nodes or links are removed from this network, it does not fragment quickly. Chances are slim that the number of failures necessary to fragment the network will actually occur within the same time period. Although we optimized our network centralization for this small “toy” network, we cannot expect this for most real networks. Yet, the goal remains to keep this metric as small as possible.

This topology in Figure 5 was built starting with a Ring topology—a simple architecture. A link was added and the network was remeasured. Was this structure better than the previous? If so, the current structure was kept and another link was added and the network was remeasured. This iterative process was continued until no further improvements happened after several changes. This process does not guarantee an *optimum* solution, yet it quickly converges on a *good* solution—even large networks improve quickly with just a few added links.

Figure 5: Routers in Partial Mesh Topology



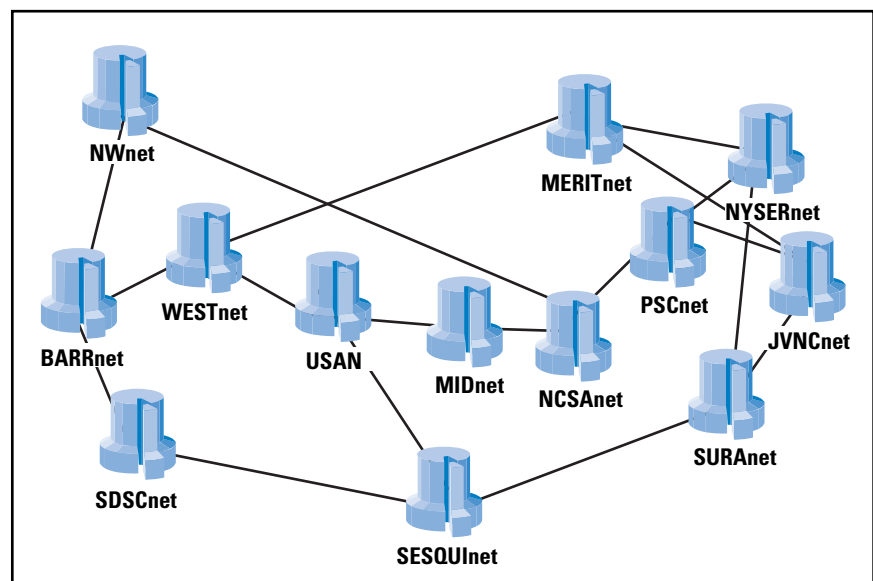
A quirky aspect of networks is that sometimes you can subtract by adding—add a link to a network and reduce the average path length. The opposite also works, sometimes. You can add by subtracting—remove a tie and watch the average hop count grow. Yet, you never know for certain what effect adding or removing a link will have—it is neither a linear nor a local phenomenon. The size and direction of these changes depend upon the existing topology of the network and the location of the added or removed tie. It is key to have a model that allows quick what-if calculations.

Let's experiment with removing random ties—a situation similar to links between routers failing. If we remove the link between Router A and Router H in Figure 5, the number of geodesics in the network increases from 72 to 76, and the average path length increases to 1.815. Yet, removing a different link, G to F, reduces the the number of geodesics in the network from 72 to 66, while the average path length increases only to 1.727. If we are concerned about too many paths in the network, we can remove another link, B to C. This further decreases the number of shortest paths to 60, while reducing physical links to 10. This is very near the 56 paths in the very efficient star topology. Whereas the star is very vulnerable because of its single point of failure, this partial mesh, with the two links removed, is still robust. While the number of geodesics drops, the average path length creeps up slightly to 1.80 with the removal of the second link. Figure 5 has no paths greater than two hops. With the two links (G to F, B to C) removed, we now have 8 geodesics of three hops, while at the same time 12 fewer geodesics to load into routing tables, and two fewer physical links. It is a constant trade-off.

NSFnet Backbone

The NSFnet Backbone network, shown in Figure 6, connected the supercomputing centers in the USA in 1989. It is a partial mesh design that functions as a real-life example to test our social network algorithms.

Figure 6: NSFnet in 1989



We remember our three competing goals for good internetwork design.

- Reducing hop count: average path length in steps/hops
- Reducing available paths: total geodesics in the network
- Increasing the number of failures the network can withstand: network centralization

What happens to these goals as we experience failures in the links or the nodes of the network? Table 1 shows the base metrics for Figure 6 and then shows what happens to the metrics, and our three goals, when five different failures occur.

Table 1: Possible Link and Node Failures

Scenario	Number of Geodesics in the Network	Network Centralization	Longest Path (hops)	Average Path Length (hops)
Original Design (Figure 6)	200	0.062	4	2.370
1) Node failure: NCSA	180	0.208	5	2.689
2) Node failure: MID	180	0.083	4	2.489
3) Node failure: JVNC	148	0.046	4	2.324
4) Link failure: NCSA-PSC	230	0.167	6	2.974
5) Link failure: USAN-MID	212	0.123	5	2.660
6) Link failure: MERIT-JVNC	192	0.069	4	2.458

The most damaging was link failure 4—the link failure between NCSA and PSC. This link is between two of the most central nodes in the network. If the flows between nodes are distributed somewhat evenly, then this link is one of the most traveled in the network.

The least damaging is node failure 3—the node failure at JVNC. In fact, this failure improved most metrics! By removing this node from the network, the number of network paths drops significantly, network centralization decreases, path length decreases slightly, and the longest path is still four hops.

The original NSFnet topology design is very efficient. I tried two different strategies to improve the network. The first strategy involved moving existing links to connect different pairs of routers. No obviously better topology was found by rearranging links among the routers. I was not able to find a better design that reduced both the number of geodesics and the average path length without significantly increasing the number of physical links in the network.

The second strategy is counter-intuitive, yet often networks respond well to this approach. It is the “subtracting by adding” approach described above. By adding new links in the right place in the network, we not only reduce the distance between nodes, we also decrease the number of geodesics in the network.

Because the NSFnet nodes had a maximum limit of three direct neighbors, I started connecting the nodes of Degree = 2. Options 1 through 3 show the various combinations and their effect on the total network. The improvements are minimal, yet each option offers specific strengths.

Option 2 offers more improvements than the others.

- The longest geodesic was reduced to three hops.
- The average path length was reduced throughout the network.
- The number of paths for the routers to remember was reduced slightly.
- Network centralization did not increase enough to noticeably affect the number of failures the network could withstand.

Table 2: Possible Network Improvements

Scenario	Number of Geodesics in the Network	Network Centralization	Longest Path (hops)	Average Path Length (hops)
Original Design (Figure 6)	200	0.062	4	2.370
Option 1 (add link: SDSC–MID)	202	0.071	4	2.287
Option 2 (add link: NW–DSC)	198	0.074	3	2.273
Option 3 (add link: NW–MID)	202	0.050	4	2.356

The improvement in Option 2 (add link: NW–SDSC) was actually implemented in the 1991 version of NSFnet—an excellent example of the “subtracting by adding” network dynamic. Networks are complex systems. How the network responds to change is based on the distribution and pattern of connections throughout the network.

Conclusion

In the real world we may not have the flexibility to experiment with our network model as we have with these examples. There will be more constraints. The information flows in your organization may require that specific pairs of routers have direct links—even if those connections would not be recommended by the algorithms we have been examining. Yet, when we have our “must-have” connections in place, we can experiment with the placement of the remaining connections using these social network metrics to indicate when we are getting close to a robust, yet efficient topology.

Given “initial conditions,” social network methods can model our computer networks and suggest link changes^[7] to form an effective topology that has a short average hop count, not too many paths, and just enough redundancy.

References

- [1] Krebs V., “Visualizing Human Networks,” *Release 1.0*, Esther Dyson’s Monthly Report, February 1996.
- [2] Watts D., Strogatz S., “Collective Dynamics of Small World Networks,” *Nature*, 4 June 1998.
- [3] Freeman L., “Centrality in Social Networks: A Conceptual Clarification,” *Social Networks*, No. 1, 1979.
- [4] Krackhardt D., “Assessing the Political Landscape: Structure, Cognition, and Power in Organizations,” *Administrative Science Quarterly*, No. 35, 1990, page 351.
- [5] Burt, Ronald S., *Structural Holes—The Social Structure of Competition*, ISBN 0674843711, Harvard University Press, 1992.
- [6] Retana, A., Slice, D., White, R., *Advanced IP Network Design*, ISBN 1578700973, Cisco Press, 1999.
- [7] Hagen G., Discussions with fellow network researcher, Guy Hagen, regarding combinatorial algorithms and models for recommending changes to improve the overall topology of a network.

VALDIS E. KREBS leads his own management consulting firm—orgnet.com. He holds an undergraduate degree in Mathematics & Computer Science and a graduate degree in Human Resources. Since 1988 he has applied organizational network analysis to improve knowledge work within and between Fortune 500 firms such as IBM, Lucent, TRW, and supported consulting firms such as Ernst & Young, PricewaterhouseCoopers, and Booz-Allen-Hamilton. In addition to knowledge networks, he has applied these methodologies to mapping, measuring, and molding strategic alliances, communities of interest, emergent structures on the WWW, and internetworks. His work has been referenced in many publications, including the *Wall Street Journal*, *Entrepreneur*, *Training*, *PC Magazine*, *ZDNet*, *Corporate Leadership Council’s Best Practices Reports*, *Knowledge Management*, *Across the Board*, *Business Week*, *HR Executive*, *Personnel Journal*, *FORTUNE*, and Esther Dyson’s influential information industry newsletter, *Release 1.0*. He writes a regular column, “Working in the Connected World,” for the *IHRIM Journal*. His Web site is at: www.orgnet.com and his e-mail is: valdis@orgnet.com

New Frontiers for Research Networks in the 21st Century

by Robert J. Aiken, Cisco Systems, Inc.

A famous philosopher, Yogi Berra, once said, “Prediction is hard. Especially the future.”^[1] In spite of this sage advice, we will still make an attempt at identifying the frontiers for research networks. By first examining and then extrapolating from the evolution and history of past research networks, we may be able to get an idea about the frontiers that face research networks in the future. One of the initial roles of the research network was to act as a testbed for network research on basic network protocols, mostly focusing on network Layers 1 through 4 (that is, the physical, data link, network, transport, and network management layers), but also including basic applications such as file transport and e-mail. During the early phases of the Internet, the commercial sector could not provide the network infrastructure sought by the research and education communities. Consequently, research networks evolved and provided backbone and regional network infrastructures that provided production-quality access to important research and education resources such as supercomputer centers and collaboratories^[2]. Recent developments show that most research networks have moved away from being testbeds for network research and have evolved into production networks serving their research and education communities. It’s time to make the next real evolutionary step with respect to research networks, and that is to shift our research focus toward maximizing the most critical of resources—*people*.

Given the growth and maturity of commercial service providers today, there may no longer be a pressing technical need for governments to continue to support pan-national backbone networks, or possibly even production-like national infrastructures, for Internet-savvy countries. Since commercially available *Virtual Private Networks* (VPNs) can now easily support many of the networked communities that previously required dedicated research networks, government and other supporting organizations can now support their research and education communities by providing the funding for backbone network services much as it does for telephony, office space, and computing capabilities; that is, as part of their research award. However, there may be valid social, political, and long-term economical reasons for continuing the support for such networks. For instance, a nation may decide that in order to ensure its economic survival in the future it wishes to accelerate the deployment and use of Internet technologies among its people, and thus the nation may decide to subsidize national research networks. In addition, it should be noted that VPNs often recreate the “walled” separation of communities, a scenario that was previously accomplished through the hard multiplexing of circuits.

But, in order to make technical advances in the e-economy, governments should now focus on supporting the evolution of intelligent and adaptable edge and access networks. These, in turn, will support the *Ubiquitous Computing* (UC) and persistent presence environments that will soon be an integral part of our future Internet-based economies.

The United States's recently expanded *National Science Foundation* (NSF)^[3] research budget and the *Defense Advanced Projects Agency's* (DARPA's)^[4] prior support of middleware research are good examples of moving in the right direction. The Netherland's Gigaport^[5] project, which incorporates network and application research as well as an advanced technology access and backbone network infrastructure, is a good example of how visionary research networks are evolving.

Just as Internet technologies and network research have matured and evolved, so should the policies concerning the support of research networks. Policies need to be developed to again encourage basic network research and the development of new technologies. In addition, research networks need to encourage and accentuate new network capabilities in edge networks, on campus infrastructures, and in the end systems to support the humans in these new environments. This article focuses mainly on the future of research networks in e-developed nations; but, this is not to diminish the need or importance for e-developed nations to help encourage the same development in network-challenged nations.

Context and Definitions

Before delving into our discussion, we first need to define a few terms. These definitions will not only aid in our discussion, but may also help to highlight the role and function of various types of research networks. The most important terms to define are “network research” and “research network,” both of which often get interchanged during discussions concerning policy, funding, and technology.

In this article, the term “network research” means long-term basic research on network protocols and technologies. The many types of network research can be categorized into three classes. The first category covers research on network transport infrastructure and generally includes research on the *Open System Interconnection* (OSI) Model Layers 1 through 4 (that is, the physical, data link, network, and transport layers) as well as research issues relating to the interconnection and peering of these layers and protocols. We will refer to this class of research as “transport services.”

The second class consists of research covering what can nominally be referred to as “middleware”^[6]. Middleware basically includes many of the services that were originally identified as network Layers 4 through 6. Layer 4 is included because of the need for interfaces to the network layer (sockets, TCP, and so on).

In addition, it nominally includes some components, such as e-mail gateways or directory services, which are normally thought of as being network applications, but which have subcomponents that may also be included in middleware. Given that the definition of middleware is far from an exact science, we shall say that middleware depends on the existence of the network transport services and supports applications.

The third area covers research on the real applications (for example, e-commerce, education, health care, and so on), network interfaces, network applications (for example, e-mail, Web, file transfer, and so on), and the use of networks and middleware in a distributed heterogeneous environment. Applications depend on both the middleware and transport layers. Advanced applications include *Electronic Persistence Presence* (EPP) and UC. EPP, or e-presence, describes a state of a person or application as always being “on the network” in some form or another. The concept of session-based network access will no longer apply. EPP assumes that support for UC and both mobile and nomadic networking exists. UC refers to the pervasive presence of computing and networking capabilities throughout all of our environments; that is, in automobiles, homes, and even on our bodies.

A “research network,” on the other hand, is a production network; that is, one aspiring to the goal of 99.99999-percent “up time” at Layers 1 through 3, which supports various types of domain-specific application research. This application research is most often used to support the sciences and education, but can also be used in support of other areas of academic and economic endeavor. These networks are often referred to as *Research Networks* (RNs) or *Research and Education* (R&E) *Networks*. In this article, we further classify these RNs based on their general customer base. *Institutional Research Networks* (IRNs) support universities, institutes, libraries, data warehouses, and other “campus”-like networks. *National Research Networks* (NRNs)^[7], such as the Netherland’s Gigaport or Germany’s DFN networks, support IRNs or affinity-based networks. *Pan National Research Networks* (PNRNs) interconnect and support NRNs. An example of a couple of current production PNRNs are Dante’s Ten-155 and the NORDUNET^[8] networks. In this article we will also classify the older *National Science Foundation Network’s* (NSFNET’s), *very-high-performance Backbone Network Service* (vBNS), CANARIE’s CA*NET 3^[9], and the Internet 2^[10] Abilene networks as PNRNs because in terms of scale and policy they address the same issues of interconnecting a heterogeneous set of regionally autonomous networks (for example, NSFNET’s regionals and Internet 2’s Gigapops) as do the PNRNs.

A hybrid state of RN also exists. When we introduce one or more advanced technologies into a production system, we basically inject some amount of chaos into the system. The interplay between the new technologies and other existing technologies at various levels of the infrastructure, as well as scaling issues, can cause unanticipated results.

Research quality systems engineering and design is then required to address these anomalies. An example of this phenomenon is the problem encountered with ATM cell discard and its effect on TCP streams and subsequent retransmissions (that is, early packet discard and partial packet discard). The term *Virtual Private Network* (VPN) is used in this article in the classical sense; that is, a network tunneled within another network (for example, IP within IP, ATM virtual circuits [VCs], and so on), and it is not necessarily a security-based network VPN. *Acceptable Use Policy* (AUP) refers to the definition of the type of traffic or use that is allowed on a network infrastructure. *Conditions of Use* (COU) is basically another version of AUP.

Background

During the early phases of the evolution of research networks and the Internet, national research networks were building and managing backbone networks because there was a technical reason to do so. Governments supported these activities, because at the time the commercial sector Internet Service Providers (ISPs) could not do it and the expertise to do so resided within the R&E community. Much of the research or testing of this time still focused on backbone technologies as well as aggregation networks and architectures. Research networks started out by supporting longer-term risky network research and quickly evolved to support shorter-term no-risk production infrastructure.

The research during the *Advanced Research Projects Agency Network* (ARPANET) and early NSFNET phases of the Internet focused on basic infrastructure protocols and technologies. Now commodity services, these services are both easily and cost-effectively available from the commercial sector. We have come a long way since then. Except for a few universities and research centers, the commercial sector now dominates R&D in the backbone technology space. Commercially provided VPNs can now cost-effectively support most of the requirements of the R&E communities. Given the current domination of R&D in backbone technologies by the commercial sectors, as well as the need to address true end-to-end services, it is time that network research and research networks realign their focus onto the research and development of end-system and campus and edge network technologies. Most of the intelligence of the network (for example, *Quality of Service* [QoS], security, content distribution and routing, and so on) will live at the edges, and in some way will be oblivious to the backbone service over which it will operate. In addition, in order for applications to be able to make use of this network, intelligent RNs need to be able to provide the middleware and services that exist between the application and the transport systems. The real future for most RNs is in helping to analyze and identify, not necessarily run and manage, advanced network infrastructures for their R&E communities.

One of the problems faced by the R&E community is how to obtain support from their governments and other supportive organizations (both for-profit and nonprofit). In attempts to support advanced applications and end-user research, organizations and governments may be convinced into supporting RNs, which end up providing commodity services and competing with the commercial sector. One reason that this can occur is that governments often wish to see results very quickly in order to justify their support of the research community; but, by doing so they drive the recipient researchers and research network providers to focus on short-term results and abandon basic long-term research. This pressure from the supporting organizations can also force researchers to compete in a space—that is, transport layers—for which industry may be better suited and adapted in both scale and time. Another issue facing today's research networks is that many of the R&E community, who once would endure downtime and assume some risk in trade for being part of an experimental network, are now demanding full production-quality services from those same R&E networks. Subsequently, the RNs are then being precluded from aggressively pursuing and using really advanced technologies that may pose a risk. And finally, many times research networks, science communities, and researchers claim they are doing network research, when in reality they are not, because they wish to have decent network connectivity, and they assume that this is the only way to get funding and support for good network connectivity with which to support their real research objectives. All of these issues have driven RNs at all levels into difficult positions. RNs need to be able to again take risks if they are to push the envelope in adopting new technology. Likewise, it is also valid to provide production-quality network transport services to support research for middleware, network application (for example collaborative technologies), and R&E application (for example, medical, sciences, education, and so on) research. All of these requirements need to be addressed in the manner most expedient and cost-effective to the government or organization providing the support.

All research carries with it a certain amount of risk. There is theoretical and experimental research. Some research is subject to validation; some is *retrospective*—for example, examining packet traces to verify the existence of nonlinear synchronization—but some is *prospective* and involves reprogramming network resources, and any reprogramming is susceptible to bugs. The amount of risk often depends on the area of research undertaken. The lower down in the network structure that one performs experimental research, the more difficult it is to support this research and still maintain a production-like environment for the other researchers and applications; yet we need to provide support for all levels of experimental research, as described in MORPHNET^[11]. The ideal environment would support applications that could easily migrate from a production network to one prototyping recent network research, and then back again if the experiment fails. Recent advances in optical networking show promise in realizing this goal, but many technical and policy-based challenges are yet to be addressed.

ARPANET and Early NSFNET Phase: 1980s

The ARPANET, one of the many predecessors of today's Internet, was a research project run by researchers as a sandbox where they could develop and test many of the protocols that are now integral components of the Internet. Because this was a research network that supported network research, there were times the network would “go down” and become unavailable. Although that was certainly not the goal, it was a reality when performing experimental network research. This was acceptable to all involved and allowed for the quick “research-to-production” cycle, now associated with the Internet, to develop. The management of the network with respect to policy was handled by the *Internet Activities Board* (IAB), which has since been renamed the *Internet Architecture Board*, and revolved around the actual use of the network as a research vehicle. The research focused mainly on Layers 1 through 4, and application research was secondary and used to demonstrate the underlying technologies.

At the end of the 1980s, the Internet and its associated set of protocols rapidly gained speed in deployment and use among the research community. This started the major shift away from research networks supporting experimental network protocols toward RNs supporting applications via production research networks; for example, the mission agencies' (that is, those agencies whose mission was fairly well focused in a few scientific areas) networks at the *Department of Energy* (DoE) (ESnet^[12]) and NASA (NSInet). At the same time, the NSFNET was still somewhat experimental with the introduction and use of “home-grown” T1 and T3 routers, as well as with pioneering research on peering and aggregation issues associated with the hierarchical NSFNET backbone. It also focused on issues relating to the interconnection of the major agency networks and international networks at the *Federal Internet Exchanges* (FIXes), as well as the policy landscape of interconnecting commercial e-mail (MCIMail) with the Internet. The primary policy justification for supporting these networks (for example ESnet, NSInet, NSFNET) in the late 1980s was to provide access to scarce resources, such as supercomputer centers, although the NSFNET still supported network research, albeit on peering and aggregation.

In addition, the NSFNET was first in pioneering research on network measurement and characterization, leading to today's *Cooperative Association for Internet Data Analysis* (CAIDA) as well as to Surveyor installations on Abilene. As researchers became dependent on the network to support their research, the ability to introduce new and risky technologies into the network became more difficult, as shown by the second-phase T3 router upgrade for the NSFNET when many researchers vehemently complained about any “downtime.”

At this time, there were still no commercial service providers from which to procure IP services to connect the numerous and varied sites of the NSFNET and other research networks. Hence there were still valid technical reasons for NRNs and R&E networks to exist and provide backbone services.

The policy decisions affecting the interconnection of the agency networks at the FIXes, as well as engineering international interconnectivity, were loosely coordinated by an ad hoc group of agency representatives called the *Federal Research Internet Coordinating Committee* (FRICC). The FRICC became the *Federal Networking Council* (FNC) in the early 1990s, and then became the *Large-Scale Network* (LSN) working group by the mid-1990s.

The FNC wisely left the management of the Internet protocols to the IAB, the *Internet Engineering Task Force* (IETF), and the *Internet Engineering Steering Group* (IESG); however, the FNC did not completely relinquish its responsibility, as evidenced by its prominent role in prodding the development of *Classless Interdomain Routing* (CIDR) and originating the work that led to new network protocols (for example, IPv6).

The Next-Generation NSFNET: Early 1990s

During the early 1990s, the Internet evolved and grew larger. It could no longer remain undetected on the government policy radar screen. Many saw the NSFNET and agency networks as competing with commercial *Internet Service Providers* (ISPs). Because of the charters of the agencies of the U.S.-based RNs (for example NSF, DoE, NASA), all traffic crossing their networks had to adhere to their respective AUPs. These AUPs prohibited any “commercial entity-to-commercial entity traffic” to use a U.S. government supported network as transit. In addition, the demand for generic Internet support for all types of research and education communities became much stronger, and at the same time there was growing support among the U.S. Congress and Executive branches to end the U.S. Federal Government support of the U.S. Internet backbone.

In response to these pressures and the responses to a NSF draft “New NSFNET” proposal, the NSF elected to get out of the business of being the Internet backbone within the United States. This policy change was the nexus for the design of the vBNS, *Network Access Points* (NAPs), and *Routing Arbiter* (RA) described in the ABF paper^[13] by early 1992. The vBNS was meant to provide the NSF supercomputer sites a research network that was capable of providing the high-end network services required by the sites for their Metacenter, as well as to provide the capability for their researchers to perform network research because the centers were still the locus for network expertise. The NAPs were designed to enhance the AUP free interconnectivity of both commercial and R&E ISPs and to further evolve the interconnection of the Internet started by the FIXes and the *Commercial Internet eXchange* (CIX).

The research associated with NRNs is already evolving from dealing with mainly IP and transport protocol research to research addressing the routing and peering issues associated with a highly interconnected mesh of networks. Research was an integral part of the NAP and RA design, but it was now focused on peering of networks as opposed to the transport layer protocols themselves. Although this network was not official until 1995, commercial prototype AUP free NAPs (for example, MAE-EAST) immediately sprang up and hastened the transition to a commercial network. The network was transformed from a hierarchical network topology to a decentralized and distributed peer-to-peer model. It no longer existed for the sole purpose of connecting a large aggregation of R&E users to supercomputer centers and other “one-of-a-kind” resources. The NAPs and the “peering” advances associated with the NAPs constituted a very crucial step for the success of applications such as the *World Wide Web* (WWW) and the subsequent commercialization of the Internet because they provided the required seamless interconnected infrastructure. Although some ISPs, for example UUNET and PSInet, were quickly building out their infrastructure at that time, there still existed the need for PNRNs to act as brokers for acquiring and managing end-to-end IP services for their R&E customer base; it would not be much longer, however, before the ISPs had the necessary infrastructure in place to do this themselves.

The Internet 2 Phase: 1996–2000

The transition to the vBNS, NAP, and RA architecture became official early in 1995 and, as a result, the United States university community lost its government-subsidized production backbone. NSF-supported regionals had lost their support years earlier, and many had already transitioned to become commercial service providers, and the NSF “connections” program for tier 2 and lower schools persisted because it was felt (policy wise) that it was still valid to support such activities. The result of this set of affairs led to the creation of the Internet 2. Many of the top research universities in the United States felt that the then-current set of ISPs could not affordably provide adequate end-to-end services and bandwidth for the academic community’s perceived requirements. As a result, the NSF decided to again support production-quality backbone network services for an elite set of research institutions. This was clearly a policy decision by NSF that had support from the U.S. Congress and Executive branches of government, even though in the early 1990s both Congress and the Executive branches were fairly vocal about not supporting such a network.

The initial phase was to expand to the vBNS and connect hundreds of research universities. The vBNS again changed from a research network, connecting a few sites and focusing on network and Metacenter research, back into a production research network. The vBNS is soon eclipsed by the OC-48 Abilene network. Gigapops, which are localized evolutions of NAPs, are used to connect the top R&E institutions to the Internet 2 backbones (that is, vBNS and Abilene).

These backbones were subject to COU as a way to restrict the traffic to that in direct support of R&E, much like the NSFNET was subject to its AUP.

The ISPs who complained so bitterly about unfair competition in the early 1990s no longer cared, because they had more business than they could handle in selling to corporate customers. An ironic spin on this scenario is that the business demands placed on the commercial ISPs by the late 1990s drove them to aggressively adopt new technologies to remain competitive. Not only were they willing to act as testbeds, they paid for that privilege since it gave them a competitive edge. The result is that in a lot of cases regarding the demonstration and testing of backbone-class technologies, the R&E community was time-wise behind the commercial sector. This situation is further aggravated by the fact that many, but not all, backbone network-savvy R&E folks went to work in industry. Another side effect of this transition is the loss of available network monitoring data. The data used by CAIDA, *The National Laboratory for Applied Network Research* (NLNAR), and other network monitoring researchers had been gathered at the FIXes where most traffic used to pass. With the transition to a commercially dominated infrastructure, meaningful data becomes harder to obtain. In addition, as a result of the COU of the Internet 2 network, and the type of applications it supports (for example, trying to set bandwidth speed records), the traffic passing over its networks can no longer be assumed to be representative Internet data, and its value in this regard is diminished.

Another milestone is reached. ISPs have grown or merged so that they are offering both wide- and local-area network services, and anyone can now easily acquire national and international IP and transport services. The deployment and use of VPNs allows the commercial service providers (SPs) to provide and support various acceptable policy networks with differing AUP/COU on the same infrastructure. The technical need for most PNRNs or NRNs to exist to fulfill this function fades away. Researchers should now be able to specify wide-area network support as a line item in their research proposal budgets, just as they do for telephony and computing support. Most governments do not support separate research “Plain Old Telephone Service” (POTS) networks so that researchers can talk with one another. They provide funding in the grants to allow the researchers to acquire this from the commercial sector. However, valid technical reasons for selectively supporting some research networks still exist. A prime example is the CA*Net 3 network in Canada, which has been extremely aggressive in the adoption and use of preproduction optical networking technologies and infrastructure and has been instrumental in advancing our knowledge on this area.

During this evolution of research networks capabilities, network research is also going through its own evolution. DARPA starts focusing its research on optics, wireless, mobility, and network engineering as part of its Next-Generation Internet program. In addition, the research moves up the food chain of network layers. DARPA and DoE start supporting research on middleware. Globus^[14], along with Legion^[15], Condor^[16], and POLDER^[17], are major middleware research efforts that become the main impetus for GRIDs; and although they are focused mainly on seeking the holy grail of distributed computing, many of the middleware services they are developing are of value in a broader research and infrastructure context. The focus of network research and research networks now starts moving away from backbone transport services to research on advanced collaborative, ubiquitous computing, mobile, nomadic, and EPP environments.

The policy management of the Internet now becomes an oxymoron and reflects the completion of the transition of the Internet to a distributed commercial Internet. Many organizations are now vying for a say in how the Internet evolves. Even the IETF is suffering from its own success. It now faces many of the same political challenges the ITU faced, that is, some commercial companies now try to affect the standards process for their own benefit by introducing standards contributions and only later disclosing the fact that they have filed patents on the technology in question. It is now much more difficult to make policy decisions regarding the future of Internet protocols, technologies, and architectures.

Future Frontiers

UC and EPP are the paradigm shifts at the user level that are already drastically altering our concept and understanding of networks. The scale, number, and complexity of networks supporting these new applications will far exceed anything we have experienced or managed in the past. Users will “be on the net” all the time, either as themselves or indirectly through agents and “bots.” They will be mobile and nomadic. There will be “n” multiple instances of a user active on a network at the same time, and not necessarily from the same logical or geographical location. The frontiers associated with this new focus are many times more complex from a systems integration level than any work we have done in the past with backbone networks. This new frontier will provide new technical challenges at the periphery of the network; that is, the intelligent access and campus networks necessary to support these new environments. EPP and UC will drastically affect our research networks and application environments, much as the Web and its protocols drastically changed Internet and traffic patterns in the 1990s.

The frontiers faced by research networks of the future will depend upon many technical and sociopolitical factors on a variety of levels. The sociopolitical frontiers can be divided into two different classes, one for e-developed nations who have already gone through the learning process

of building an Internet-based infrastructure, and another for the e-challenged nations who still face the challenges of building a viable network transport infrastructure. The developed nations need to now grapple with how they can encourage the next evolutionary phase of their Internet-based economies. Because of the fast evolution of technology, the technical need for subsidizing transport-based network infrastructure is no longer the pressing need it was in the 1990s. The future research network will most likely be nothing more than a VPN based on a commercial ISP “cloud” service that interconnects researchers. The *High Energy Physicists* (HEPs) have already proved that life as a VPN-based affinity group overlaid on production network services is a viable solution to providing for their network requirements. The *High-Energy Physics Network* (HEPnet)^[18] is a virtual set of users and network experts using ESnet and other ISP VPN-based network services to support the HEP scientists. Although we still have some technical challenges associated with backbone network technology (for example, optics), there are now only a very small number of institutions and organizations capable of working with industry and making substantial contributions in this area.

The new technical challenges that need to be addressed now include how to build and deploy intelligent edge and campus networks, content delivery and routing, mobile/nomadic/wireless access to the Internet, and the support for both UC and EPP. The latter two require major advancements and will require a whole bevy of middleware that is both network aware and an integral component of an intelligent network infrastructure. This includes, but is not limited to, directories, locators, presence servers, call admission control services, self-configuring services, mobility, media servers, policy servers, bandwidth brokers, intrusion-detection servers, accounting, authentication, and access control. IRNs and RNs can contribute to our knowledge and growth of these new areas by acting as leaders in areas that tend to be more difficult for the commercial sector to address, for instance, the development and deployment of advanced end-to-end services that operate over one or more ISP-provided clouds. Examples include interdomain bandwidth broker services, multi *Public Key Infrastructure* (PKI) trust models, defining multisite policies and schemas for directory-based policy services, and developing scalable naming conventions.

In order for policy makers to make informed decisions on the evolution and support of Internet technologies and architectures, they will need access to a generic mix of real backbone network data. There still exists a dire need at this point for such data. Innovative solutions that respect the privacy and business concerns of all types of ISPs and RNs, while at the same time making available “scrubbed” data, need to be developed. In addition, with the new focus on edge and metro networks, we might be able to shift our monitoring attentions to this area as well in order to better understand traffic demands and patterns on these scales of networks. Network monitoring is only one of the challenges facing us.

As the scale and complexity of networks grows, even at the pico and body area network level, we will need to develop new techniques to support network modeling, simulation, and experimentation. The University of Utah is developing a test facility^[19] comprising a large number of networked processors, the network equivalent of a supercomputer center, to be used experimentally in the design and development of new transport layer protocols.

Summary

“Being on the net” will change our way of doing e-everything, and the evolution of the underlying infrastructure will need to change in order to support this paradigm shift. The intelligence of the network will not only move to the periphery, but even beyond, to the personal digital assistant and body area network. Therefore, it is important that the goals and focus of the research networks also evolve. Leave the R&D associated with backbone networks mainly with the commercial sector because this is their *raison d’être*. The research networks of the future will be mostly VPNs, with a few exceptions, as noted earlier in this article. Research networks need to focus on the new technologies at the periphery as well as the middleware necessary to support the advanced environments that will soon be commonplace. Many research networks will themselves become virtual, for example, HEPnet, providing expertise but not necessarily a network service.

Policy makers must adapt to address not only these substantial technical and architectural changes but also second-order policy issues such as security and privacy and how to ensure that we don’t end up with a bifurcated digital economy of e-savvy and e-challenged communities.

E-developed nations have already been through the technology learning curve of implementing and deploying a transport infrastructure. The e-challenged nations, with respect to network infrastructure, still face these same challenges, and they have the benefit of taking advantage of the knowledge of the nations who have successfully made the transition. In order to speed up the deployment of Internet technologies and infrastructure in the e-challenged nations, it may be best to first create technologically educated people and then to provide them an economic and social environment where they can apply their knowledge and build the infrastructure. E-savvy nations should help by providing the “know-how.” The *North Atlantic Treaty Organization* (NATO) has a joint program with the *Trans-European Research and Education Networking Association* (TERENA) to provide for the instruction of Eastern European nations on the use and deployment of Internet technology (that is, how to configure and manage routers).

In lieu of subsidizing networks in these nations, NATO and TERENA are providing the basic knowledge that these people need to build, manage, and evolve their own networks and infrastructure. This should be the model to consider for e-developing nations. This is not to diminish the challenges of building network infrastructure in some areas where there is no such infrastructure, and perhaps in some of these areas working with other utility infrastructure providers might advance this cause.

Disclaimer

The ideas, comments, and projections proffered in this article are the sole opinions of the author, and in no way represent or reflect official or unofficial positions or opinions on the part of Cisco Systems, Inc. This article is based on my experience designing and managing operational international research networks, as well as being a program manager for network research, during the formative years of the Internet (that is, my tenure as a program manager for the United States Government's National Science Foundation and the Department of Energy), and my recent experience within Cisco working with next-generation Internet projects and managing its University Research Program. Many of the examples that I cite in this work are based on the development and deployment of the U.S.-based Internet and research networks, although the lessons learned in the United States may also be illuminating elsewhere.

Gratitude

I would like to thank my friend and colleague, Dr. Stephen Wolff, of the Office of the CTO, Cisco Systems Inc., for many good suggestions with respect to improving the content and presentation of this article; but, mostly for his good-humored authentication of my history and facts.

References

- [0] This article was presented at the third Global Research Village Conference organized jointly by the Organization for Economic Cooperation and Development (OECD) and the Netherlands in Amsterdam, December 6–8, 2000.
- [1] This is also attributed to the famous Physicist Niels Bohr.
- [2] Wulf, William A. 1988. "The National Collaboratory—A white paper," Appendix A. In "Towards a National Collaboratory," Unpublished report of a National Science Foundation invitational workshop. Rockefeller University, New York, March 17–18, 1989.
- [3] <http://www.nsf.gov/>
- [4] <http://www.darpa.mil/>
- [5] <http://www.gigaport.nl/>
- [6] `Draft-aiken-middleware-reqndef-01.txt`, Internet Draft, Work in Progress, May 1999, <http://www.anl.gov/ECT/Public/research/morphnet.html>

- [7] See <http://www.dante.org/> and <http://www.terena.nl/> for full lists of European research networks.
- [8] <http://www.nordu.net/>
- [9] <http://www.canarie.ca/>
- [10] <http://www.internet2.org/>
- [11] “Architecture of the Multi-Modal Organizational Research and Production Heterogeneous Network (MORPHnet),” Aiken, et al, ANL-97/1 technical report, and 1997 Intelligent Network and Intelligence in Networks Conference.
<http://moat.nlanr.net/Papers/iinren.ps>
- [12] <http://www.es.net/>
- [13] “NSF Implementation Plan for an Interagency Interim NREN,” (aka Architecture for vBNS, NAPs and RAs), Aiken, Braun, and Ford, GA A21174, May 1992.
- [14] <http://www.globus.org/>
- [15] <http://www.cs.virginia.edu/~legion/>
- [16] <http://www.cs.wisc.edu/condor/>
- [17] <http://www.science.uva.nl/projects/polder/>
- [18] <http://www.hep.net/hepnrc.html>
- [19] <http://www.cs.utah.edu/flux/testbed/>

ROBERT J. AIKEN has an MS in Computer Science from Temple University. He is the Manager of the Cisco University Research Program. Prior to joining Cisco, Bob was the network and security research program manager for DoE's HPCC program and Next-Generation Internet (NGI) initiative. He was a program manager at the National Science Foundation (NSF), and with colleagues Peter Ford and Hans-Werner Braun coauthored the conceptual design and architecture of the second-generation National Science Foundation Network (NSFNET) (vBNS, Network Access Points [NAPs], and the Routing Arbiter [RA]), which enabled the commercialization of the then-U.S.-federally supported Internet. Before his NSF tenure, he served as DoE's ESnet program manager and was the creator and manager of the ESnet Network Information and Services group. Prior to his career in networking, Bob was responsible for managing supercomputers and coding their operating systems. His academic experience includes being an Assistant Professor of Computer Science at Hood College in Maryland, an adjunct Professor at California State University, Hayward, and the Manager of Technology Services at Gettysburg College in Pennsylvania. E-mail: raiken@cisco.com

Book Review

Intrusion Detection *Network Intrusion Detection—An Analyst's Handbook*, by Stephen Northcutt, ISBN 0735708681, New Riders Publishers, 1999.

Network security and the ability to detect intrusion attempts has become extremely important in today's networks, regardless of size. I was looking for a book that would get technical on the details in these matters. Laura Chappell, the guru of packet-level information (www.packet-level.com), recommended this book to me. I should have realized what I was getting into at that point. I purchased the book, which was a bit expensive for its size at \$39.99, and eagerly began reading it.

Mr. Northcutt starts out with a good discussion on how Kevin Mitnick conducted his famous attack. The book presents some very good information on a variety of topics, intermixed with personal observations and opinion. This made for an enjoyable read. If you are considering getting an *Intrusion Detection System* (IDS), then this book will provide you with some valuable insight and guidelines to consider from a recognized industry expert in this field. Mr. Northcutt is affiliated with The *System Administration, Networking, and Security* (SANS) Institute (www.sans.org).

Be aware that this book is not for the faint of heart. You will dive into the depths of packets and intrusion detection rather quickly, and never look back. This is both good and bad. I prefer an easy-to-read technical book, but the level of technical knowledge required to make sense of many of the examples is rather extensive. This includes how the many trace examples are presented in rather specialized fashion; in addition, the touted "detailed" explanations varied in usefulness quite a bit.

The book was marketed as a training aid; however, I suspect most readers need to be quite experienced to benefit from it. I admit I had to read many sections more than once in order to grasp the finer points being conveyed. I am confident that many readers have already echoed this sentiment to the author and publisher, since the second edition of this book was published in September 2000 and the page count has doubled, with only a modest price increase. I put it on my Christmas list!

—Tom Thomas, Mentor Technologies Group
tothomas@mentortech.com

Would You Like to Review a Book for IPJ?

We receive numerous books on computer networking from all the major publishers. If you've got a specific book you are interested in reviewing, please contact us and we will make sure a copy is mailed to you. The book is yours to keep if you send us a review. We accept reviews of new titles, as well as some of the "networking classics." Contact us at ipj@cisco.com for more information.

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

New Top-Level Domains

On November 16, 2000 The board of directors of the *Internet Corporation for Assigned Names and Numbers*, (ICANN) announced its selections for registry operators for new top level domains. The applications selected for further negotiation are the following:

.aero	Societe Internationale de Telecommunications Aeronautiques SC, (SITA)
.biz	JVTeam, LLC
.coop	National Cooperative Business Association, (NCBA)
.info	Afilias, LLC
.museum	Museum Domain Management Association, (MDMA)
.name	Global Name Registry, LTD
.pro	RegistryPro, LTD

The ICANN staff will now work through the end of the year to negotiate registry agreements with the applicants selected. The proposed schedule for completion of negotiations is December 31, 2000. The negotiated registry agreements must then be approved by the board of directors. Following that approval, the ICANN board will forward its recommendations to the U.S. Department of Commerce for implementation. For more on the history of ICANN's new TLD application process, please see <http://www.icann.org/tlds/> Multimedia archives of the annual meeting can be reviewed at <http://cyber.law.harvard.edu/icann/1a2000/>

ICANN is a technical coordination body for the Internet. Created in October 1998 by a broad coalition of the Internet's business, technical, academic, and user communities, ICANN is assuming responsibility for a set of technical functions previously performed under U.S. government contract by IANA and other groups. Specifically, ICANN coordinates the assignment of the following identifiers that must be globally unique for the Internet to function: Internet domain names, Internet Protocol address numbers, and protocol parameter and port numbers. In addition, ICANN coordinates the stable operation of the Internet's root server system. As a non-profit, private-sector corporation, ICANN is dedicated to preserving the operational stability of the Internet; to promoting competition; to achieving broad representation of global Internet communities; and to developing policy through private-sector, bottom-up, consensus-based means. ICANN welcomes the participation of any interested Internet user, business, or organization. See <http://www.icann.org>

ISOC Launches Platinum Membership Level

The *Internet Society* (ISOC) is pleased to announce its *Platinum Sponsorship Program*. The Platinum program, which is in addition to and distinct from ISOC's standard organizational membership categories, provides interested organizations with the ability to designate support for specific areas of ISOC's work.

The initial participants, who also helped define the program, included Cisco, IBM, Microsoft, Nortel, RIPE NCC and SoftComca.com. AP-NIC has since joined the list of Platinum sponsors. Platinum level sponsors contribute \$100,000 annually, with non-profit organizations eligible for funding at half that amount.

The Platinum program was initially developed to bolster support for the standards activities of ISOC, specifically ISOC's support of the *Internet Engineering Task Force* (IETF). Recently the program was expanded beyond Standards to include the three remaining areas of ISOC activities: Education & Training, Public Policy, and Member Services. As a result, participants in the Platinum program can now earmark their contribution for any of these four functional areas, or choose to allocate support for multiple areas, should they so desire.

ISOC is dependent upon individual and organizational members for its funding. ISOC believes that allowing contributors to designate where their money will be spent through the Platinum program enhances the Society's ability to undertake activities in these four areas, and, at the same time, provides an attractive support option for many organizations. ISOC will provide a report on the use of funds to each Platinum-Level sponsor at the end of each year. More information on the Platinum-Level Support Program can be found at:

<http://www.isoc.org/isoc/membership/platinum.shtml>

More information on ISOC's standard membership categories is available from: <http://www.isoc.org/orgs/benefits.shtml>

100 Million Internet Hosts

The Internet reached 100,000,000 hosts on 2 November 2000, according to John S. Quarterman, founder of Matrix.Net, a provider of Internet performance, measurement and intelligence. From its humble beginnings of 4 sites in the western United States in December 1969, the Internet has now reached over 150 countries and is nearly pole to pole. "This is an impressive achievement," said Quarterman. "We have been tracking the growth and development of the Internet for this entire decade. If this kind of growth continues, we will hit 1,000,000,000 hosts in 2006." For more information, see <http://www.matrix.net/>

This publication is distributed on an "as-is" basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Technology
WorldCom, USA

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Edward R. Kozel, Member of The Board of Directors
Cisco Systems, Inc., USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

Copyright © 2000 Cisco Systems Inc.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-10/5
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

Bulk Rate Mail U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol Journal

March 2001

Volume 4, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
The BGP Routing Table	2
LAN QoS.....	16
Book Reviews	24
Call for Papers	29
Fragments	30

FROM THE EDITOR

The rapid growth of the Internet has led to numerous changes to the underlying technologies. In the early days, host names and their corresponding IP addresses were kept in a flat text file ("**HOSTS.TXT**"), updated weekly by the Network Information Center at SRI International. In the mid 1980s it became clear that this method of name/address mapping would not scale, and a new distributed lookup mechanism was designed and deployed. This new method, known as the *Domain Name System* (DNS), has proven successful even in the face of millions of Internet hosts.

Another result of Internet growth is the potential for depletion of the IP Version 4 (IPv4) 32-bit address space. In the early 1990s, this became a matter of great focus for the Internet Engineering Task Force (IETF). The "short-term" fix for this problem was to abandon the original concept of A, B and C address classes and introduce *Classless Interdomain Routing* (CIDR), which consumes addresses in a much more efficient manner—that is to say, more slowly. Address consumption has also been slowed by the use of *Network Address Translation* (NAT) and private address space. Predictions for when the Internet will finally run out of IPv4 addresses varies. The long-term solution is to replace IPv4 with IPv6 which uses 128 bits for addressing.

One area of Internet growth that is currently causing some concern among ISPs is the growing size of the routing table that each router participating in the *Border Gateway Protocol* (BGP) must keep in memory. Our first article, by Geoff Huston, is a detailed look at this problem. Geoff takes an historical look at the BGP routing table, and discusses ways to address some of the issues.

In our March 2000 issue, Geoff Huston wrote an article entitled "Quality of Service—Fact or Fiction?" that discussed the prospects for achieving QoS on an Internet-wide scale. In this issue, Bill Stallings looks at QoS in the LAN environment, which is generally easier to control than the Internet as a whole. LAN QoS has been standardized in IEEE 802.1D which is the subject of this article.

We apologize for the delay in getting our online subscription system up and running. It should be available in the very near future. Meanwhile, please continue to use ipj@cisco.com for any subscription questions or to give feedback on anything you read in this journal.

—Ole J. Jacobsen, Editor and Publisher

ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Analyzing the Internet BGP Routing Table

by Geoff Huston, Telstra

The Internet continues along a path of seemingly inexorable growth, at a rate that has, at a minimum, doubled in size each year. How big it needs to be to meet future demands remains an area of somewhat vague speculation. Of more direct interest is the question of whether the basic elements of the Internet can be extended to meet such levels of future demand, whatever they may be. To rephrase this question, are there inherent limitations in the technology of the Internet—or its architecture of deployment—that may impact the continued growth of the Internet to meet ever-expanding levels of demand?

Numerous potential areas can be searched for such limitations, including the capacity of transmission systems, the switching capacity of routers, the continued availability of addresses, and the capability of the routing system to produce a stable view of the overall topology of the network. This article examines the Internet routing system and the longer-term growth trends that are visible within this system.

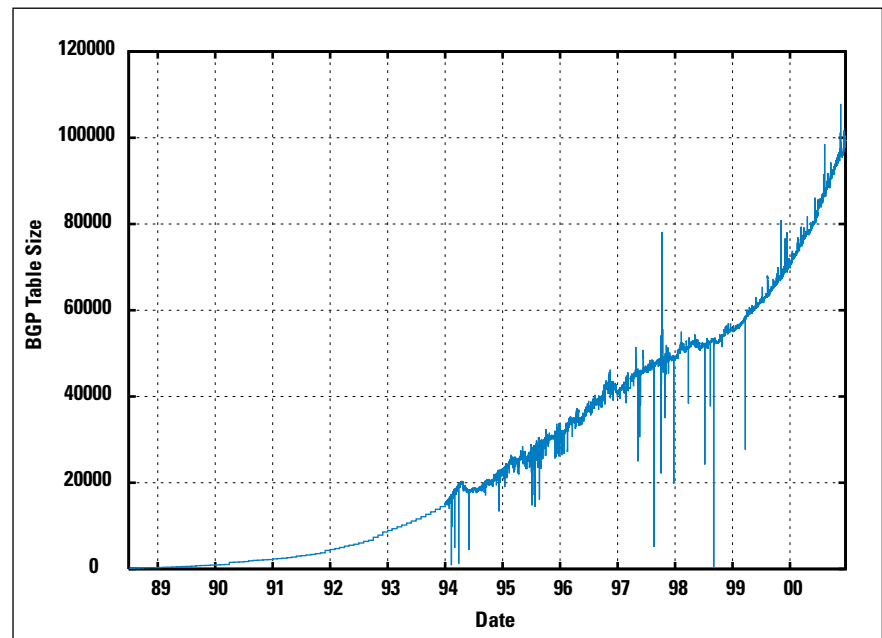
The structure of the global Internet can be likened to a loose coalition of semi-autonomous constituent networks. Each of these networks operates with its own policies, prices, services, and customers. Each network makes independent decisions about where and how to secure the supply of various components that are needed to create the network service. The cement that binds these networks into a cohesive whole is the use of a common address space and a common view of routing. Integrity of routing within each constituent network, or *Autonomous System* (AS), is maintained through the use of an interior routing protocol (or *Interior Gateway Protocol*, or IGP). The collection of these networks is joined into one large routing domain through the use of an inter-network routing protocol (or *Exterior Gateway Protocol*, or EGP).

When the scaling properties of the Internet were studied in the early 1990s, two critical factors identified in the study were, not surprisingly, routing and addressing^[1]. As more devices connect to the Internet, they consume addresses, and the associated function of maintaining reachability information for these addresses implies ever-larger routing tables. The work in studying the limitations of the 32-bit IPv4 address space produced many outcomes, including the specification of IPv6, as well as the refinement of techniques of *Network Address Translation* (NAT) intended to allow some degree of transparent interaction between two networks using different address realms. Growth in the routing system is not directly addressed by these approaches, because the routing space is the cross product of the complexity of the topology of the network, multiplied by the number of autonomous domains of connectivity policy multiplied by the base size of a routing-table entry. When a network advertises a block of addresses into the exterior routing space, this entry is generally carried across the entire exterior routing domain of the

Internet. To measure the characteristics of the global routing table, it is necessary to establish a point in the default-free part of the exterior routing domain and examine the *Border Gateway Protocol* (BGP) routing table that is visible at that point.

Measurements of the size of the routing table were somewhat sporadic in the beginning, and many measurements were taken at approximately monthly intervals from 1988 until 1992 at Merit^[2]. This effort was resumed in 1994 by Erik-Jan Bos at Surfnets in the Netherlands, who commenced measuring the size of the BGP table at hourly intervals at the start of that year. This measurement technique was adopted by the author in 1997, using a measurement point located at the edge of AS 1221 in Australia, again using an hourly interval for the measurement^[6]. The result of these efforts is that we now have a detailed view of the dynamics of the Internet routing-table growth that spans 13 years (Figure 1).

Figure 1: BGP Table Growth 1988–2000



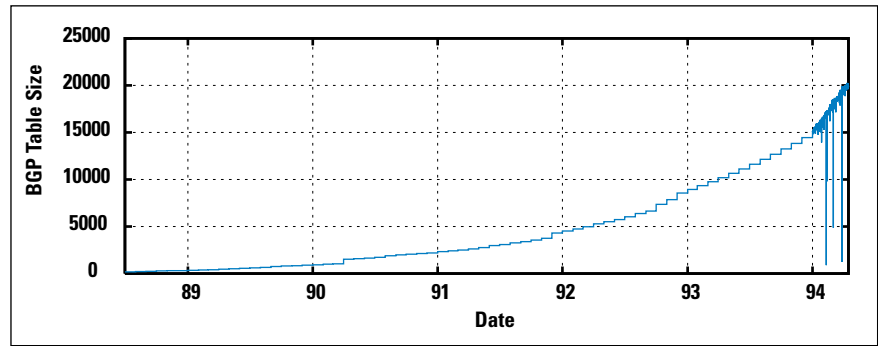
BGP Table Growth

At a gross level, there appear to be four distinct phases of growth visible in this data.

Pre-CIDR Growth

The initial characteristics of the routing-table size from 1988 until April 1994 show definite characteristics of exponential growth (Figure 2). Much of this growth can be attributed to the growth in deployment of the historical Class C address space (/24 address prefixes). Unchecked, this growth would have led to saturation of the BGP routing tables in nondefault routers within a few years. Estimates of the time at which this would have happened vary somewhat, but the overall observation was that the growth rates were exceeding the growth in hardware and software capability of the deployed network at that time.

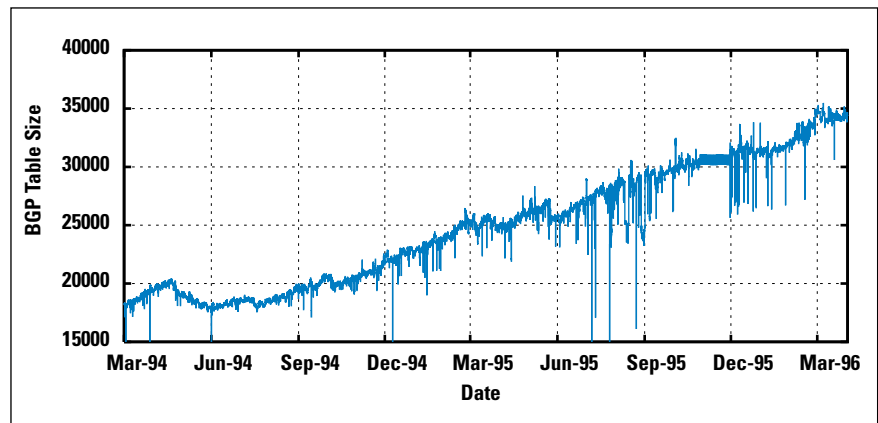
Figure 2: BGP Table Growth 1988–1994



CIDR Deployment

The response from the engineering community was the introduction of routing software that dispensed with the requirement for the Class A, B, and C address delineation, replacing this scheme with a routing system that carried an address prefix and an associated prefix length. A concerted effort was undertaken in 1994 and 1995 to deploy *Classless Interdomain Routing* (CIDR), based on encouraging deployment of the CIDR-capable version of the BGP protocol, BGP4. The effects of this effort are visible in the routing table (Figure 3). Interestingly enough, the efforts of the *Internet Engineering Task Force* (IETF) CIDR Deployment Working Group are visible in the table, with downward movements in the size of the routing table following each IETF meeting.

Figure 3: BGP Table Growth 1994–1995

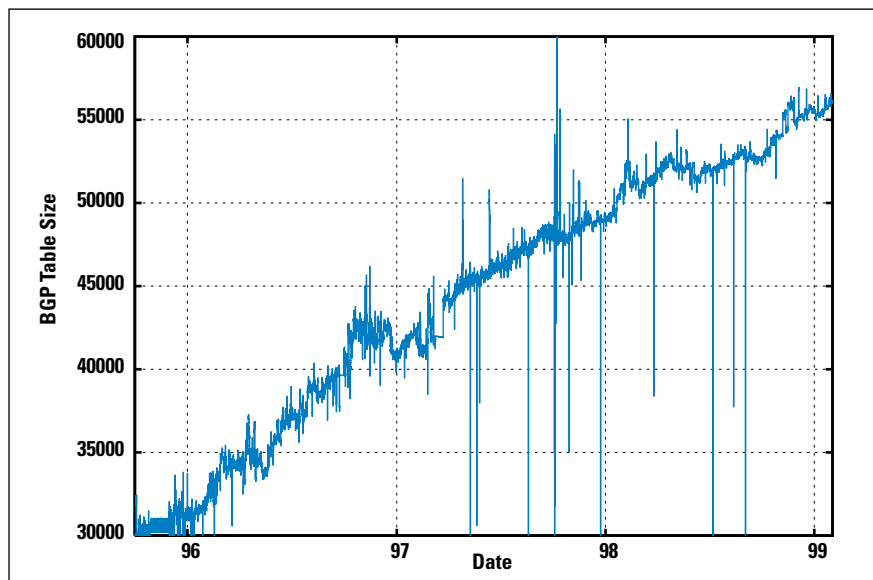


The intention of CIDR was one of supporting an address architecture termed “provider address aggregation,” where a network provider is allocated an address block from the address registry, and announces this entire block into the exterior routing domain. Customers of the provider use a suballocation from this address block, and these smaller routing elements are aggregated by the provider and not directly passed into the exterior routing domain. During 1994, the size of the routing table remained relatively constant at approximately 20,000 entries as the growth in the number of providers announcing address blocks was matched by a corresponding reduction in the number of address announcements as a result of CIDR aggregation.

CIDR Growth

For the next four years until the start of 1998, CIDR proved remarkably effective in damping unconstrained growth in the BGP routing table. While other metrics of Internet size grew exponentially during this period, the BGP table grew at a linear rate, adding about 10,000 entries per year. (Figure 4). Growth in 1997 and 1998 was even lower than this linear rate. Although the reasons behind this are somewhat speculative, it is relevant to note that this period saw intense aggregation within the *Internet Service Provider* (ISP) industry, and in many cases this aggregation was accompanied by large-scale renumbering to fit within provider-based aggregated address blocks. During this period, credit for this trend also must be given to Tony Bates, whose weekly reports of the state of the BGP address table, including listings of further potential for route aggregation, provided considerable incentive to many providers to improve their levels of route aggregation^[4].

Figure 4: BGP Table Growth 1995–1998

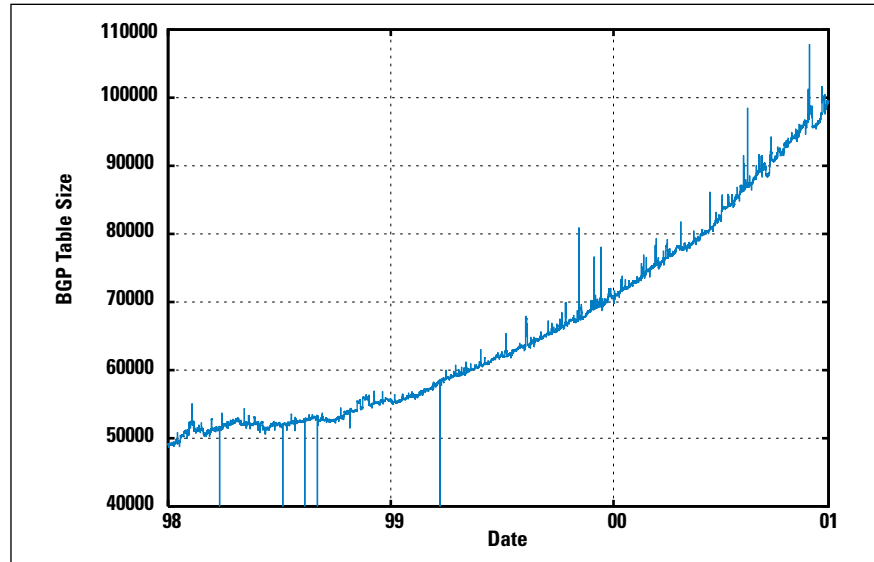


A close examination of the table reveals a greater level of stability in the routing system at this time. The short-term (hourly) variation in the number of announced routes decreased, both as a percentage of the number of announced routes and in absolute terms. One of the other benefits of using large aggregate address blocks is that an instability at the edge of the network is not immediately propagated into the routing core. The instability at the last hop is absorbed at the point at which an aggregate route is used in place of a collection of more specific routes. This, coupled with widespread adoption of BGP route flap damping, has been every effective in reducing the short-term instability in the routing space. It has been observed that whereas the absolute size of the BGP routing table is one factor in scaling, another is the processing load imposed by continually updating the routing table in response to individual route withdrawals and announcements. The encouraging picture from this table is that the levels of such dynamic instability in the network have been reduced considerably by a combination of route flap damping and CIDR.

Current Growth

In late 1998, the trend of growth in the BGP table size changed radically, and the growth for the past two years is again showing all the signs of a reestablishment of exponential growth. It appears that CIDR has been unable to keep pace with the levels of growth of the Internet. (Figure 5). Once again the concern is that this level of growth, if sustained, will outstrip the capability of hardware, or current capability of the BGP routing protocol, or possibly both.

Figure 5: BGP Table Growth 1998–2000



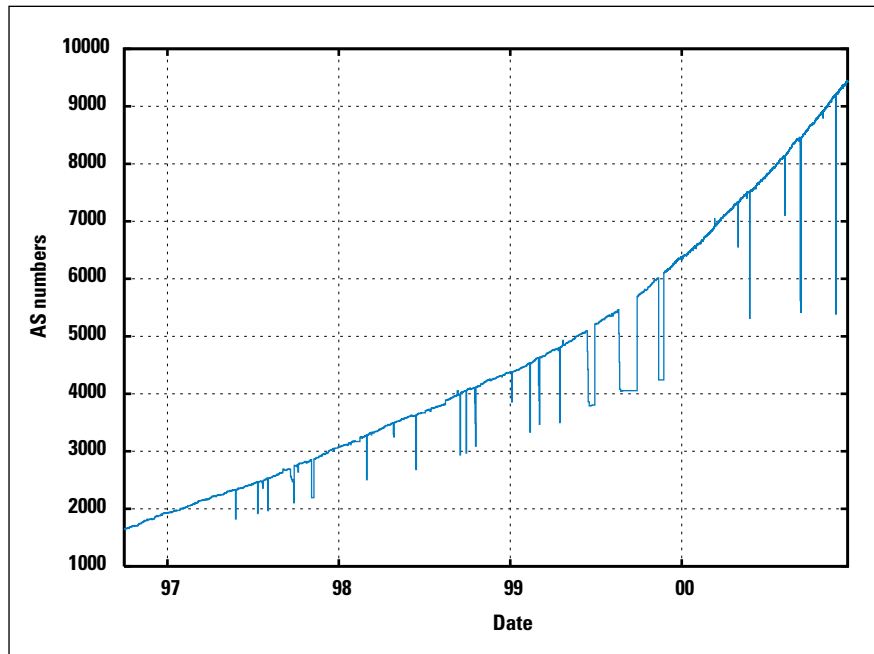
Related Measurements Derived from BGP Table

The level of analysis of the BGP routing table has been extended in an effort to identify the reasons for this resumption of exponential growth. Current analysis includes measuring the number of ASs in the routing system, and the number of distinct AS paths, the range of addresses spanned by the table, and the average span of each routing entry.

AS Number Consumption

Each network that is multihomed within the topology of the Internet and wishes to express a distinct external routing policy must use an AS to associate its advertised addresses with such a policy. In general, each network is associated with a single AS, and the number of ASs in the default-free routing table tracks the number of entities that have unique routing policies. There are some exceptions to this, including large global transit providers with varying regional policies, where multiple ASs are associated with a single network, but such exceptions are relatively uncommon. The trend of AS number deployment over the past four years is also exponential (Figure 6). The growth in the number of ASs can be correlated with the growth in the amount of address space spanned by the BGP routing table. At the end of 2000, the span of advertised addresses is growing at an annual rate of 7 percent, while the number of ASs is growing by 51 percent. Each AS is, on average advertising smaller address ranges. This points to increasingly finer levels of routing detail being announced into the global routing domain, a trend that causes some level of concern.

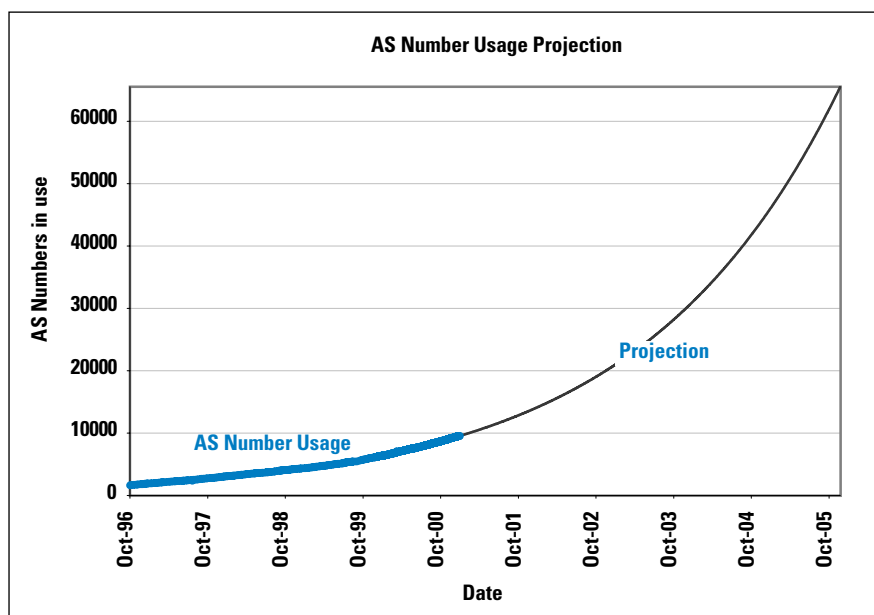
Figure 6: AS Number
Deployment



This is a likely result of an increasingly dense interconnection mesh, where an increasing number of networks are moving from a single-homed connection into multihoming and peering. The spur for this may well be the declining unit costs of communications bearer services.

If this rate of growth continues, the 16-bit AS number set will be exhausted by late 2005 (Figure 7). Work is under way within the IETF to modify the BGP protocol to carry AS numbers in a 32-bit field^[5]. Although the protocol modifications are relatively straightforward, the major responsibility rests with the operations community to devise a transition plan that will allow gradual transition into this larger AS number space.

Figure 7: AS Number
Projections

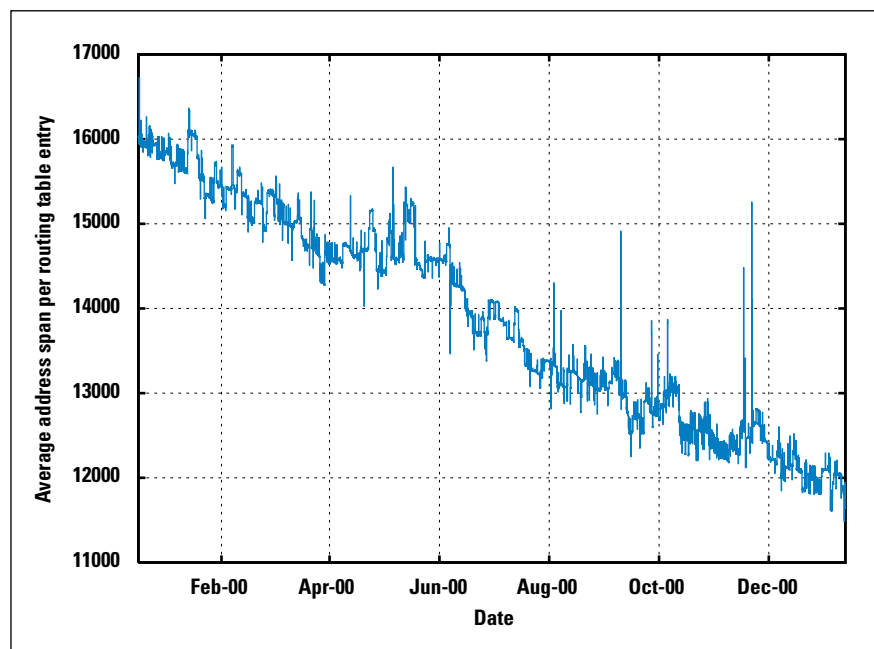


Average Prefix Length of Advertisements

The intent of CIDR aggregation was to support the use of large aggregate address announcements in the BGP routing table. To check whether this is still the case, researchers have tracked the average span of each BGP announcement for the past 12 months. The data indicates a decline in the average span of a BGP advertisement from 16,000 individual addresses in November 1999 to 12,100 in December 2000 (Figure 8). This corresponds to an increase in the average prefix length from /18.03 to /18.44. Separate observations of the average prefix length used to route traffic in operation networks in late 2000 indicate an average length of 18.1^[8]. Again, this trend is cause for concern because it implies the increasing spread of traffic over greater numbers of increasingly finer forwarding-table entries. This, in turn, has implications for the design of high-speed core routers, particularly when extensive use is made of cached forwarding entries within the switching subsystem.

One potential scenario is that the size of the advertisement continues to decrease. With the widespread use of address translation gateway systems, such as NAT, and the continued concern over the finite nature of the IPv4 address pool, this is certainly a highly likely scenario. Projections of the average prefix length of advertisements using current trends in the number of BGP table entries and the total address span advertised in the BGP table indicate a lengthening of the average prefix length of advertisements by 1 bit length every 29 months. This has implications in the lookup algorithms used in routing design, depending on the space/time trade-offs used in the lookup algorithm design. This trend implies that either lookups need to search deeper through the prefix chain to find the necessary forwarding entry, requiring faster memory subsystems to perform each lookup, or the lookup table needs to be both larger and more sparsely populated, increasing the requirements for high-speed memory within the router forwarding subsystem.

Figure 8: Average Span of BGP Advertisement



Prefix Length Distribution

In addition to looking at the average prefix length, the analysis of the BGP table also includes an examination of the number of advertisements of each prefix length.

An extensive effort was introduced in the mid-1990s to move away from extensive use of the Class C space and to encourage providers to advertise larger address blocks. This has been reinforced by the address registries who have used provider allocation blocks of /19 and, more recently, /20. These measures were introduced when there were approximately 20,000 to 30,000 entries in the BGP table. It is interesting to note that five years later, of the 96,000 entries in the routing table, about 53,000 entries have a /24 prefix. In absolute terms, the /24 prefix set is the fastest-growing prefix set in the entire BGP table.

The routing entries of these smaller address blocks also show a much higher level of change on an hourly basis. Although a large number of BGP routing points perform route flap damping, there is still a very high level of announcements and withdrawals of these entries in this particular area of the routing table when viewed using a perspective of route updates per prefix length. Given that the number of these small prefixes is growing rapidly, there is cause for some concern that the total level of BGP flux, in terms of the number of announcements and withdrawals per second, may be increasing, despite the pressures from flap damping. This concern is coupled with the observation that, in terms of BGP stability under scaling pressure, it is not the absolute size of the BGP table that is of prime importance, but the rate of dynamic path recomputations that occur in the wake of announcements and withdrawals. Withdrawals are of particular concern because of the number of transient intermediate states that the BGP distance-vector algorithm explores in processing a withdrawal. Current experimental observations indicate a typical convergence time of about 2 minutes to propagate a route withdrawal across the BGP domain^[7]. An increase in the density of the BGP mesh, coupled with an increase in the rate of such dynamic changes, does have serious implications in maintaining the overall stability of the BGP system as it continues to grow.

The registry allocation policies also have had some impact on the routing-table prefix distribution. The original registry practice was to use a minimum allocation unit of a /19, and the 10,000 prefix entries in the /17 to /19 range are a consequence of this policy decision. More recently, the allocation policy now allows for a minimum allocation unit of a /20 prefix, and the /20 prefix is used by about 4000 entries; in relative terms, this is one of the fastest-growing prefix sets.

The number of entries corresponding to very small address blocks (smaller than a /24), although small in number as a proportion of the total BGP routing table, is the fastest growing in relative terms. The number of /25 through /32 prefixes in the routing table is growing faster, in terms of percentage change, than any other area of the routing table. If prefix length filtering were in widespread use, the practice of announcing a very small address block with a distinct routing policy would have no particular beneficial outcome, because the address block would not be passed throughout the global BGP routing domain and the propagation of the associated policy would be limited in scope. The growth of the number of these small address blocks, and the diversity of AS paths associated with these routing entries, points to a relatively limited use of prefix-length filtering in today's Internet. In the absence of any corrective pressure in the form of widespread adoption of prefix-length filtering, the very rapid growth of global announcement of very small address blocks is likely to continue.

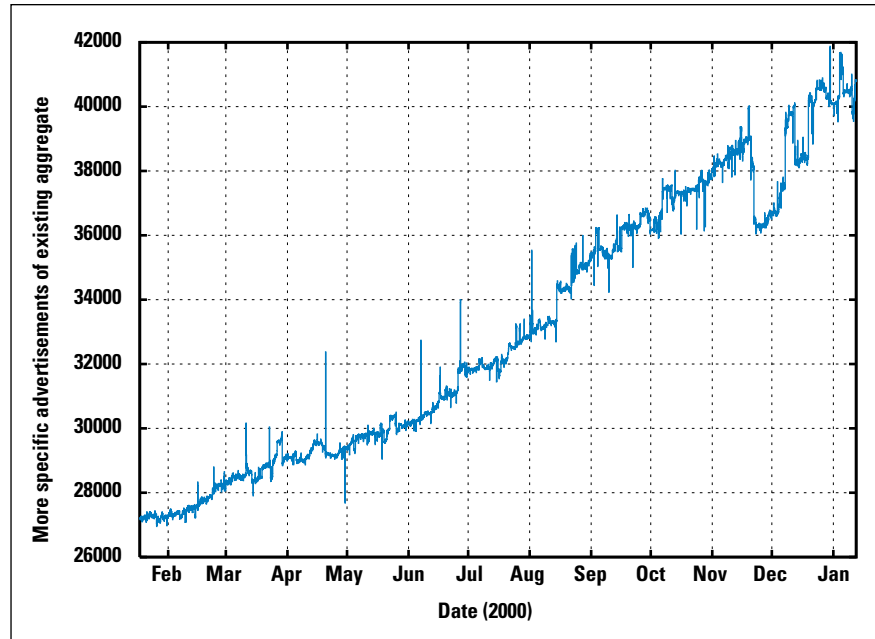
Aggregation and Holes

With the CIDR routing structure, it is possible to advertise a more specific prefix of an existing aggregate. The purpose of this more specific announcement is to punch a "hole" in the policy of the larger aggregate announcement, creating a different policy for the specifically referenced address prefix. Another use of this mechanism is not to promulgate a different connectivity policy, but to perform some rudimentary form of load balancing and mutual backup for multihomed networks. In this model, a network may advertise the same aggregate advertisement along each connection, but then advertise a set of specific advertisements for each connection, altering the specific advertisements such that the load on each connection is approximately balanced. The two forms of holes can be readily discerned in the routing table—while the approach of policy differentiation uses an AS path that is different from the aggregate advertisement, the load balancing and mutual backup configuration uses the same AS path for both the aggregate and the specific advertisements.

Although it is difficult to understand whether the use of such specific advertisements was intended to be an exception to a more general rule or that it was not intended to be within the original intent of CIDR deployment, there appears to be very widespread use of this mechanism within the routing table. Approximately 37,500 advertisements, or 37 percent of the routing table, is being used to punch policy holes in existing aggregate announcements (Figure 9). Of these, the overall majority of about 30,000 routes use distinct AS paths, so that once more we are seeing a consequence of finer levels of granularity of connection policy in a densely interconnected space.

Although long-term data is not available for the relative level of such advertisements as a proportion of the full routing table, the growth level does strongly indicate that policy differentiation at a fine level within existing provider aggregates is a significant driver of overall table growth.

Figure 9: More Specific Advertisements

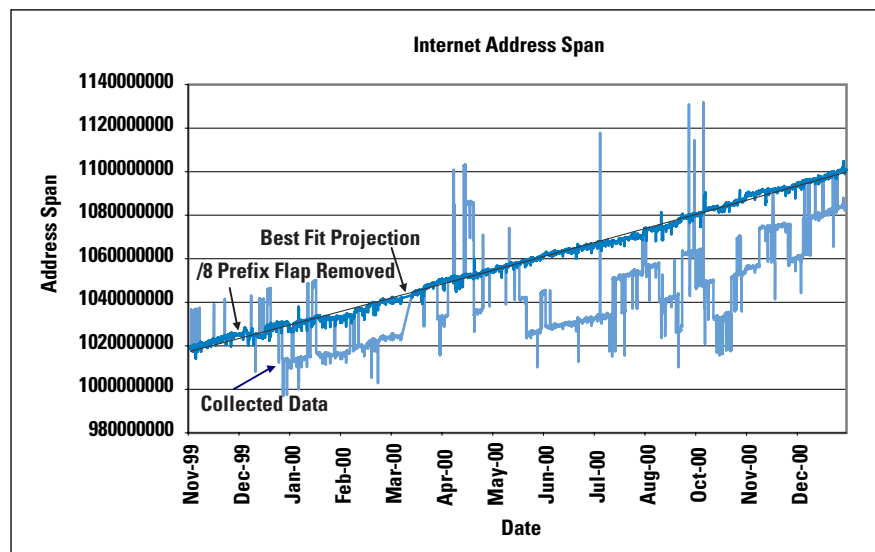


Address Consumption

A decade ago there were two major concerns over scaling of the Internet, and of the two, the consumption of address space was considered to be the more immediate and compelling threat to the continued viability of the network to sustain growth.

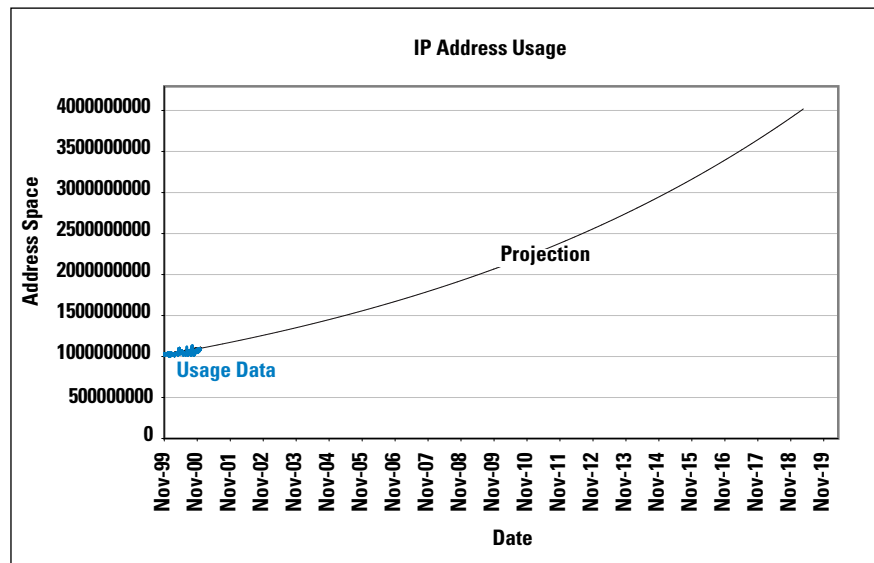
Within the scope of this exercise, it has been possible to track the total span of address space covered by BGP routing advertisements. Over the period from November 1999 until December 2000, the span of address space has grown from 1.02 billion addresses to 1.06 billion. However, numerous /8 prefixes are periodically announced and withdrawn from the BGP table, and if the effects of these prefixes are removed, the final value of addresses spanned by the table is approximately 1.09 billion addresses (Figure 10).

Figure 10: Total Address Space



This is an annual growth rate of a little less than 7 percent, and at that rate of address deployment, the IP Version 4 address space will be able to support another 19 years of such growth (Figure 11). Compared to the 42-percent growth in the number of routing advertisements, it would appear that much of the growth of the Internet in terms of growth in the number of connected devices is occurring behind various forms of NATs. In terms of solving the perceived finite nature of the address space identified just under a decade ago, the Internet appears so far to have embraced the approach of using NATs, irrespective of their various perceived functional shortcomings^[3]. This observation also supports the observed increase of smaller address fragments supporting distinct policies in the BGP table, because such small address blocks encompass arbitrarily large networks located behind one or more NAT gateways.

Figure 11: Address Space Projection



Anomalies

A common space such as the inter-provider domain is not actively managed by any single entity, and various anomalies appear in the routing table from time to time.

One notable event occurred in late 1997, when some large prefixes were deconstructed into a massive set of /24 prefixes and this set was inadvertently passed into the inter-provider BGP domain. The BGP table graphs show a sudden upswing in the number of routing table entries from 50,000 entries to about 78,000 entries. It could have been higher, except that a commonly used routing hardware platform at the time ran into table memory exhaustion at that number of table entries, and further promulgation of additional routing entries ceased. Numerous other anomalies also exist in the table, including the presence of a /31 prefix and several hundred /32 prefixes.

Although many of these anomalies can be attributed to configuration errors of various forms, the underlying observation is that there are no universally used strong filters on what can broadcast into the BGP routing space. Considering the distributed nature of this table and the critical role that it plays in supporting the global Internet, this can be considered a significant current vulnerability. One potential response is to make more use of authentication measures. A validity check could be a precondition to accepting any route advertisement, allowing the receiver of the advertisement a means to check that the origin AS intended to advertise this route. This would create greater resiliency against inadvertent leaks of large sets of advertisements into the broader inter-domain space. It would also improve the resiliency of the BGP domain against some forms of deliberate attack.

Conclusions

There are strong parallels between the BGP routing space and the condition commonly referred to as “The Tragedy Of The Commons.” The BGP routing space is simultaneously everyone’s problem, because it impacts the stability and viability of the entire Internet, and no one’s problem, in that no single entity can be considered to manage this common resource.

In other common resource domains, when the value of the resource is placed under threat because of damaging exploitative practices, the most typical form of corrective action is through the imposition of a consistent set of policies and practices intended to achieve a particular outcome. The vehicle for such an imposition of policies and practices is most commonly that of regulatory fiat. In a globally distributed space such as the BGP table, it is a challenging task to identify the source and authority of such potential regulatory activity.

Multihomed Small Networks

It would appear that one of the major drivers of the recent growth of the BGP table is that of small networks multihoming with numerous peers and numerous upstream providers. In the appropriate environment where numerous networks are in relatively close proximity, using peer relationships can reduce total connectivity costs, as compared to using a single upstream service provider. Equally significantly, multihoming with numerous upstream providers is seen as a means of improving the overall availability of the service. In essence, multihoming is seen as an acceptable substitute for upstream service resiliency.

This has a potential side effect: When multihoming is seen as a preferable substitute for upstream provider resiliency, the upstream provider cannot command a price premium for proving resiliency as an attribute of the provided service, and, therefore, has little incentive to spend the additional money required to engineer resiliency into the network. The actions of the multihomed network clients then become self-fulfilling.

One way to characterize this behavior is that service resiliency in the Internet is becoming the responsibility of the customer, not the service provider.

In such an environment resiliency still exists, but rather than being a function of the bearer or switching subsystem, resiliency is provided through the function of the BGP routing system. The question is not whether this is feasible or desirable in the individual case, but whether the BGP routing system can scale adequately to continue to undertake this role.

A Denser Interconnectivity Mesh

The decreasing unit cost of communications bearers in many part of the Internet is creating a rapidly expanding market in exchange points and other forms of inter-provider peering. The deployment model of a single-homed network with a single upstream provider is rapidly being supplanted by a model of extensive interconnection at the edges of the Internet. The underlying deployment model assumed by CIDR assumed a different structure, more akin to a strict hierarchy of supply providers. The business imperatives driving this denser mesh of interconnection in the Internet are irresistible, and the casualty in this case is the CIDR-induced dampened growth of the BGP routing table.

Traffic Engineering via Routing

Further driving this growth in the routing table is the use of selective advertisement of smaller prefixes along different paths in an effort to undertake traffic engineering within a multihomed environment. Although considerable effort is being undertaken to develop traffic-engineering tools within a single network using *Multiprotocol Label Switching* (MPLS) as the base flow management tool, inter-provider tools to achieve similar outcomes are considerably more complex when using such switching techniques. At this stage, the only tool being used for inter-provider traffic engineering is that of the BGP routing table, further exacerbating the growth and stability pressures being placed on the BGP routing domain.

The effects of CIDR on the growth of the BGP table have been outstanding, not only because of their initial impact in turning exponential growth into a linear growth trend, but also because CIDR was effective for far longer than could have been reasonably expected in hindsight. The current growth factors at play in the BGP table are not easily susceptible to another round of CIDR deployment pressure within the operator community. It may well be time to consider how to manage a BGP routing table that has millions of small entries, rather than the expectation of tens of thousands of larger entries.

We started this journey over ten years ago when considering the scaling properties of addressing and routing. It is perhaps fitting that we tie the two concepts back together again as we consider the future of the BGP inter-provider routing space. The observation that the BGP growth pressures are largely due to an uptake in multihoming and the associated advertisement of discrete connectivity policies by increasingly smaller networks at the edge of the network has a corollary for address allocation policy. In such a ubiquitous environment of multihomed networks, we will also need to review how address blocks are allocated to network providers, because the concept of provider-based address allocation that assumes a relatively strict hierarchical supply structure is becoming less and less relevant in today's Internet.

References

- [1] D. Clark, L. Chapin, V. Cerf, R. Braden, R. Hobby, "Towards the Future Internet Architecture," RFC 1287, December 1991.
- [2] V. Fuller, T. Li, J. Yu, and K. Varadhan, "Supernetting: an Address Assignment and Aggregation Strategy," RFC 1338, June 1992.
- [3] T. Hain, "Architectural Implications of NAT," RFC 2993, November 2000.
- [4] T. Bates, "The CIDR Report," updated weekly at:
<http://www.employees.org/~tbates/cidr-report.html>
- [5] E. Chen, Y. Rekhter, "BGP Support for Four-Octet AS Number Space," work in progress, currently published as an Internet Draft:
[draft-chen-as4bytes-00.txt](#), November 2000.
- [6] "BGP Table Report" updated hourly at
<http://www.telstra.net/ops/bgp>
- [7] C. Labovitz, A. Ahuja, "The Impact of Internet Policy and Topology on Delayed Routing Convergence—Update to This Work," ISMA Winter 2000 Workshop, CAIDA, December 2000.
- [8] Peter Lothberg, personal communication.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. Huston is currently the Chief Scientist in the Internet area for Telstra. He is also a member of the Internet Architecture Board, and is the Secretary of the Internet Society Board of Trustees. He is author of *The ISP Survival Guide*, ISBN 0-471-31499-4, *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, and coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. All three books are published by John Wiley & Sons. E-mail: gih@telstra.net

LAN QoS

by William Stallings

A typical organization's on-premise network configuration has multiple *Local-Area Networks* (LANs) connected by bridges or Layer 2 switches. The LANs may all be of one type (for example, Ethernet) or may be of mixed types (for example Ethernet, Token Ring, wireless). In either case, the issue of *Quality of Service* (QoS) arises.

User Priority and Access Priority

The first attempt to deal with LAN QoS in a standardized fashion appears in the original version of IEEE 802.1D, which is a specification that defines the protocol architecture for bridges and Layer 2 switches, which operate at the *Media Access Control* (MAC) level. IEEE 802.1D deals with the interconnection of LANs with the same MAC protocol and with LANs with different MAC protocols. In addition to passing MAC frames from one LAN to another across the bridge, the bridge is able to pass parameters from software that controls the incoming port to the software that controls the outgoing port. Two of these parameters are *user_priority* and *access_priority*.

The *user_priority* and *access_priority* parameters relate to the problem of how to handle priorities. In the case of IEEE 802.3 (Ethernet) and 802.11 (wireless LAN), priority is not supported. Other 802 LAN types support up to eight levels of priority. The *user_priority* value provided to the MAC-layer entity at the incoming port is derived from the incoming MAC frame; in the case of an incoming frame with no priority value, a value of *unspecified* is used. The *user_priority* value issued to the MAC entity at the outgoing port is to be placed in the outbound MAC frame for LAN types that provide a priority field. The *access_priority* refers to the priority used by a bridge MAC entity to access a LAN for frame transmission. We may not want the *access_priority* to be equal to the *user_priority* for several reasons:

- A frame that must go through a bridge has already suffered more delay than a frame that does not have to go through a bridge; therefore, we may wish to give such a frame a higher access priority than the requested user priority.
- It is important that the bridge not become a bottleneck. Therefore, we may wish to give all frames being transmitted by a bridge a relatively high priority.

The rules for handling priorities can now be summarized. The *user_priority* is determined from the priority field of the incoming frame and placed in the priority field of the outbound frame. Priorities are not used to transmit 802.3 and 802.11 MAC frames, and the frames themselves have no priority field. Therefore, if the outbound frame is 802.3 or 802.11, any incoming priority field (from a frame that has such a field) is ignored. If the incoming frame is 802.3 or 802.11 and the outbound frame requires a priority field, then the priority field in the outbound frame is set to a default *user_priority* value. If both incoming and outbound frames carry a priority field, then the priority field in the outbound MAC frame is set equal to the priority field in the inbound MAC frame.

The *access_priority* is also determined from the priority field of the incoming frame. For incoming 802.3 and 802.11 frames, a *user_priority* of 0 (lowest priority) is assumed. Table 1 shows the access priorities assigned to outgoing MAC frames for each of the LAN types, as a function of incoming user priority value. For 802.3 and 802.11, there is no access priority mechanism and, therefore, a priority of 0 is used. For 802.4 and 802.6, there are eight available access priorities, so the incoming user priority is mapped to the outgoing access priority using equality. IEEE 802.12 permits only two priority levels; half of the possible user priority values are mapped into each of these levels. For the two Token Ring types (802.5 and Fiber Distributed Data Interface [FDDI]), although eight priority levels are available, the highest priority (level 7) is not used in bridge forwarding. The reason for this restriction is that the token-passing protocol reserves priority 7 for its use in transmitting frames needed to manage the token-passing process, such as recovering from a frame loss.

Table 1: Outbound Access Priorities

User Priority	Outbound Access Priority per MAC Method						
	802.3	802.4	802.5	802.6	802.11	802.12	FDDI
0	0	0	0	0	0	0	0
1	0	1	1	1	0	0	1
2	0	2	2	2	0	0	2
3	0	3	3	3	0	0	3
4	0	4	4	4	0	4	4
5	0	5	5	5	0	4	5
6	0	6	6	6	0	4	6
7	0	7	6	7	0	4	6

802.3 = CSMA/CD 802.11 = Wireless LAN
802.4 = Token bus 802.12 = Demand priority (100VG-AnyLAN)
802.5 = Token ring FDDI = Fiber Distributed Data Interface (token ring)
802.6 = DQDB (Distributed Queue, Dual Bus) MAN

Traffic Classes

These rules, summarized in Table 1, are effective in communicating a priority requested by a user and in obtaining access to a LAN in competition with other devices also attempting to transmit on that LAN. However, the rules do not directly provide guidance concerning the relative priority with which frames are to be handled by a bridge. For example, consider a bridge connected to a Token Ring on one side and an Ethernet on the other, and suppose that the bridge receives a large volume of traffic from the Token Ring so that a number of frames are buffered waiting to be transmitted onto the Ethernet. Should the bridge transmit these frames in the order in which they were received, or should the bridge account for the user priority of all waiting frames in determining which frame to transmit next? Consideration of this issue led to the development of a new concept, *traffic class*, which is incorporated in the 1998 version of IEEE 802.1D. This new material is sometimes referred to as 802.1p in the literature. This was the designation when the traffic-class standard was in draft form. In the 802 scheme, a lowercase letter refers to a supplement to an existing standard and an uppercase letter refers to a base standard. Thus 802.1D is a base standard defining bridge operation, and 802.1p is a supplement to the earlier version of 802.1D. With the publication of the 1998 version, the traffic-class supplement was incorporated into 802.1D, and the designation 802.1p is no longer used.

The goal of the traffic-class addition to 802.1D is to enable Layer 2 switches and bridges to support time-critical traffic, such as voice and video, effectively. In the remainder of this article, we begin with an overview of the use of traffic classes in bridges. Next, we examine the mapping of user priorities into traffic classes. Finally, we look at the larger issue of QoS in an internet that includes bridges as well as routers and other Layer 3 switches.

The 1998 version of IEEE 802.1D distinguishes three concepts:

- *User priority*: The user priority is a label carried with the frame that communicates the requested priority to downstream nodes (bridges and end systems). Typically, the user priority is not modified in transit through bridges, unless a mapping is needed for the use of a different number of priority levels by different MAC types. Thus, the user priority has end-to-end significance across bridged LANs.
- *Access priority*: The access priority is used, on LANs that support priority, to compete for access to the shared LAN with frames from other devices (end systems and other bridges) attached to the same LAN. For example, the token-passing discipline in a Token Ring network enables higher-priority frames to gain access to the ring ahead of lower-priority frames when frames from multiple stations are waiting to gain access. When both the incoming and outbound LAN are of the same MAC type, the bridge assigns an access priority equal to the incoming user priority. Otherwise, the bridge must perform a mapping as defined in Table 1.

- *Traffic class*: A bridge can be configured so that multiple queues are used to hold frames waiting to be transmitted on a given outbound port, in which case the traffic class is used to determine the relative priority of the queues. All waiting frames at a higher traffic class are transmitted before any waiting frames of a lower traffic class. As with access priority, traffic class is assigned by the bridge on the basis of incoming user priority.

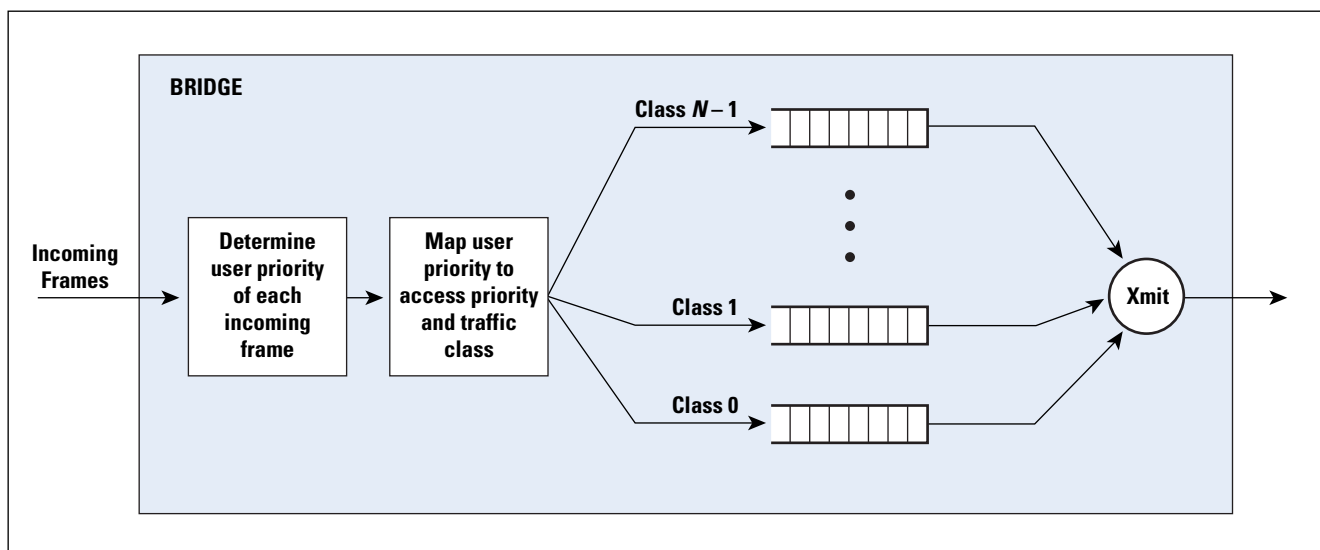
The significance of traffic classes can be seen by recognizing that a frame experiences two types of delay at a bridge:

- *Queuing delay*: The time that a frame waits until it becomes first in line for transmission on the outbound port. This delay is determined by the queuing discipline used by the bridge. The simplest scheme is first-in, first-out (FIFO). Traffic classes permit more sophisticated schemes.
- *Access delay*: The delay that a frame experiences waiting for permission to transmit on the LAN, in competition with frames from other stations attached to the same LAN. This delay is determined by the MAC protocol used (for example Token Ring, Carrier Sense Multiple Access Collision Detect [CSMA/CD]).

The total delay experienced by a frame at a bridge is the sum of its queuing delay and its access delay.

Figure 1 illustrates the mechanism used to support traffic classes at a bridge. A bridge may support up to eight different traffic classes on any outbound port by implementing up to eight distinct queues, or buffers, for that port. A traffic-class value is associated with each queue, ranging from a low of 0 to a high of $N - 1$, where N is the number of traffic classes associated with a given outbound port ($N \leq 8$).

Figure 1: IEEE 802.1 D
Traffic Class Operation



On a given output port with multiple queues, the rules for transmission follow:

1. A frame may be transmitted from a queue only if all queues corresponding to numerically higher values of traffic class are empty. For example, if there is a frame in queue 0, it can be transmitted only if all the other queues at that port are currently empty.
2. Within a given queue, the order of frame transmission must satisfy the following: The order of frames received by this bridge and assigned to this outbound port shall be preserved for:
 - Unicast frames with a given combination of destination address and source address
 - Multicast frames for a given destination address

In practice, a FIFO discipline is typically used. Thus, a strict priority mechanism is used. It follows that during times of congestion, lower-priority frames may be stuck indefinitely at a bridge that devotes its resources to moving out the higher-priority frames.

Mapping of User Priority to Traffic Class

IEEE 802.1D provides guidance on the mapping of user priorities into traffic classes. Table 2 shows the recommended mapping. We can make two comments immediately:

1. The mapping is based on the user priority associated with the frame, which, as was mentioned earlier, has end-to-end significance. However, the 802.3 and 802.11 frame formats do not include a priority field, meaning that this end-to-end information could be lost. To address this issue, the bridge is able to reference the priority field contained in a tag header defined in IEEE 802.1Q, which deals with virtual LANs. The 802.1Q specification defines a tag header of 32 bits that is inserted after the source and destination address fields of the frame header. This tag header includes a 3-bit priority field. Thus, if 802.1Q is in use by Ethernet and wireless LAN sources, a user priority can be defined that stays with the frame from source to destination.
2. Outbound ports associated with MAC methods that support only a single access priority, such as 802.3 and 802.11, can support multiple traffic classes. Recall that the traffic class deals with queuing delay, while the access priority deals with access delay.

To understand the reason for the mappings recommended in Table 2, we need to consider the types of traffic that are associated with each traffic class. IEEE 802.1D provides a list of traffic types, each of which can benefit from simple segregation from the others. In descending importance, these types include:

- Network control (7): Both time critical and safety critical, consisting of traffic needed to maintain and support the network infrastructure, such as routing protocol frames.

- Voice (6): Time critical, characterized by less than 10-ms delay, such as interactive voice.
- Video (5): Time critical, characterized by less than 100-ms delay, such as interactive video.
- Controlled load (4): Non-time-critical but loss sensitive, such as streaming multimedia and business-critical traffic. A typical use is for business applications subject to some form of reservation or admission control, such as capacity reservation per flow.
- Excellent effort (3): Also non-time-critical but loss sensitive, but of lower priority than controlled load. This is a best-effort type of service that an information services organization would deliver to its most important customers.
- Best effort (2): Non-time-critical and loss insensitive. This is LAN traffic handled in the traditional fashion.
- Background (0): Non-time-critical and loss insensitive, but of lower priority than best effort. This type includes bulk transfers and other activities that are permitted on the network but that should not impact the use of the network by other users and applications.

Only seven traffic types are defined in IEEE 802.1D. The standard leaves as spare an eighth type, which could be used for traffic of more importance than background but less importance than best effort. The numbers in parentheses in the preceding list are the traffic-class values corresponding to each traffic type if there are eight queues and hence eight traffic classes available at a given output port.

Table 2: Recommended User Priority to Traffic Class Mapping

		Number of Available Traffic Classes							
		1	2	3	4	5	6	7	8
User Priority	0 (default)	0	0	0	1	1	1	1	2
	1	0	0	0	0	0	0	0	0
	2	0	0	0	0	0	0	0	1
	3	0	0	0	1	1	2	2	3
	4	0	1	1	2	2	3	3	4
	5	0	1	1	2	3	4	4	5
	6	0	1	2	3	4	5	5	6
	7	0	1	2	3	4	5	6	7

We can now address the issue of the mapping between user-priority and traffic-class value. If eight traffic class values are available (eight queues at this output port), the obvious mapping would be equality; that is, a user priority of K would map into traffic class K for $0 \leq K < 7$. This obvious mapping is not desirable because of the treatment of default priorities. For 802.3 and 802.11, which do not use priorities, the de-

fault user priority is 0. For other MAC types, such as 802.5, if the user does not specify a priority, the MAC level assigns a default value of 0. The 802.1D standard points out that using a different default value would result in some confusion and probably a lack of interoperability. However, the logical default traffic type is best effort. The solution proposed by 802.1D is to map a user priority of 0 to traffic-class value 2. When there are eight traffic class values available, then user-priority values 1 and 2 map to traffic-class values 0 (background) and 1 (spare value), respectively.

This solution is reflected in Table 2, which shows the mapping of user priority to traffic class when there are eight available traffic classes. The table also shows the mapping when there are fewer traffic classes. To understand the entries in this table, we need to consider the way in which 802.1D recommends grouping traffic types when fewer than eight queues are configured at a given output port. Table 3 shows this grouping. The first row in the table shows that if there is only one queue, then all traffic classes are carried on that queue. This is obvious. If there are two queues (second row), 802.1D recommends assigning network control, voice, video, and controlled load to the higher-priority queue, and excellent effort, best effort, and background to the lower-priority queue. The reasoning supplied by the standard follows: To support a variety of services in the presence of bursty best-effort traffic, it is necessary to segregate time-critical traffic from other traffic. In addition, further traffic that is to receive superior service and that is operating under admission control also needs to be separated from the uncontrolled traffic. The allocation of traffic types to queues for the remaining rows of the table can be explained similarly.

Table 3: Suggested Traffic Types

		Traffic Types							
Number of Queues	1	BE (EE, BK, VO, CL, VI, NC)							
	2	BE (EE, BK)				VO (CL, VI, NC)			
	3	BE (EE, BK)				CL (VI)		VO (NC)	
	4	BK		BE (EE)		CL (VI)		VO (NC)	
	5	BK		BE (EE)		CL	VI	VO (NC)	
	6	BK		BE	EE	CL	VI	VO (NC)	
	7	BK		BE	EE	CL	VI	VO	NC
	8	BK	—	BE	EE	CL	VI	VO	NC
		1	2	0	3	4	5	6	7
		User Priority							

Note: In each entry, the boldface type is the traffic type that has driven the allocation of types to classes.

BK = Background VI = Video (<100 ms latency and jitter)
 BE = Best Effort VO = Voice (<10 ms latency and jitter)
 EE = Excellent Effort NC = Network Control
 CL = Controlled Load

Internet Traffic Quality of Service

The user-priority and traffic-class concepts enable MAC-level bridges and Layer 2 switches to implement a traffic-handling policy within a bridged collection of LANs that gives preference to certain types of traffic. These concepts are needed because these bridges and switches cannot see “above” the MAC layer and hence cannot recognize or utilize QoS indications in higher layers such as IP. However, it is often the case that traffic from a bridged set of LANs must cross Wide-Area Networks (WANs) that make use of QoS functionality. An example of this is an ATM network, which provides for user-specified QoS. Another example is an IP-based internet, which can provide IP-level QoS. Some means is needed for mapping between traffic classes and QoS for such configurations. This is an evolving area of technology and standardization, but a general picture can be provided.

In the case of IP-based internets, the IP *Type-of-Service* (ToS) field provides a way to label traffic with different QoS demands. The ToS field is preserved along the entire path from source to destination through, potentially, multiple routers. Fortunately, the mapping from traffic class to ToS is straightforward. The ToS field includes a 3-bit Precedence subfield. A router connecting a LAN to an internet can be configured to read the Layer 2 Traffic-Class field and copy that into the ToS Precedence field in one direction, and copy the 3-bit Precedence field into the User Priority field in the other direction.

In the case of an ATM connection, a bridge or Layer 2 switch might be connected to a LAN on one side and an ATM network on the other, using the ATM network to link to other remote LANs. For local LAN traffic arriving at the bridge, the bridge must match the user priority level with the appropriate ATM service class and other ATM parameters. For this purpose, the bridge can consult a mapping table whose settings have been predefined through the policy controls of network management software. An appropriate virtual connection is used to carry the traffic. If the traffic exits the ATM network at another LAN, the bridge on that end can map incoming traffic from each virtual connection into the appropriate traffic class and user priority.

References

A more detailed discussion of bridges, Layer 2 switches, and IEEE 802.1D is contained in [1]. The IEEE 802.1 working group is at <http://grouper.ieee.org/groups/802/1/index.html>.

- [1] Stallings, W., *Local and Metropolitan Area Networks, Sixth Edition*, Prentice Hall, 2000.

WILLIAM STALLINGS is a consultant, lecturer, and author of over a dozen books on data communications and computer networking. He has a PhD in computer science from M.I.T. His latest book is *Local and Metropolitan Area Networks, Sixth Edition* (Prentice Hall, 2000). His home in cyberspace is WilliamStallings.com and he can be reached at ws@shore.net

Book Reviews

E-mail Books *Essential Email Standards: RFCs and Protocols Made Practical* by Pete Loshin, ISBN 0-471-34597-0, John Wiley & Sons, Inc., 2000. www.wiley.com

Internet Email Protocols: A Developer's Guide, by Kevin Johnson, ISBN 0-201-43288-9, Addison-Wesley, 1999. www.awl.com

Deciding when to write a book about an exciting new technology is pretty easy. At first issuance of the standards for it, or emergence of a market for it, out will come the requisite texts. In 1993, when the commercial Internet started to surface, Marshall Rose produced *The Internet Message: Closing The Book With Electronic Mail* [Prentice Hall, 1993]; it's an excellent introduction to the core e-mail services. As the market grew, Rose and David Strom issued a more operations-oriented effort, *Internet Messaging: From Desktop to the Enterprise* [Prentice Hall, 1998]. For anyone serious about e-mail technology and operations, it remains required reading.

But what about straight technology exposition when the standards that have been in use for more than 20 years keep getting modified? In the case of Internet mail, this dilemma has been exacerbated by an extended recent effort to coalesce documentation for the service, compiling and clarifying the contents of many independent *Internet Engineering Task Force* (IETF) documents into two, one for the transfer service and one for the mail object definition. The best time to publish a book on the subject would be at the issuance of the two revisions. Unfortunately, the IETF effort has taken perhaps 3 years longer than expected, and Wiley and Addison-Wesley decided the market needed these books earlier. Hence the authors were faced with a juggling act, referring to original specifications, with appropriate nods to the new—but unstable—drafts.

Comprehensive Introductions

This tactical caveat notwithstanding, Peter Loshin's *Essential Email Standards: RFC and Protocols Made Practical* and Kevin Johnson's *Internet Email Protocols: A Developer's Guide* are credible and reasonably thorough. They introduce the reader to the technical details of Internet mail. Loshin adds detail about the standards culture that produced the specification. Johnson adds a bit of programming detail. No textbook on a technology should be used as the primary reference by someone building products, of course; and these are no exception. These are comprehensive introductions.

With such books, the criteria are simple. I look for helpful overall organization, clear language, and accurate content. These two books qualify. They summarize and restate the basic descriptions of services, data formats, protocol commands, and responses associated with the various standards.

Extra points are assigned when a book comes with commentary that provides some insight into the technical philosophy or operational pragmatics of the technology. Pleasantly, both books have a bit of these extras, too. Such texts typically also have minor technical errors; and these fit that profile, too. Since the reader is not using the book as an implementation reference, the occasional, small errors cause no harm.

Loshin's effort is 330 hardbound pages. Johnson's is about a third longer, softbound. Both books cover the core services of *Submit*, *Simple Mail Transfer Protocol Service Extensions* (ESMTP), the *Post Office Protocol* (POP), the *Internet Mail Access Protocol* (IMAP), RFC 822, and *Multipurpose Internet Mail Extensions* (MIME), that is, posting, relaying, and accessing e-mail, as well as description of the e-mail object. Both also discuss security. Submit is a recent spinoff from SMTP, for local user-relay posting. It began as a clone of ESMTP, but on a different port, and will permit service-to-service relaying functionality to diverge from the local, first-hop posting process. The market treats POP and IMAP as essentially competitive protocols, and both books explain their details adequately. I wish they had made the very simple architectural point that POP does last-hop delivery, to the user's PC-based message store, whereas IMAP is primarily for user access to a message store on a remote system. That is, one is for simply dumping an entire message queue onto the waiting user machine, whereas the other is for ongoing and interaction with portions of message data. On the other hand, an example of Loshin's extra credit is for noting that ISPs are reticent to support IMAP—they have not yet discovered that they could make money being a small business' back-office data store—whereas corporations like IMAP because it is an open standard that permits replacing proprietary workgroup message stores.

E-mail address resolution can be a bit tricky, requiring general understanding of the Domain Name Service and specific cleverness with MX "routing" records. Johnson devotes a useful, but very terse 2+ pages to the topic. Loshin allocates a 8+ pages.

Security

As with every other aspect of Internet standards making, e-mail security is problematic because no IETF-originated security protocol has yet gained wide deployment and use. Oddly continuing the peculiarity of security as a topic, both books are a little off-beat, albeit differently. Johnson provides a relatively extensive introduction to basic security technology, including descriptions of various algorithms, as well as a listing of the types of security attacks that can occur. He also discusses enhancements to the basic e-mail protocols for invoking security mechanisms. Loshin has a more functional systems orientation concerning overall e-mail security architecture. Although Loshin does not usually spend much time on ancient history, for some reason in this chapter he discusses two IETF failures of *Privacy Enhanced Mail* (PEM) and *MIME Object Security Services* (MOSS).

Both discuss *Pretty Good Privacy* (PGP), and PGP is certainly the long-standing popular choice among the technical community. Johnson discusses it in some detail; Loshin's coverage is minimal. *Secure MIME* (S/MIME) has support from major industry software vendors. Loshin treats it equally as tersely as he treats PGP. Johnson barely mentions it.

Standards

Loshin spends the first 50 pages on the Internet standards community, process, and documents. His book also covers Internet News (NNTP) and some work involving standard data for business cards (vCard) and calendaring and scheduling (iCalendar). Besides being interesting topics, these last two were probably included because the Internet Mail Consortium acquired intellectual property rights to the precursor work and highlights the topics on its Web page. Loshin also ends with a chapter about the future, where he adds the topics of instant messaging and message tracking, based on continuing IETF standards work. An included CD-ROM contains a copy of the book, with Web links to cited documents such as RFCs.

Johnson's forays beyond the core services discuss messaging filtering and mailing-list processing, UNIX file issues, and generic, terse descriptions of some programming languages. He also discusses the *Internet Message Support Protocol* (IMSP), the *Application Configuration Access Protocol* (ACAP), and the *Lightweight Directory Access Protocol* (LDAP), protocols for accessing user configuration data. Obviously he intends that the reader take seriously the "Developer's" reference in the book title.

The Differences

Perhaps it is the programmer's orientation that caused Johnson to be so thorough with his discussions. This includes discussion of e-mail protocols that are not standards and not in use. Loshin is far more selective and reflective. And therein lies the easy distinction between the two efforts. Loshin gives an understanding of a portion of application space, providing the basic technical details tidbits of useful insight. Johnson is more mechanical and more detailed; in effect he chooses to be less selective and more detailed in what he dumps on the reader, letting the reader decide what is useful.

—Dave Crocker, *Brandenburg Internet Working*
dcrocker@brandenburg.com

Paging through this book, my first impressions are that it uses very little math and that it is a comprehensive standards-based overview of practical wireless systems. The authors' multidisciplinary tack—systems, networks, and services—is evidenced by their conceptual approach to engineering design issues and their straightforward explanations of implementation issues. The primary concern of the book as a whole is: “How does it all fit together?”

Organization

The authors divide the book into five major units. The first three units covered their topics well and enhanced my understanding of wireless communications. However, the final two units fell short of my expectations. Coverage of the *Wireless Application Protocol* (WAP) and other up-and-coming issues in wireless networking was patchy and unbalanced.

The “PCS Network Management” section provides an overview of the concepts, definitions, and procedures used in current wireless network implementations. Basic roaming concepts including handoff geometry, detection, and queuing schemes are briefly discussed. An understanding of foundational engineering concepts is assumed as the authors provide detailed algorithmic descriptions of hard and soft handoff message flows.

The “IS-41 Mobile Systems” section provides an introductory overview of *Signaling System 7* (SS7) as a supporting protocol for the IS-41 mobile communications protocol. The importance of integration between these two protocols is presented in practical example format. Intersystem handoff and authentication techniques applicable to IS-41 are then discussed. Included in this section is a functional overview of network signaling for *Personal Access Communications* (PACS) networks as related to IS-41. However, a general understanding of the PACS radio system is assumed.

GSM

Global System for Mobile Communication (GSM) systems are the largest focus of this book. A full ten chapters are dedicated to the concepts and applications of this technology. The section appropriately starts with a high-level overview of the GSM system architecture and moves through mobility management and roaming. Here, the authors present several alternative roaming concepts aimed at reducing the cost of roaming service. Additionally, mobile number portability mechanisms and costs are also addressed. Likewise, significant attention is given to the technical aspects of GSM networks and their integration with data networks. Full chapters are dedicated to describing the GSM network signaling software platform (MAP), operations, administration, and management functions, Voice over IP integration, and General Packet Radio Service over GSM.

For the student, *Wireless and Mobile Network Architectures* is a capstone reference that ties together several courses worth of technical information with a practical focus toward real-world applications. For professional IT managers, engineers, and software developers, it is a practical and handy tutorial for getting up-to-speed on second-generation wireless and mobile technologies.

Questions

Each chapter ends with a set of very open-ended and thought-provoking analysis and design questions. Reading the chapter does not necessarily prepare you to do in-depth design; rather, you gain enough knowledge to sketch out a basic approach to solving the problem. It is obvious that many of the problems would require interdisciplinary collaboration to arrive at a tenable solution. Members of such a team would contribute different perspectives based on their particular area of expertise.

Worthwhile Reference

This book assumes that the reader has mastered the basics in the field of mobile communications and is seeking to implement a practical design. Throughout the book are many easy-to-follow algorithmic or flow-chart explanations of various wireless communications processes. However, the information gleaned from these treatments tended to be more about functionality than design. Although a worthwhile reference, this book is by no means “all you need to design and implement a mobile services network.”

—Albert C. Kinney
kinney@ieee.org

Would You Like to Review a Book for IPJ?

We receive numerous books on computer networking from all the major publishers. If you’ve got a specific book you are interested in reviewing, please contact us and we will make sure a copy is mailed to you. The book is yours to keep if you send us a review. We accept reviews of new titles, as well as some of the “networking classics.” Contact us at ipj@cisco.com for more information.

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

ICANN Launches At-Large Membership Study

The *Internet Corporation for Assigned Names and Numbers* (ICANN) recently announced that it was commencing a comprehensive study of the structure of its At Large membership. The study will be conducted by an *At Large Membership Study Committee* that will make recommendations to ICANN's Board of Directors on how individuals can effectively participate in ICANN's policy development, deliberations and actions for technical coordination of the Internet.

Mr. Carl Bildt, the former Prime Minister of Sweden and noted United Nations envoy, will serve as Chair of the nine member Study Committee. An international statesman and information technology advisor, Bildt's current duties include Special Envoy of the Secretary General of the United Nations to the Balkans, Member of Parliament of Sweden, and Advisor and Board Member of several Internet and technology-related corporations.

"The Board's approval of the Study Committee and Carl Bildt's selection as Chair is a demonstration of ICANN's commitment to finding an effective way for the perspectives of individuals in every country to be heard and given due consideration," said Vint Cerf, Chairman of the ICANN Board of Directors. "We are extremely fortunate to have someone with Carl Bildt's international consensus building experience to lead this critical effort."

The Committee, which is chartered to seek input from all interested parties and to work toward a broad consensus on ICANN's At Large membership, will use multiple mechanisms for input, including public forums, mailing lists, and a public website. The Committee will encourage the participation of organizations and individuals worldwide, including the development of independent studies and analyses from across the global Internet's constituencies.

"ICANN's actions affect the whole world's Internet users, and I look forward to the challenging task of forging a consensus on the best method for representing this ever-growing constituency," said Bildt. "This will be an international cooperative effort, and I am counting on the participation of a diversity of Internet stakeholders that have an interest in ICANN to help us deliver a workable solution."

The Board invited Charles Costello and Pindar Wong to serve as the Committee's Vice-Chairs. Costello is director of the Carter Center's Democracy Program, and served as an outside monitor for ICANN's At Large elections held last year. Wong served as an ICANN Director and Vice Chairman of the Board during 1999–2000. He also is an active Internet policy leader in the Asia Pacific Region, and Chairman of Verifi (Hong Kong) Ltd., an Internet infrastructure consultancy. The remaining members of the committee are Pierre Dandjinou, Esther Dyson, Oliver Iteanu, Ching-Yi Lu, Thomas Niles, and Oscar Robles.

ICANN also announced the appointment of Denise Michel as the Committee's Executive Director. Ms. Michel has extensive experience in both private and public sector technology policy development, having served previously on the staff of the U.S. National Science Foundation, the American Electronics Association and the U.S. Department of Commerce. From 1993–1995, she was Sr. Technology Advisor to the Secretary of Commerce, Mr. Ronald Brown.

Following public comment, the Board also adopted a charter for the study to ensure a consistent base of expectations on the scope and details of the study committee's work. ICANN has posted the charter on its website at:

<http://www.icann.org/committees/at-large-study/charter-22jan01.htm>

For more information about the At Large Membership Study Committee, see: <http://www.atlargestudy.org/>

Correction

In the article "The Trouble with NAT," which appeared in our previous issue, a table of private nonroutable IP addresses taken from RFC 1918 was shown. The table contained an error, as pointed out by a couple of our readers. The correct table appears below.

Class	Private Address Range
A	10.0.0.0 ... 10.255.255.255
B	172.16.0.0 ... 172.31.255.255
C	192.168.0.0 ... 192.168.255.255

Upcoming Events

The Internet Society (ISOC) will hold its annual conference INET in Stockholm, Sweden, June 5–8, 2001. For more information, see:

<http://www.isoc.org/inet2001/>

Just before INET, The Internet Corporation for Assigned Names and Numbers (ICANN) will hold its meeting in the same venue. The dates are June 1–4, 2001 and you can find more information at:

<http://www.icann.org/calendar.htm>

The Internet Engineering Task Force (IETF) will next meet in London, England, August 5–10. For more information, see:

<http://www.ietf.org>

This publication is distributed on an "as-is" basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Technology
WorldCom, USA

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Edward R. Kozel, Member of The Board of Directors
Cisco Systems, Inc., USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

Copyright © 2001 Cisco Systems Inc.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-10/5
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSR STD U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol Journal

June 2001

Volume 4, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Mobile IP	2
Goodbye DES, Welcome AES	15
The Middleware Muddle.....	22
Book Review.....	28
Fragments	30

FROM THE EDITOR

A user of a laptop computer “on the road” typically connects to the Internet in one of two ways. The oldest, and most common method, is to dial into an ISP’s network and obtain an IP address using the *Point-to-Point Protocol* (PPP). The other method involves attaching the laptop to a local network (usually via Ethernet) and obtaining an IP address through the *Dynamic Host Configuration Protocol* (DHCP). The “local network” could be anything from the high-speed connection provided in some hotels, to an enterprise network at some corporation or other institution. In all cases, the IP address is fixed for the duration of the network session, and the routing of packets from the laptop back to its “home” network remains a relatively straight-forward task (ignoring NATs, firewalls and other complexities for the moment). Suppose however, the mobile computer is using a wireless connection and traveling between several networks over a short period of time. In this scenario one would still like to maintain network connectivity in a seamless manner. The IETF has been working on Mobile IP to address this problem. Mobile IP is the subject of our first article by Bill Stallings.

The art of cryptography is certainly not new, but its use in computer-communications is a more recent phenomena. The *Data Encryption Standard* (DES) has been widely used since it was standardized in 1977. The strength of a particular encryption scheme depends on the key length and the sophistication of the mathematics involved in transforming the so-called cleartext to the encrypted form. As computers have become more powerful it is now possible to systematically “guess” the 56-bit DES keys in a matter of hours, thus a new encryption standard is needed. This new standard, known as the *Advanced Encryption Standard* (AES), is described by Edgar Danielyan.

Many aspects of computer networking can be described as “controversial,” that is, there are strongly held opinions about a particular technology or its use. In this issue we begin a new series of articles labelled “Opinion,” hoping to bring out some of the different views held by members of the networking community. We hope you will take issue with some of these columns and send us your own opinion piece. We begin the series with an article by Geoff Huston entitled “The Middleware Muddle.” Let us know what you think by sending your comments to ipj@cisco.com

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Mobile IP

by William Stallings

In response to the increasing popularity of palm-top and other mobile computers, Mobile IP was developed to enable computers to maintain Internet connectivity while moving from one Internet attachment point to another. Although Mobile IP can work with wired connections, in which a computer is unplugged from one physical attachment point and plugged into another, it is particularly suited to wireless connections.

The term “mobile” in this context implies that a user is connected to one or more applications across the Internet, that the user’s point of attachment changes dynamically, and that all connections are automatically maintained despite the change. This scenario is in contrast to a user, such as a business traveler, with a portable computer of some sort who arrives at a destination and uses the computer notebook to dial into an *Internet Service Provider* (ISP).

In this latter case, the user’s Internet connection is terminated each time the user moves, and a new connection is initiated when the user dials back in. Each time an Internet connection is established, software in the point of attachment (typically an ISP) is used to obtain a new, temporarily assigned IP address. For each application-level connection (for example, *File Transfer Protocol* [FTP], Web connection), this temporary IP address is used by the user’s correspondent. A better term for this kind of use is “nomadic.”

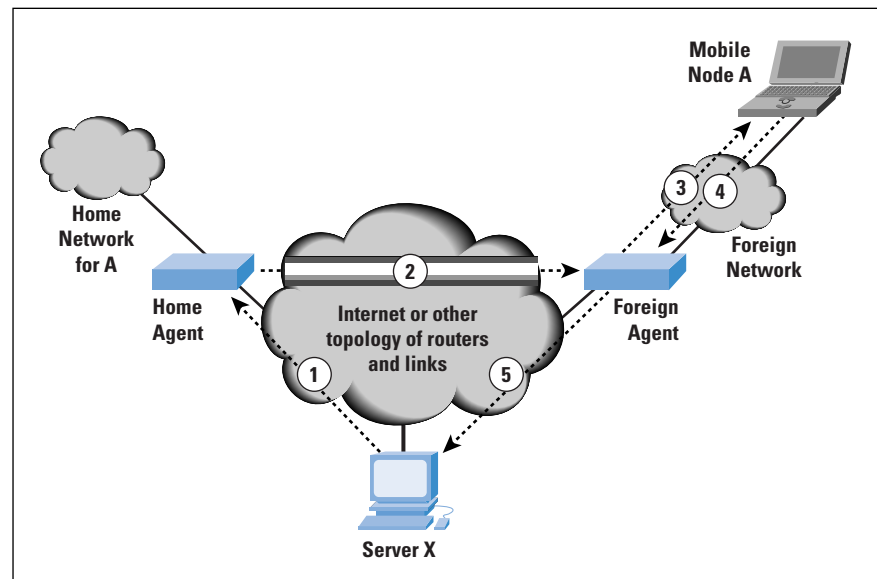
We begin with a general overview of Mobile IP and then look at some of the details.

Operation of Mobile IP

Routers make use of the IP address in an IP datagram to perform routing. In particular, the *network portion* of an IP address is used by routers to move a datagram from the source computer to the network to which the target computer is attached. Then the final router on the path, which is attached to the same network as the target computer, uses the *host portion* of the IP address to deliver the IP datagram to the destination. Further, this IP address is known to the next higher layer in the protocol architecture. In particular, most applications over the Internet are supported by *Transmission Control Protocol* (TCP) connections. When a TCP connection is set up, the TCP entity on each side of the connection knows the IP address of the correspondent host. When a TCP segment is handed down to the IP layer for delivery, TCP provides the IP address. IP creates an IP datagram with that IP address in the IP header and sends the datagram out for routing and delivery. However, with a mobile host, the IP address may change while one or more TCP connections are active.

Figure 1 shows in general terms how Mobile IP deals with the problem of dynamic IP addresses. A mobile node is assigned to a particular network, known as its *home network*. Its IP address on that network, known as its *home address*, is static. When the mobile node moves its attachment point to another network, that is considered a *foreign network* for this host. When the mobile node is reattached, it makes its presence known by registering with a network node, typically a router, on the foreign network known as a *foreign agent*. The mobile node then communicates with a similar agent on the user's home network, known as a *home agent*, giving the home agent the *care-of address* of the mobile node; the care-of address identifies the foreign agent's location. Typically, one or more routers on a network will implement the roles of both home and foreign agents.

Figure 1: Mobile IP Scenario



When IP datagrams are exchanged over a connection between the mobile node (A) and another host (server X in Figure 1), the following operations occur:

1. Server X transmits an IP datagram destined for mobile node A, with A's home address in the IP header. The IP datagram is routed to A's home network.
2. At the home network, the incoming IP datagram is intercepted by the home agent. The home agent encapsulates the entire datagram inside a new IP datagram, which has the A's care-of address in the header, and retransmits the datagram. The use of an outer IP datagram with a different destination IP address is known as *tunneling*.
3. The foreign agent strips off the outer IP header, encapsulates the original IP datagram in a network-level *Protocol Data Unit* (PDU) (for example, a LAN *Logical Link Control* [LLC] frame), and delivers the original datagram to A across the foreign network.

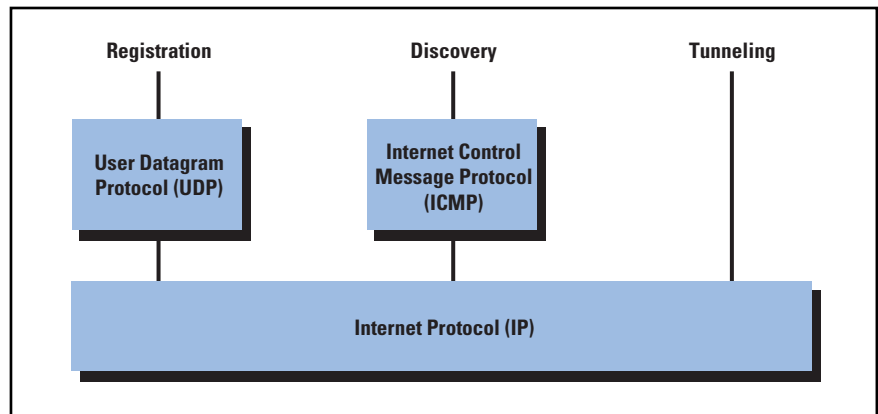
4. When A sends IP traffic to X, it uses X's IP address. In our example, this is a fixed address; that is, X is not a mobile node. Each IP datagram is sent by A to a router on the foreign network for routing to X. Typically, this router is also the foreign agent.
5. The IP datagram from A to X travels directly across the Internet to X, using X's IP address.

To support the operations illustrated in Figure 1, Mobile IP includes three basic capabilities:

- *Discovery*: A mobile node uses a discovery procedure to identify prospective home agents and foreign agents.
- *Registration*: A mobile node uses an authenticated registration procedure to inform its home agent of its care-of address.
- *Tunneling*: Tunneling is used to forward IP datagrams from a home address to a care-of address.

Figure 2 indicates the underlying protocol support for the Mobile IP capability. The registration protocol communicates between an application on the mobile node and an application in the home agent, and hence uses a transport-level protocol. Because registration is a simple request/response transaction, the overhead of the connection-oriented TCP is not required, and, therefore, the *User Datagram Protocol* (UDP) is used as the transport protocol. Discovery makes use of the existing *Internet Control Message Protocol* (ICMP) by adding the appropriate extensions to the ICMP header. ICMP is a connectionless protocol well suited for the discovery operation. Finally, tunneling is performed at the IP level.

Figure 2: Protocol Support for Mobile IP



Discovery

The discovery process in Mobile IP is very similar to the router advertisement process defined in ICMP. Accordingly, agent discovery makes use of ICMP router advertisement messages, with one or more extensions specific to Mobile IP.

The mobile node is responsible for an ongoing discovery process. It must determine if it is attached to its home network, in which case IP datagrams may be received without forwarding, or if it is attached to a foreign network.

Because handoff from one network to another occurs at the physical layer, a transition from the home network to a foreign network can occur at any time without notification to the network layer (that is, the IP layer). Thus, discovery for a mobile node is a continuous process.

For the purpose of discovery, a router or other network node that can act as an agent periodically issues a router advertisement ICMP message with an advertisement extension. The router advertisement portion of the message includes the IP address of the router. The advertisement extension includes additional information about the role of the router as an agent, as discussed subsequently. A mobile node listens for these *agent advertisement messages*. Because a foreign agent could be on the home network of the mobile node (set up to serve visiting mobile nodes), the arrival of an agent advertisement does not necessarily tell the mobile node that it is on a foreign network. The mobile node must compare the network portion of the router IP address with the network portion of its own home address. If these network portions do not match, then the mobile node is on a foreign network.

The *agent advertisement extension* follows the ICMP router advertisement fields and consists of the following fields:

- *Type*: 16, indicates that this is an agent advertisement.
- *Length*: $(6 + 4N)$, where N is the number of care-of addresses advertised.
- *Sequence number*: The count of agent advertisement messages sent since the agent was initialized.
- *Lifetime*: The longest lifetime, in seconds, that this agent is willing to accept a registration request from a mobile node.
- *R*: Registration with this foreign agent is required (or another foreign agent on this network). Even those mobile nodes that have already acquired a care-of address from this foreign agent must reregister.
- *B*: Busy. The foreign agent will not accept registrations from additional mobile nodes.
- *H*: This agent offers services as a home agent on this network.
- *F*: This agent offers services as a foreign agent on this network.
- *M*: This agent can receive tunneled IP datagrams that use minimal encapsulation, explained subsequently.
- *G*: This agent can receive tunneled IP datagrams that use *Generic Routing Encapsulation* (GRE), explained subsequently.
- *Y*: This agent supports the use of Van Jacobson header compression, an algorithm defined in RFC 1144 for compressing fields in the TCP and IP headers.
- *Care-of address*: The care-of address or addresses supported by this agent on this network. There must be at least one such address if the F bit is set. There may be multiple addresses.

There may also be an optional *prefix-length extension* following the advertisement extension. This extension indicates the number of bits in the router address that define the network number. The mobile node uses this information to compare the network portion of its own IP address with the network portion of the router. The fields include the following:

- *Type*: 19, indicates that this is a prefix-length advertisement.
- *Length*: N , where N is the value of the Num Addrs field in the ICMP router advertisement portion of this ICMP message. In other words, this is the number of router addresses listed in this ICMP message.
- *Prefix length*: The number of leading bits that define the network number of the corresponding router address listed in the ICMP router advertisement portion of this message. The number of prefix length fields matches the number of router address fields (N).

Foreign agents are expected to periodically issue agent advertisement messages. If a mobile node needs agent information immediately, it can issue an ICMP router solicitation message. Any agent receiving this message will then issue an agent advertisement.

As was mentioned, a mobile node may move from one network to another because of some handoff mechanism, without the IP level being aware of it. The agent discovery process is intended to enable the agent to detect such a move. The agent may use one of two algorithms for this purpose:

- *Use of Lifetime field*: When a mobile node receives an agent advertisement from a foreign agent that it is currently using or that it is now going to register with, it records the Lifetime field as a timer. If the timer expires before the agent receives another agent advertisement from the agent, then the node assumes that it has lost contact with that agent. If, in the meantime, the mobile node has received an agent advertisement from another agent and that advertisement has not yet expired, the mobile node can register with this new agent. Otherwise, the mobile node should use agent solicitation to find an agent.
- *Use of network prefix*: The mobile node checks whether any newly received agent advertisement is on the same network as the current care-of address of the node. If it is not, the mobile node assumes that it has moved and may register with the agent whose advertisement the mobile node has just received.

The discussion so far has involved the use of a care-of address associated with a foreign agent; that is, the care-of address is an IP address for the foreign agent. This foreign agent will receive datagrams at this care-of address, intended for the mobile node, and then forward them across the foreign network to the mobile node. However, in some cases a mobile node may move to a network that has no foreign agents or on which all foreign agents are busy.

As an alternative, the mobile node may act as its own foreign agent by using a *colocated care-of address*. A colocated care-of address is an IP address obtained by the mobile node that is associated with the current interface to a network of that mobile node.

The means by which a mobile node acquires a colocated address is beyond the scope of Mobile IP. One means is to dynamically acquire a temporary IP address through an Internet service such as *Dynamic Host Configuration Protocol* (DHCP). Another alternative is that the colocated address may be owned by the mobile node as a long-term address for use only while visiting a given foreign network.

Registration

When a mobile node recognizes that it is on a foreign network and has acquired a care-of address, it needs to alert a home agent on its home network and request that the home agent forward its IP traffic. The registration process involves four steps:

1. The mobile node requests the forwarding service by sending a registration request to the foreign agent that the mobile node wants to use.
2. The foreign agent relays this request to the home agent of that mobile node.
3. The home agent either accepts or denies the request and sends a registration reply to the foreign agent.
4. The foreign agent relays this reply to the mobile node.

If the mobile node is using a colocated care-of address, then it registers directly with its home agent, rather than going through a foreign agent.

The registration operation uses two types of messages, carried in UDP segments. The *registration request message* consists of the following fields:

- *Type*: 1, indicates that this is a registration request.
- *S*: Simultaneous bindings. The mobile node is requesting that the home agent retain its prior mobility bindings. When simultaneous bindings are in effect, the home agent will forward multiple copies of the IP datagram, one to each care-of address currently registered for this mobile node. Multiple simultaneous bindings can be useful in wireless handoff situations to improve reliability.
- *B*: Broadcast datagrams. Indicates that the mobile node would like to receive copies of broadcast datagrams that it would have received if it were attached to its home network.
- *D*: Decapsulation by mobile node. The mobile node is using a colocated care-of address and will decapsulate its own tunneled IP datagrams.
- *M*: Indicates that the home agent should use minimal encapsulation, explained subsequently.

- *V*: Indicates that the home agent should use Van Jacobson header compression, an algorithm defined in RFC 1144 for compressing fields in the TCP and IP headers.
- *G*: Indicates that the home agent should use GRE encapsulation, explained subsequently.
- *Lifetime*: The number of seconds before the registration is considered expired. A value of zero is a request for deregistration.
- *Home address*: The home IP address of the mobile node. The home agent can expect to receive IP datagrams with this as a destination address, and must forward those to the care-of address.
- *Home agent*: The IP address of the mobile node home agent. This informs the foreign agent of the address to which this request should be relayed.
- *Care-of address*: The IP address at this end of the tunnel. The home agent should forward IP datagrams that it receives with the mobile node home address to this destination address.
- *Identification*: A 64-bit number generated by the mobile node, used for matching registration requests to registration replies and for security purposes, as explained subsequently.
- *Extensions*: The only extension so far defined is the authentication extension, explained subsequently.

The *registration reply message* consists of the following fields:

- *Type*: 3, indicates that this is a registration reply.
- *Code*: Indicates result of the registration request.
- *Lifetime*: If the code field indicates that the registration was accepted, the number of seconds before the registration is considered expired. A value of zero indicates that the mobile node has been deregistered.
- *Home address*: The home IP address of the mobile node.
- *Home agent*: The IP address of the mobile node home agent.
- *Identification*: A 64-bit number used for matching registration requests to registration replies.

The only extension so far defined is the authentication extension, explained subsequently.

A key concern with the registration procedure is security. Mobile IP is designed to resist two types of attacks:

1. A node may pretend to be a foreign agent and send a registration request to a home agent so as to divert traffic intended for a mobile node to itself.
2. A malicious agent may replay old registration messages, effectively cutting the mobile node from the network.

The technique that is used to protect against such attacks involves the use of message authentication and the proper use of the identification field of the registration request and reply messages.

For purposes of message authentication, each registration request and reply contains an *authentication extension* with the following fields:

- *Type*: Used to designate the type of this authentication extension.
- *Length*: 4 plus the number of bytes in the authenticator.
- *Security parameter index (SPI)*: An index that identifies a security context between a pair of nodes. This security context is configured so that the two nodes share a secret key and parameters relevant to this association (for example, authentication algorithm).
- *Authenticator*: A code used to authenticate the message. The sender inserts this code into the message using a shared secret key. The receiver uses the code to ensure that the message has not been altered or delayed. The authenticator protects the entire registration request or reply message, any extensions prior to this extension, and the type and length fields of this extension.

The default authentication algorithm uses keyed MD5 to produce a 128-bit message digest. For Mobile IP, a “prefix+suffix” mode of operation is used. The MD5 digest is computed over the shared secret key, followed by the protected fields from the registration message, followed by the shared secret key again. Three types of authentication extensions are defined:

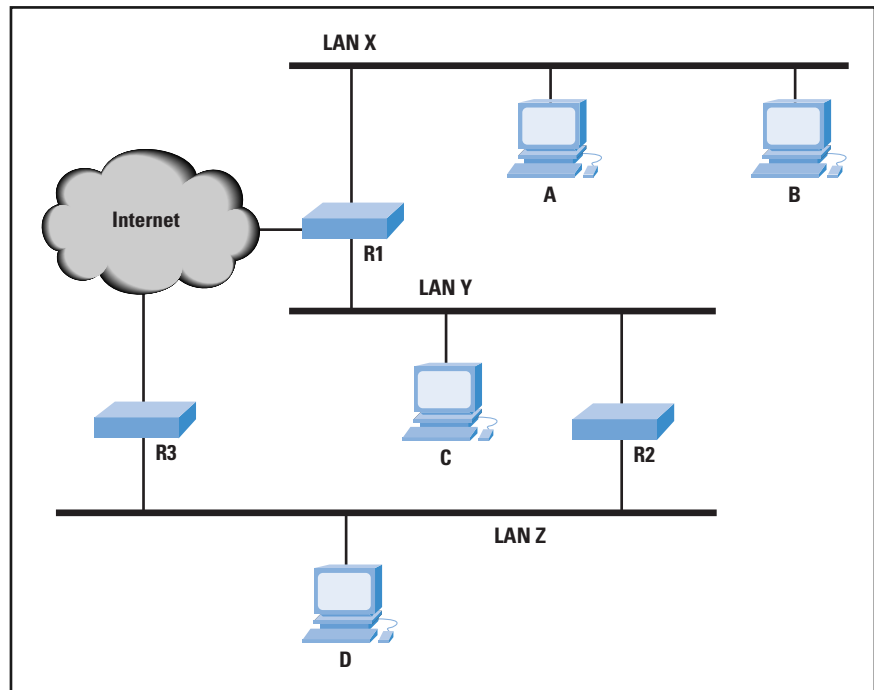
- *Mobile-home*: This extension must be present and provides for authentication of the registration messages between the mobile node and the home agent.
- *Mobile-foreign*: The extension may be present when a security association exists between the mobile node and the foreign agent. The agent will strip this extension off before relaying a request message to the home agent and add this extension to a reply message coming from a home agent.
- *Foreign-home*: The extension may be present when a security association exists between the foreign agent and the home agent.

Note that the authenticator protects the identification field in the request and reply messages. As a result, the identification value can be used to thwart replay types of attacks. As was mentioned, the identification value enables the mobile node to match a reply to a request. Further, if the mobile node and the home agent maintain synchronization so that the home agent can distinguish a reasonable identification value from a suspicious one, then the home agent can reject suspicious messages. One way to do this is to use a timestamp value. As long as the mobile node and home agent have reasonably synchronized values of time, the timestamp will serve the purpose. Alternatively, the mobile node could generate values using a pseudorandom number generator. If the home agent knows the algorithm, then it knows what identification value to expect next.

Tunneling

When a mobile node is registered with a home agent, the home agent must be able to intercept IP datagrams sent to the mobile node home address so that these datagrams can be forwarded via tunneling. The standard does not mandate a specific technique for this purpose but references *Address Resolution Protocol* (ARP) as a possible mechanism. The home agent needs to inform other nodes on the same network (the home network) that IP datagrams with a destination address of the mobile node in question should be delivered (at the link level) to this agent. In effect, the home agent steals the identity of the mobile node in order to capture packets destined for that node that are transmitted across the home network.

Figure 3: A Simple Internetworking Example



For example, suppose that R3 in Figure 3 is acting as the home agent for a mobile node that is attached to a foreign network elsewhere on the Internet. That is, there is a host H whose home network is LAN Z that is now attached to some foreign network. If host D has traffic for H, it will generate an IP datagram with H's home address in the IP destination address field. The IP module in D recognizes that this destination address is on LAN Z and so passes the datagram down to the link layer with instructions to deliver it to a particular *Media Access Control* (MAC)-level address on Z. Prior to this time, R3 has informed the IP layer at D that datagrams destined for that particular address should be sent to R3. Thus, the MAC address of R3 is inserted by D in the destination MAC address field of the outgoing MAC frame. Similarly, if an IP datagram with the mobile node home address arrives at router R2, it recognizes that the destination address is on LAN Z and will attempt to deliver the datagram to a MAC-level address on Z. Again, R2 has previously been informed that the MAC-level address it needs corresponds to R3.

For traffic that is routed across the Internet and arrives at R3 from the Internet, R3 must simply recognize that for this destination address, the datagram is to be captured and forwarded.

To forward an IP datagram to a care-of address, the home agent puts the entire IP datagram into an outer IP datagram. This is a form of encapsulation, just as placing an IP header in front of a TCP segment encapsulates the TCP segment in an IP datagram. Three options for encapsulation are allowed for Mobile IP and we will review the first two of the following options:

- *IP-within-IP encapsulation*: This is the simplest approach, defined in RFC 2003.
- *Minimal encapsulation*: This approach involves fewer fields, defined in RFC 2004.
- *Generic routing encapsulation (GRE)*: This is a generic encapsulation procedure, defined in RFC 1701, that was developed prior to the development of Mobile IP.

In the IP-within-IP encapsulation approach, the entire IP datagram becomes the payload in a new IP datagram (Figure 4a). The inner, original IP header is unchanged except to decrement *Time To Live* (TTL) by 1. The outer header is a full IP header. Two fields (indicated as unshaded in the figure) are copied from the inner header. The version number is 4, the protocol identifier for IPv4, and the type of service requested for the outer IP datagram is the same as that requested for the inner IP datagram.

Figure 4a: Mobile IP Encapsulation

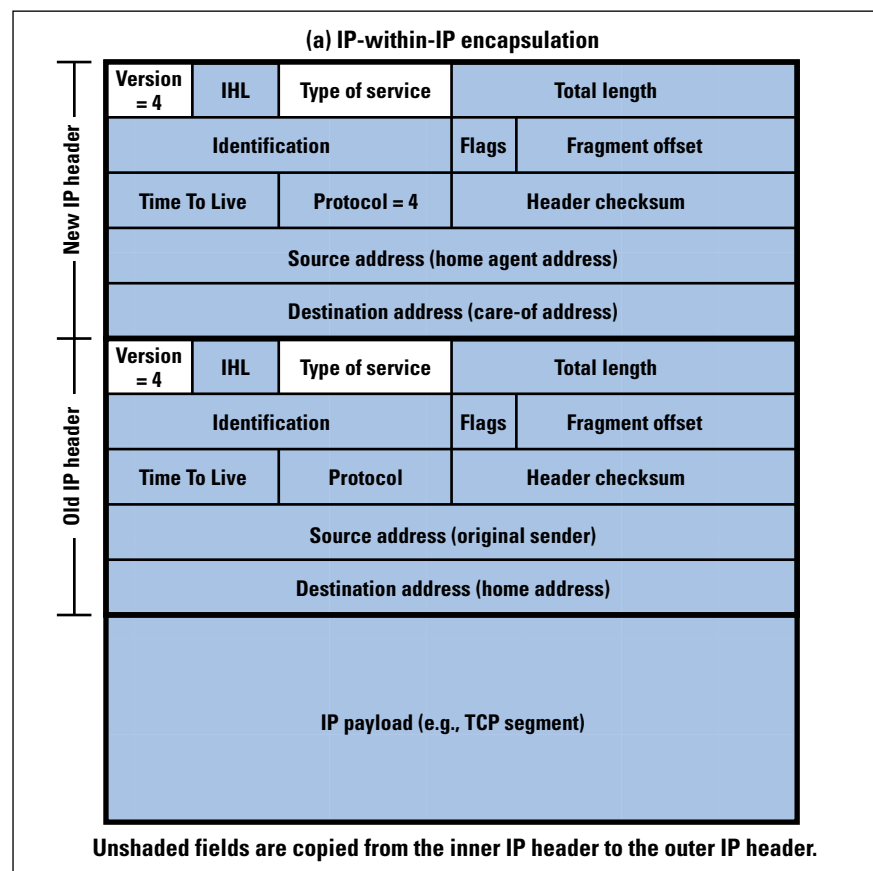
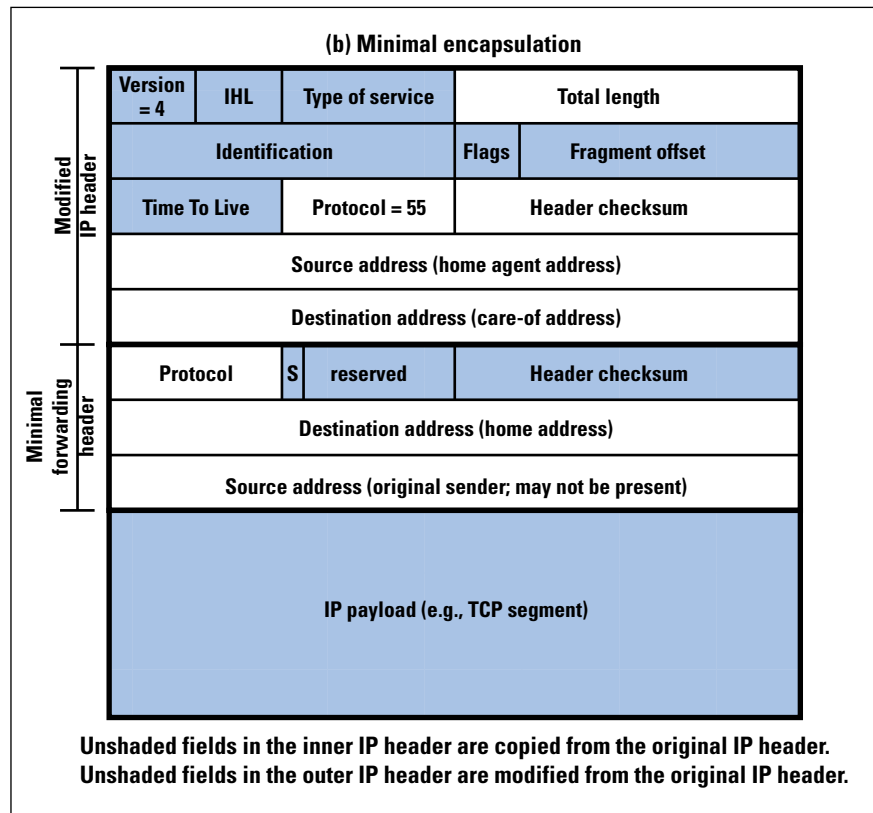


Figure 4b: Mobile IP
Encapsulation

In the inner IP header, the source address refers to the host that is sending the original datagram, and the destination address is the home address of the intended recipient. In the outer IP header, the source and destination addresses refer to the entry and exit points of the tunnel. Thus, the source address typically is the IP address of the home agent, and the destination address is the care-of address for the intended destination.

Example: Consider an IP datagram that originates at server X in Figure 1 and that is intended for mobile node A. The original IP datagram has a source address equal to the IP address of X and a destination address equal to the IP home address of A. The network portion of A's home address refers to A's home network, so the datagram is routed through the Internet to A's home network, where it is intercepted by the home agent. The home agent encapsulates the incoming datagram with an outer IP header, which includes a source address equal to the IP address of the home agent and a destination address equal to the IP address of the foreign agent on the foreign network to which A is currently attached. When this new datagram reaches the foreign agent, it strips off the outer IP header and delivers the original datagram to A.

Minimal encapsulation results in less overhead and can be used if the mobile node, home agent, and foreign agent all agree to do so. With minimal encapsulation, the new header is inserted between the original IP header and the original IP payload (Figure 4b). It includes the following fields:

- *Protocol*: Copied from the Destination Address field in the original IP header. This field identifies the protocol type of the original IP payload and thus identifies the type of header that begins the original IP payload.
- *S*: If 0, the original source address is not present, and the length of this header is 8 octets. If 1, the original source address is present, and the length of this header is 12 octets.
- *Header checksum*: Computed over all the fields of this header.
- *Original destination address*: Copied from the Destination Address field in the original IP header.
- *Original source address*: Copied from the Source Address field in the original IP header. This field is present only if the S bit is 1. The field is not present if the encapsulator is the source of the datagram (that is, the datagram originates at the home agent).

The following fields in the original IP header are modified to form the new outer IP header:

- *Total length*: Incremented by the size of the minimal forwarding header (8 or 12).
- *Protocol*: 55; this is the protocol number assigned to minimal IP encapsulation.
- *Header checksum*: Computed over all the fields of this header; because some of the fields have been modified, this value must be recomputed.
- *Source address*: The IP address of the encapsulator, typically the home agent.
- *Destination address*: The IP address of the exit point of the tunnel. This is the care-of address and may be either the IP address of the foreign agent or the IP address of the mobile node (in the case of a colocated care-of address).

The processing for minimal encapsulation is as follows. The encapsulator (home agent) prepares the encapsulated datagram with the format of Figure 4b. This datagram is now suitable for tunneling and is delivered across the Internet to the care-of address. At the care-of address, the fields in the minimal forwarding header are restored to the original IP header and the forwarding header is removed from the datagram. The total length field in the IP header is decremented by the size of the minimal forwarding header (8 or 12) and the header checksum field is recomputed.

References

Reference [1] is a good survey article on mobile IP; a somewhat less technical, more business-oriented description from the same author is [2]. For greater detail, see [3]. The August 2000 issue of *IEEE Personal Communications* contains numerous articles on enhancements to the current Mobile IP standard. The Web site of the IETF Working Group on Mobile IP, which contains current RFCs and Internet Drafts is at:

<http://ietf.org/html.charters/mobileip-charter.html>

- [1] Perkins, C., "Mobile IP," *IEEE Communications Magazine*, May 1997.
- [2] Perkins, C., "Mobile Networking through Mobile IP," *IEEE Internet Computing*, January-February 1998.
- [3] Perkins, C., *Mobile IP: Design Principles and Practices*, ISBN 0-201-63469-4, Prentice Hall PTR, 1998.
- [4] Solomon, J., *Mobile IP: The Internet Unplugged*, ISBN 0138562466, Prentice Hall PTR, 1998.

WILLIAM STALLINGS is a consultant, lecturer, and author of over a dozen books on data communications and computer networking. He also maintains a computer science resource site for CS students and professionals at **WilliamStallings.com/StudentSupport.html**. He has a PhD in computer science from M.I.T. His latest book is *Wireless Communications and Networks* (Prentice Hall, 2001). His home in cyberspace is **WilliamStallings.com** and he can be reached at **ws@shore.net**

Goodbye DES, Welcome AES

by Edgar Danielyan

Much has changed since introduction of the *Data Encryption Standard* (DES)^[2] in 1977. Our hardware is faster, we have more memory, and the use of computer networks in all areas of human activity is increasing. The widely used DES has, on several occasions, been proven to be inadequate for many applications—especially those involving the transmission of sensitive information over public networks such as the Internet, where the entire transmission may be intercepted and cryptanalyzed. Specialized hardware has been built that can determine the 56-bit DES key in a few hours. These considerations, and others, have signaled that a new standard algorithm and longer keys are necessary.

Fortunately, in January 1997, the U.S. *National Institute of Standards and Technology* (NIST) announced that it's time for a new encryption standard: the *Advanced Encryption Standard* (AES). They formalized their requirements and issued a call for candidate algorithm nominations in September 1997. The deadline for submissions was June 1998, when a total of 15 algorithms were submitted for consideration. This article shows why DES is outdated and should not be used for any purposes that require serious encryption. It also provides a brief description of the soon-to-come replacement of DES, the Advanced Encryption Standard.

Data Encryption Standard

Published as the U.S. Federal Information Processing Standard 46 in 1977, DES is still widely used, despite being proven inadequate for use in many applications. It is a symmetric block cipher (shared secret key), with its block size fixed at 64 bits. There are four defined modes of operation, with the *Electronic Code Book* (ECB) mode being the most widely used^[1]. Additionally, DES has been incorporated into numerous other standards, such as American Bankers Association's *Protection of Personal Identification Numbers in Interchange Standard*, *Management and Use of Personal Identification Numbers Standard*, *Key Management Standard*, and three ANSI standards, *Data Encryption Algorithm* (DEA), *Standard for Personal Identification Number (PIN) Management and Security*, and *Standard for Financial Institution Message Authentication*^[3]. In particular, DES is also specified as an approved algorithm in the *IP Security Architecture* (IPSec) standard^[9], which is used in the equipment from many different suppliers.

Key Length

Key length is one of the two most important security factors of any encryption algorithm—the other one being the design of the algorithm itself. DES uses a 64-bit block for the key; however, 8 of these bits are used for odd parity and are, therefore, not counted in the key length. The effective key length is then calculated as 56 bits, giving 2^{56} possible keys. A true 64-bit key has 256 times as many keys, whereas a 128-bit key is 2^{72} times “better” than a 56-bit key. As if this was not enough, DES also has so-called *weak* and *semi-weak* keys. During the encryption process, the key is used to generate two values that are used for separate purposes during the process. These 16 weak and semi-weak keys will produce values that don’t appear to be random. They will give outputs of all-ones, all-zeros, or distinguishable patterns of ones and zeros. It is generally recognized that these 16 key values should not be used. The key length was known to be a factor in trusting DES soon after DES was published. For this reason, people started exploring the use of multiple encryption passes and multiple keys. *Triple DES* (3DES) is a way of using DES encryption three times.

The most common method is to first encrypt the data block with one key. The output of this operation is run through the decryption process with a second key, and the output of that operation is run through the encryption process again with the first key. This process makes the effective key length 112 bits long. Again, the problem with weak and semi-weak keys remains. The disadvantage of Triple DES is that it is about one-third as fast as DES when processing data. This effort just slightly extended the life of DES while a suitable alternative could be found.

Breaking the DES

In addition to the brute-force key search (for example, trying every possible key in order to recover the plaintext—for DES that would be 2^{56} keys), there is also a technique known as *cryptanalysis*, which may be used to find the key or the plaintext. Essentially, there are two publicized ways to cryptanalyze DES: *differential* and *linear*. Discovered by Biham and Shamir in 1990, differential cryptanalysis was previously unknown to the public. In short, differential cryptanalysis looks at the difference between pairs of ciphertext and uses the information about these differences to find the key. Linear cryptanalysis, discovered by M. Matsui, on the other hand, uses a method called *linear approximations* to analyze block ciphers (not only DES). Because some internal structures used in DES are not designed to be strong against linear cryptanalysis, it is quite effective when used against DES. To show that the DES is inadequate and should not be used in important systems anymore, RSA Data Security^[7] sponsored a challenge to see how long it would take to decrypt successively more difficult algorithms (see <http://www.rsasecurity.com/rsalabs/challenges> for more information). Two organizations played key roles in breaking the DES: the distributed.net and the *Electronic Frontier Foundation* (EFF).

distributed.net

distributed.net^[6] is a worldwide distributed computing network. Started in 1997, the company now has thousands of participants who are contributing their idle computing power to provide an equivalent of about 160,000 Pentium II computers working in parallel. The company's mission statement says, in particular:

"We will deploy our software to form an immense, globally distributed computer that solves large-scale problems and provides an accessible pool of computational power to projects that need it. This deployment will also demonstrate the real-world utility of both distributed computing in general and our software in particular."

It may be said that they are doing well: projects undertaken and successfully completed by distributed.net include the CS Cipher, DES III, DES II 2, and RC5-56 challenges. At the time of writing, distributed.net is working on two projects: breaking RC5 with a 64-bit key and finding *Optimal Golomb Rulers* (OGRs). The idea behind distributed.net is that it is possible to distribute chunks of data over the Internet to be processed in parallel by participating computers during their idle time. The results of these calculations are then sent to a central computer that coordinates the distributed computation. The same principle is used by the SETI (Search for Extraterrestrial Intelligence) @ Home project.

Electronic Frontier Foundation

The EFF's DES cracking computer was designed by Cryptography Research, Advanced Wireless Technologies, and the EFF^[5]. The design was based upon theoretical work by Michael Wiener^[10]. It checked 90 billion keys per second, was assembled in six Sun 2 cabinets, and had 27 boards and 1800 custom chips. Built for less than \$250,000, it found the key in approximately 56 hours of brute-force search.

DES I

The DES I contest was the first attempt to prove that DES is no longer fit for any serious use. It was completed on June 17, 1997, by R. Verser in a collaborative effort, after checking about 14 percent (10,178,478, 175,420,416 keys) of the key space. It took 84 days.

DES II

There were, in fact, two DES II challenges. distributed.net participated in the first one, which began on January 13, 1998, and completed it on February 23, 1998. About 63 quadrillion keys were checked. At the end, the participants of distributed.net were checking 28 gigakeys per second. The decrypted text was "The unknown message is: Many hands make light work." The EFF won the second challenge on July 15, 1998, in less than three days, with distributed.net coming in second. This time the plaintext read "It's time for those 128-, 192-, and 256-bit keys."

DES III

The DES III contest, announced by RSA Data Security on December 12, 1998, to start on January 18, 1999, was also a success. In an official press release, RSA said:

“First adopted by the federal government in 1977, the 56-bit DES algorithm is still widely used by financial services and other industries to protect sensitive on-line applications, despite growing doubts about its vulnerability to hackers. It has been widely known that 56-bit keys, such as those offered by the government’s DES standard, offer marginal protection against a committed adversary.”

It took 22 hours and 15 minutes for Electronic Frontier Foundation’s Deep Crack computer and distributed.net’s worldwide distributed computing network to find out the 56-bit DES key, decipher the message, and win the \$10,000 contest. The decrypted message read “See you in Rome (Second AES Conference, March 22–23, 1999)” and was found after checking about 30 percent of the key space. This latest exercise finally proved that DES belongs to the past.

AES Timeline

In April 1997, NIST organized a workshop to consider criteria and submission guidelines of candidate algorithms; later in September, an official call for nominations was published in the U.S. Federal Register. By June 1998, 15 algorithms were submitted to the NIST for consideration:

- CAST-256 (Entrust Technologies)
- CRYPTON (Future Systems)
- DEAL (Richard Outerbridge, Lars Knudsen)
- DFC (National Centre for Scientific Research, France)
- E2 (NTT)
- FROG (TecApro Internacional)
- HPC (Rich Schroeppe)
- LOKI97 (Lawrie Brown, Josef Pieprzyk, Jennifer Seberry)
- MAGENTA (Deutsche Telekom)
- Mars (IBM)
- RC6 (RSA)
- Rijndael (Joan Daemen, Vincent Rijmen)
- Safer+ (Cylink)
- Serpent (Ross Anderson, Eli Biham, Lars Knudsen)
- Twofish (Bruce Schneier, John Kelsey, Doug Whiting, David Wagner, Chris Hall, Niels Ferguson)

NIST asked for public comments on these 15 algorithms and set the date for the second AES candidate conference to March 1999, to be held in Rome, Italy. The candidate algorithms were tested from both cryptological and performance viewpoints. One of the original NIST requirements for the algorithm was that it had to be efficient both in software and hardware implementations. (DES was originally practical only in hardware implementations.) Java and C reference implementations were used to do performance analysis of the algorithms. A few months later, a NIST press release announced the selection of 5 out of 15 algorithms that survived rigorous testing and cryptanalysis. This fact is not to say that the algorithms that were not selected were broken or were without merit. Those algorithms either were not as efficient, or were not as practical to implement.

The selected algorithms were Mars, RC6, Rijndael, Serpent, and Twofish. These algorithms were accepted as cryptologically strong and flexible, as well as able to be efficiently implemented in software and hardware. In August 2000, the National Security Agency published the VHDL model for performance testing of algorithms when implemented in hardware. Finally, in October 2000, a NIST press release announced the selection of Rijndael as the proposed Advanced Encryption Standard.

Rijndael

Rijndael^[4] (pronounced “Reign Dahl,” “Rain Doll,” or “Rhine Dahl”) was designed by Joan Daemen, PhD (Proton World International, Belgium) and Vincent Rijmen (Catholic University of Leuven, Belgium). Both authors are internationally known cryptographers. Rijndael is an efficient, symmetric block cipher. It supports key and block sizes of 128, 192, and 256 bits. The main design goals for the algorithm were simplicity, performance, and strength (that is, resistance against cryptanalysis). When used in *Cipher Block Chaining Message Authentication Code* (CBC MAC) mode, Rijndael can be used as a MAC algorithm; it also may be used as a hash function and as a pseudo random number generator (both are special mathematical functions widely used in cryptography; an example of a hash function is *Message Digest 5* (MD5)—a popular message digest algorithm by Ron Rivest). In their specification of the algorithm, the authors specifically state the strength of Rijndael against differential, truncated differential, linear, interpolation, and Square attacks. Although Rijndael is not based on Square^[8], some ideas from the Square algorithm design are used in Rijndael.

Square is a 128-bit symmetric iterated block cipher designed by Daemen, Rijmen, and Knudsen. Its primary design goal was strength against both linear and differential cryptanalyses; the high degree of parallelism of the Square algorithm allows efficient implementation on parallel computers.

Of course, the length of the key is also very important, especially because the most efficient known attack against Rijndael is an exhaustive key search. It would take 2^{255} runs of Rijndael to find a key 256 bits long. To the credit of the authors, Rijndael does not use “parts” or tables from other algorithms, making it easy to implement alone.

Table 1: Comparing DES and AES

	DES	AES
Key Length	56 bits	128, 192, or 256 bits
Cipher Type	Symmetric block cipher	Symmetric block cipher
Block Size	64 bits	128, 192, or 256 bits
Developed	1977	2000
Cryptanalysis resistance	Vulnerable to differential and linear cryptanalysis; weak substitution tables	Strong against differential, truncated differential, linear, interpolation and Square attacks
Security	Proven inadequate	Considered secure
Possible Keys	2^{56}	2^{128} , 2^{192} , or 2^{256}
Possible ASCII printable character keys*	95^7	95^{16} , 95^{24} , or 95^{32}
Time required to check all possible keys at 50 billion keys per second**	For a 56-bit key: 400 days	For a 128-bit key: 5×10^{21} years

* When a text password input by a user is used for encryption (there are 95 printable characters in ASCII).

**In theory, the key may be found after checking 1/2 of the key space. The time shown is 100% of the key space.

Summary

It is expected that AES will be officially published as a *Federal Information Processing Standard* (FIPS) in April–June 2001, and implementations of AES in various security systems probably will surface shortly thereafter. In the meantime, authoritative information on AES developments may be found on NIST’s Web site at <http://csrc.nist.gov/encryption/aes/>. The full mathematical specification of the algorithm and reference implementations in C and Java are also available from the same Web site.

References

- [1] *Applied Cryptography*, 2nd edition, by Bruce Schneier, 1996, John Wiley & Sons.
- [2] National Institute of Standards and Technology (NIST),
<http://www.nist.gov>
- [3] American National Standards Institute (ANSI),
<http://www.ansi.org>
- [4] The Rijndael Specification, <http://csrc.nist.gov/encryption/aes/rijndael/Rijndael.pdf>
- [5] Electronic Frontier Foundation, <http://www.eff.org>
- [6] distributed.net, <http://www.distributed.net>
- [7] RSA Security, <http://www.rsa.com>
- [8] Square Specification,
<http://www.esat.kuleuven.ac.be/~rijmen/square>
- [9] Kent, S., Atkinson, R., "Security Architecture for the Internet Protocol," RFC 2401, November 1998.
- [10] Michael Wiener, "Efficient DES Key Search," Proceedings of the CRYPTO'93 Conference, August 1993.
- [11] Madson, C., Doraswamy, "The ESP DES-CBC Cipher Algorithm With Explicit IV," RFC 2405, November 1998.

[A prior version of this article was published in the February 2001 issue of the *login:* magazine].

EDGAR DANIELYAN is a Cisco Certified Network, Design and Security Professional, as well as member of ACM, USENIX, SAGE, and the IEEE Computer Society. He has worked for a national telco, a bank, the United Nations, and the Ministry of Defense, among others. Currently self-employed, he consults and writes on internetworking, UNIX, and security. E-mail: edd@danielyan.com

Opinion: The Middleware Muddle

by Geoff Huston

[This occasional column is an individual soapbox on views of various aspects of the Internet. The views stated here are intended to be mildly provocative, and, if backed to the wall, the author will rapidly disclaim any responsibility for them whatsoever!]

It is not often that an entire class of technology can generate an emotive response. But, somehow, middleware has managed to excite many strong reactions. For some *Internet Service Providers* (ISPs), middleware—in the form of *Web caches*—is not only useful, it's critical to the success of their enterprise. For many corporate networks, middleware—in the form of *firewalls*—is the critical component of their network security measures. For such networks, middleware is an integral part of the network. Other networks use middleware, in the form of *Network Address Translators* (NATs), as a means of stretching a limited number of Internet public addresses to provide connectivity services to a much larger local network. For others, middleware is seen as something akin to network heresy. For them, not only does middleware often break the basic semantics of the Internet Protocol, it is also in direct contravention to the end-to-end architecture of the Internet. Middleware, they claim, breaks the operation of entire classes of useful applications, and this makes the Internet a poorer network as a result.

Emotions have run high in the middleware debate, and middleware has been portrayed as being everything from absolutely essential to the operation of the Internet as we know it, to being immoral and deceptive. Strong stuff indeed from an engineering community, even one as traditionally opinionated as Internet engineers.

So what is middleware all about and why the fuss?

It may be helpful to start with a definition of middleware. One definition of middleware is that of anything in the network that functions at a level in a network reference model above that of end-to-end transport (TCP/IP), and below that of the application environment (the *Application Programming Interface* [API])^[1]. Of course, this definition encompasses a very broad class of services that covers everything from *Authentication, Authorization, and Accounting* (AAA) servers and *Domain Name System* (DNS) servers through to various forms of information discovery services and resource management.

Another possible definition of middleware adopts the perspective of the integrity of the end-to-end model of Internet architecture^[2]. From this perspective, middleware is a class of network devices that do something other than forward or discard an IP packet onward along the next hop to the destination address of the packet—in other words, anything other than a packet-switching element that sits in the transmission path of the packet.

With such an end-to-end definition of middleware, these middleware units may intercept the packet and alter the header or payload of the packet, redirect the packet to be delivered to somewhere other than its intended destination, or process the packet as if it were addressed to the middleware device itself. From this perspective, AAA, the DNS, and related services from our first definition are simply applications that traverse the network.

There's nothing like confusion over definitions to fuel a debate, and this area is no exception. However, a debate over definitions is too often a dry one. So, in the interest of adding a little more incendiary material to the topic, let's simply use this second definition of middleware to look further at the issues.

Why would a network go to all this bother to trap and process certain packets? Surely it's easier and cheaper to simply forward the packet onward to its intended destination? The answer can be "yes" or "no," depending on how you feel about the role of middleware in TCP/IP.

An Example: Cache Middleware

Let's look at this in a bit more detail, using a specific flavor of middleware to illustrate the middleware dilemma. A common form of middleware is the *Transparent Web Cache*. Such a Web cache is constructed using two parts, an *interceptor* and a *cache system*. The interceptor is placed into the network, either as a software module added to a router or as a device, which is spliced into a point-to-point link. The interceptor takes all incoming TCP traffic addressed to port 80 (a *Hypertext Transfer Protocol* [HTTP] session) and redirects it across to the cache system. All other traffic is treated normally. The cache system accepts all such redirected packets as if they were directly addressed to the cache itself. It responds to the HTTP requestor as if it were the actual intended destination, using a source address that matches the destination address of the original request, assuming the identity of the actual intended content server. If the requested Web object is located in the local cache, it will deliver the object to the requestor immediately. If it is not in the cache, it will set up its own session with the original destination, send it the original request, and feed the response back to the requestor, while also keeping a copy for itself in its cache.

Caching of content works well in the Web world simply because so much Web traffic today is movement of the same Web page to different recipients. It is commonly reported that up to one half of all Web traffic in the Internet is a duplicate transmission of content. If an ISP locally caches all Web content as it is delivered, and checks the cache before passing through a content request, then the ISP's upstream Web traffic volume may be halved. Even a moderately good cache will be able to service about one quarter of the Web content from the cache. That amount of local caching can be translated into a significant cost saving for the ISP.

The cached Web content is traffic that is not purchased as transit traffic from an upstream ISP, representing a potential saving on the cost of upstream transit services. This saving, in turn, can allow the ISP to operate at a lower price point in the retail market. The cache is also located closer to the ISP's customers, and with appropriate tuning, the cache can also deliver cached content to the customer at a consistently much faster rate than a request to the original content server. For very popular Web sites the originating server may be operating more slowly under extreme load, while the local cache continues to operate at a more consistent service level. The combination of the potential for improved performance and lower overall cost is certainly one that looks enticing: the result is the same set of Web transactions delivered to customers, but cheaper and faster.

End-to-End Issues with Cache Middleware

But not everything is perfect in this transparent caching world. What if the Web server used a security model that served content only to certain requestors, and the identity of the requestor was based on their IP address? This is not a very good security model, admittedly, but it's simple, and because of its simplicity this practice enjoys very common usage. With the introduction of a transparent cache, the Web client sees something quite strange. The Web client can ping the Web server, the client can communicate with any other port on the server, and if the client were to query the status of the server, the Web server would be seen to be functioning quite normally. But, mysteriously, the client cannot retrieve any Web content from the server, and the server does not see any such request from the client. The middleware cache is sitting inside a network somewhere on the path between the client and the service, but it may well be the case that neither the end client or the end server are aware of the deployment of the middleware unit. It is not surprising that this is a remarkably challenging operational problem for either the client or the server to correctly diagnose.

A similar case is where a Web server wishes to deliver different content to different requestors, based on some inference gained from the source IP address of the requestor, or the time of day, or some other variable derived from the circumstances of the request. A transparent cache will not detect such variations in the response of the server and will instead deliver the same version of the cached content to all clients whose requests pass through the transparent cache. Variations of this situation of perceived abnormal service behavior abound, all clustered around the same concept that it is unwise in such an environment for a server to assume that it is always communicating with the end client. Indeed the situation is common enough that the Web application has explicit provision for instructing cache servers about whether the content can be cached and replayed in response to similar subsequent requests.

More subtle vulnerabilities also are present in such a middleware environment. A client can confidently assert that packets are being sent to a server, and the server appears to be responding, but the data appears to have been corrupted. Has the server been compromised? It may look like this is the case, but when middleware is around, looks can be deceiving. If the integrity of the cache is compromised, and different pages are substituted in the cache, then to the clients of the cache it appears that the integrity of original server has been compromised. The twist with transparent cache middleware is that the clients of the cache may be unaware that the cache exists, let alone that their requests are being redirected to the cache server. Any abnormalities in the responses they receive are naturally attributed to problems with the security of the server and the integrity of the associated service.

The common theme of these issues is that there are sets of inconsistent assumptions at play here. On the one hand, the assumption of an end-to-end architecture leads an application designer to assume that an IP session opened with a remote peer will indeed be with that remote peer, and not with some intercepting network-level proxy agent attempting to mimic the behavior of that remote peer. On the other hand, is the assumption that transactions adhere to a consistent and predictable protocol, and transactions may be intercepted and manipulated by middleware as long as the resultant interaction behaves according to the defined protocol.

Middleware Architecture

Are transparent caches good or bad? Is the entire concept of middleware good or bad?

There is no doubt that middleware can be very useful. Cache systems can create improved service quality and reduced cost. NATs can reduce the demand for public IP address space. Firewalls can be effective as security policy agents. Middleware can perform load balancing across multiple service points for a particular class of applications, such as a Web server farm. Middleware can dynamically adjust the Internal Protocol parameters of a TCP session to adapt to particular types of networks, or various forms of network service policies. Middleware can provide services within the network that relieve the end user of a set of tasks and responsibilities, and middleware can improve some aspects of the service quality. Middleware can make an Internet service faster, cheaper, more flexible, and more secure, although probably not all at the same time. But middleware comes at a steep long-term price.

The advantage of the Internet lies in its unique approach to network architecture. In a telephone network, the end device—a telephone handset—is a rather basic device consisting of a pair of transducers and a tone generator. All the functionality of the telephone service is embedded within the network itself.

The architecture of the Internet is the complete opposite. The network consists of a collection of packet switches with basic functionality. The service is embedded within the protocol stack and the set of applications that are resident on the connected device. Within this architecture, adding new services to the network is as simple as distributing new applications among those end systems that want to use the application. The network makes no assumptions about the services it supports, and network services can be added, refined, and removed without requiring any change to the network itself. This results in a cheap, flexible, and basic network, and it passes the entire responsibility for service control to the network users. The real strength of the Internet lies in its architectural simplicity and lack of complex interdependencies within the network.

Middleware cuts across this model by inserting directly into the network functionality that alters the behavior of the network. IP or TCP Packet Header fields may be altered on the fly, or, as with a transparent cache, middleware may intercept user traffic, use an application level interpreter to interpret the upper-level service request associated with the traffic, and generate a response, acting as an unauthorized proxy for the intended recipient. With middleware present in an IP network, sending a packet to an addressed destination and receiving a response with a source address of that destination is no guarantee that you have actually communicated with the addressed remote device. You may instead be communicating with a middleware box, or have had the middleware box alter your traffic in various ways that are not directly visible to the sender.

In such an environment, it's not just the end-user applications that define an Internet-deployed service, because middleware is also part of the Internet service architecture. Services may be deployed that are reliant on the existence of middleware to be effective. Streaming video services, for example, become far more viable as a scalable Internet service when the streaming video server content is replicated across a set of middleware streaming systems deployed close to end users of the service. To change the behavior of a service that has supporting middleware deployed requires the network middleware to be changed. A new service may not be deployed until the network middleware is altered to permit its deployment. Any application requiring actual end-to-end communications may have to have additional functionality to detect if there is network middleware deployed along the path, and then explicitly negotiate with this encountered middleware to ensure that its actual communication will not be intercepted and proxied or otherwise altered.

Conclusion

The cumulative outcome is that such a middleware-modified Internet service model is not consistent with an end-to-end architecture. It represents the introduction of a more muddled service architecture where the network may choose to selectively intervene in the interaction between one device and another. Such a network architecture may not have stable scaling properties. Such an architecture may not readily support entire classes of new applications and new services. Such an architecture may not be sufficiently flexible and powerful to underpin a ubiquitous global data communications system. All this middleware overhead makes applications more complex, makes the network more complex, and makes networking more expensive, more limited, and less flexible.

From this perspective, middleware is an unglamorous hack. To adapt a 350-year-old quote from Thomas Hobbes, middleware is nasty, brutish, and short-sighted. It is, hopefully, a temporary imposition on an otherwise elegant, simple, and adequate Internet architecture.^[3, 4]

References

- [1] Aiken, B. et.al, "Network Policy and Services: A Report of a Workshop on Middleware," RFC 2768, February 2000.
- [2] Carpenter, B. ed., "Architectural Principles of the Internet," RFC 1958, June 1996.
- [3] Hobbes, Thomas (1588–1679), *Leviathan*, London, 1651. Available from many sources, including, ISBN 0140431950, Penguin Press, 1982.
- [4] <http://www.orst.edu/instruct/phl302/texts/hobbes/leviathan-contents.html>

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. Huston is currently the Chief Scientist in the Internet area for Telstra. He is also a member of the Internet Architecture Board, and is the Secretary of the Internet Society Board of Trustees. He is author of *The ISP Survival Guide*, ISBN 0-471-31499-4, *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, and coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. All three books are published by John Wiley & Sons. E-mail: gih@telstra.net

Book Review

Internetworking with TCP/IP

Internetworking with TCP/IP (Vol. 1): Principles, Protocols, and Architectures, Douglas E. Comer, ISBN 0-13-018380-6, Prentice Hall, 2000.

Internetworking With TCP/IP (Vol. 1): Principles, Protocols, and Architectures (fourth edition) is the latest update to Comer's landmark work containing *Internetworking With TCP/IP (Vol. 2): Design, Implementation, and Internals* and *Internetworking With TCP/IP (Vol. 3): Client-Server Programming and Applications/BSD Socket Version*. As a recent engineering graduate, I wish I had read this book sooner; it is very concise and would have saved me a lot of time early in my studies.

Comer imparts Volume 1 in four sections. The first section provides a basic introduction to general networking including descriptions of typical network components. This section is most helpful for the entry-level student or casual reader. Advanced readers may want to skip right to the next section of the text, which continues with coverage of the TCP/IP networking environment from the host's point of view. Here, the organization and operation of local host protocols, addressing, and routing are thoroughly discussed. After reading this portion of the book, you will definitely understand how your desktop computer communicates on the network. Next, the global Internet architecture is laid out in a very comprehensible format. The reader is introduced to router-to-router protocols and algorithms that don't seem so complicated after this treatment. Lastly, application-level services and the client-server model of networking are covered in the final portion of the book.

Classic Reference

When reviewing one of the eminent texts in the field, it is of limited use to comment on the work chapter by chapter. However, I am compelled to comment on the quality of Chapter 11, Protocol Layering. This chapter is particularly interesting because Comer directly compares the ISO 7-layer reference model to the TCP/IP 5-layer model. As is par for this book, the comparison is clear and concise. Furthermore, the advantages and disadvantages of protocol layering are discussed in general and a realistic perspective is provided with reference to actual software implementation practices which may result in layer blurring. This is a very cogent presentation of the interaction between theory and reality in engineering. Although covering a specific topic, it could easily serve as an object lesson in a discussion of "real world" engineering techniques. In addition to Chapter 11, the chapters covering Internet routing (14 through 16) really shine as mainstays of this book. The Internet is viewed from the top down and "big network" protocols such as the *Border Gateway Protocol* (BGP) are given good coverage. This is an area where very few people are completely comfortable and Comer once more brings the important material forward in an easily understandable fashion. In the following paragraphs, I will highlight some of the new material included in the fourth edition.

New TCP/IP Concepts

The book's handling of *Classless Inter-Domain Routing* (CIDR) is very informative. In addition to explaining the inner-workings of the address space, Comer points out the requirement for new routing algorithms. This is an associated cost of adopting this new concept that is often overlooked when CIDR is presented.

Two new and important IP topics are also well-presented. Comer begins his treatment of IP Version 6 (IPv6) with a quick history of the protocol and a review of the logic behind this change. The new address space notation and allocation by type are explained very well. New advantages provided under IPv6 protocol structures are then discussed. Additionally, Mobile IP concepts and practicalities are introduced. Comer does a good job of bringing out both good news and bad news of this crucial new networking technology.

Coverage of *Random Early Drop* (RED) was rather brief and really needs more detail before readers can thoroughly grasp the concept. However, this would require greater mathematical sophistication on the part of the reader. Accordingly, depth of coverage is forgone in the interest of readability.

The section on *Network Address Translation* (NAT) does not adequately explain the dynamic nature of IP address assignment across hosts and data flows. An additional detailed example would help here.

Multimedia

In the application-level services section of the book, Comer offers a hasty explanation of how voice and video are sent over IP internets and how IP Telephony operates. The H.323 protocol is briefly mentioned as the low-bandwidth videoconferencing standard. However, it is not presented in its full importance as an umbrella recommendation from the *International Telecommunications Union* (ITU). A chapter explaining the roles of subordinate H.320 protocols in general would be a welcome addition to this section. *Quality of Service* (QoS) concepts such as *Resource Reservation Protocol* (RSVP), *Differentiated Services* (Diff-Serv), and *Real Time Protocol* (RTP) are likewise given short rift. However, IP Multicast is given significant treatment in one of the book's longest chapters; its concepts, mechanics, and implementation choices are thoroughly addressed.

Security

The book provides clear introductions to *Virtual Private Networks* (VPNs) and the IPsec set of protocols. The actual mechanics of IPsec are detailed thoroughly. Various required algorithms are introduced and pertinent RFC references are pointed out. Finally, firewall basics and implementation issues are covered. Overall, these sections clearly define the pertinent security concepts and make them simple.

Prerequisite Knowledge

This book thoroughly covers the fundamental principles of network design including implementation trade-offs and their associated foibles. However, understanding this text requires little more than a modest understanding of basic computer and networking concepts. An introductory programming course that covers computer organization, the binary number system, and basic data structures should suffice. From this point, the student can use the text for initial network familiarization as well as a future reference to ground the more abstract topics in network design.

A Must-Have Reference

An extensive, concept-based overview of the TCP/IP internetworking protocols makes Comer's Volume 1 the classic introduction to TCP/IP. He makes this an enjoyable read by breaking the topic into short, digestible chapters. Additionally, Comer pauses throughout the text to intersperse review material. Recurrent, italicized summaries provide a significant advantage to the student. These asides concisely summarize key points and provide a coherent set of landmarks for quick review and study.

By itself, Volume 1 is broad enough to be complete as an introduction to IP networking protocols. Comer further extends the work by pointing the reader to very specific resources for in-depth information including web pages and specific RFC numbers for applicable topics at the end of each chapter. One of life's simple treasures is found in the *Guide to RFCs* (Appendix 1). Here, the first 2728 RFCs are organized by major categories and subtopics. At last, a navigable index of RFCs has been incorporated with a superb text from which the beginner can delve the body of networking knowledge.

—Albert C. Kinney
kinney@ieee.org

Would You Like to Review a Book for IPJ?

We receive numerous books on computer networking from all the major publishers. If you've got a specific book you are interested in reviewing, please contact us and we will make sure a copy is mailed to you. The book is yours to keep if you send us a review. We accept reviews of new titles, as well as some of the "networking classics." Contact us at ipj@cisco.com for more information.

This publication is distributed on an "as-is" basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Fragments

Jonathan B. Postel Service Award for 2001 Presented to Daniel Karrenberg

Internet Society (ISOC) Chairman Brian Carpenter presented the 2001 *Jonathan Postel Service Award* to Mr. Daniel Karrenberg, one of the pioneers of the Internet's development in Europe, during the opening ceremony of the 2001 INET Conference. His early work was at the University of Dortmund creating a basic networked e-mail and USENET service. The success of this initiative was the seed on which the first pre-commercial network, EUnet, was built. As the Internet came to Europe in the late 1980s, Mr. Karrenberg was active in organizing the first RIPE meeting and in creating the RIPE NCC to serve as secretariat for the Internet community in Europe. The RIPE NCC became the first *Regional Internet Registry* as we know them, taking on address allocation as one of its core services. Daniel headed the effort from the start, working hard to maximize the benefit for the community.

Mr. Karrenberg humbly accepted the award, thanking the Internet community for this recognition and pledging to continue his work guided by the spirit of Jon Postel.

The Jonathan B. Postel Service Award was established by the Internet Society to honor a person who has made outstanding contributions in service to the data communications community. It is named for Dr. Jonathan B. Postel to recognize and commemorate the extraordinary stewardship exercised by Jon over the course of a thirty year career in networking. The Award consists of an engraved crystal globe and US \$20,000.00. The first award was presented posthumously to Jon Postel himself, accepted by his mother, Lois Postel at INET '99. Scott Bradner received the second award during INET 2000. For additional information on Jon Postel's life and contributions, please visit:

<http://www.isoc.org/postel/>

RFC 1149 Implemented

The Internet Engineering Task Force (IETF) has a long tradition of publishing humorous *Request For Comments* (RFCs) each year on April 1st. One of the more famous such RFCs is "A Standard for the Transmission of IP Datagrams on Avian Carriers," RFC 1149, by David Waitzman, published on April 1, 1990. This "carrier pigeon" RFC was recently implemented by a group in Bergen, Norway. For details see:

<http://www.blug.linux.no/rfc1149/>

Jon Crowcroft Joins IPJ Editorial Advisory Board

We are pleased to announce that Dr. Jon Crowcroft of University College London has joined the Editorial Advisory Board for the *Internet Protocol Journal* (IPJ). Dr. Crowcroft has been working in the field of internetworking and protocol design since the early days of the ARPANET. For more information, see:

<http://www.cs.ucl.ac.uk/staff/J.Crowcroft/>

We would also like to thank Edward Kozel, the creator of IPJ, for his support and advice over the last three years. Mr. Kozel has left Cisco to pursue other interests.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Technology
WorldCom, USA

Dr. Jon Crowcroft, Professor of Networked Systems
University College London, England

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

Copyright © 2001 Cisco Systems Inc.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-10/5
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSR STD U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol Journal

September 2001

Volume 4, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
MPLS	2
A Unique Root	15
Book Review	29
Call for Papers	31
Fragments	32

FROM THE EDITOR

Multiprotocol Label Switching (MPLS) is a technology that has received a great deal of attention in recent years. The IETF alone has produced over 300 Internet Drafts and numerous RFCs related to MPLS and continues its work on refining the standards. So, what is MPLS all about? We asked Bill Stallings to give us a basic tutorial.

The tragic events of September 11, 2001 have focused attention on the stability and robustness of the Internet. The Internet played an important role in the aftermath of the terrorist attacks. While popular news Web sites initially appeared overloaded, a great deal of private traffic in the form of instant messaging and e-mail took place. Companies directly or indirectly affected by the events in New York and Washington were quick to use the Web as a way to disseminate important information to their clients as well as to their employees. In many cases, the Internet was used in place of an overloaded telephone network. With this in mind, The *Internet Corporation for Assigned Names and Numbers* (ICANN) has decided to re-focus its next meeting to address issues of Internet stability and security, particularly with regard to naming and addressing. (See "Fragments," page 32.) To provide some background information, we bring you the article "A Unique, Authoritative Root for the DNS," by M. Stuart Lynn, the president and CEO of ICANN. Since this article has been posted for public comment, you are encouraged to address your feedback to: comments@icann.org

We would like to remind our readers to send us postal address updates. The computer-communications industry is one where people change jobs and locations often. While we do receive some address changes automatically when mail is returned to us, it is much more reliable to send us e-mail with the new information. In the near future, readers will be able to make address changes and select delivery options through a Web interface which will be deployed at <http://www.cisco.com/ipj>. Until then, please send your updates to ipj@cisco.com

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

MPLS

by William Stallings

Multiprotocol Label Switching (MPLS) is a promising effort to provide the kind of traffic management and connection-oriented *Quality of Service* (QoS) support found in *Asynchronous Transfer Mode* (ATM) networks, to speed up the IP packet-forwarding process, and to retain the flexibility of an IP-based networking approach.

Background

The roots of MPLS go back to numerous efforts in the mid-1990s to combine IP and ATM technologies. The first such effort to reach the marketplace was IP switching, developed by Ipsilon. To compete with this offering, numerous other companies announced their own products, notably Cisco Systems (Tag Switching), IBM (aggregate route-based IP switching), and Cascade (IP Navigator). The goal of all these products was to improve the throughput and delay performance of IP, and all took the same basic approach: Use a standard routing protocol such as *Open Shortest Path First* (OSPF) to define paths between end-points; assign packets to these paths as they enter the network; and use ATM switches to move packets along the paths. When these products came out, ATM switches were much faster than IP routers, and the intent was to improve performance by pushing as much of the traffic as possible down to the ATM level and using ATM switching hardware.

In response to these proprietary initiatives, the *Internet Engineering Task Force* (IETF) set up the MPLS working group in 1997 to develop a common, standardized approach. The working group issued its first set of Proposed Standards in 2001. Meanwhile, however, the market did not stand still. The late 1990s saw the introduction of many routers that are as fast as ATM switches, eliminating the need to provide both ATM and IP technology in the same network.

Nevertheless, MPLS has a strong role to play. MPLS reduces the amount of per-packet processing required at each router in an IP-based network, enhancing router performance even more. More significantly, MPLS provides significant new capabilities in four areas that have ensured its popularity: QoS support, traffic engineering, *Virtual Private Networks* (VPNs), and multiprotocol support. Before turning to the details of MPLS, we briefly examine each of these.

Connection-Oriented QoS Support

Network managers and users require increasingly sophisticated QoS support for numerous reasons. The following are key requirements:

- Guarantee a fixed amount of capacity for specific applications, such as audio/video conference
- Control latency and jitter and ensure capacity for voice
- Provide very specific, guaranteed, and quantifiable service-level agreements, or traffic contracts
- Configure varying degrees of QoS for multiple network customers

A connectionless network, such as in IP-based internetwork, cannot provide truly firm QoS commitments. A *Differentiated Service* (DS) framework works in only a general way and upon aggregates of traffic from numerous sources. An *Integrated Services* (IS) framework, using the *Resource Reservation Protocol* (RSVP), has some of the flavor of a connection-oriented approach, but is nevertheless limited in terms of its flexibility and scalability. For services such as voice and video that require a network with high predictability, the DS and IS approaches, by themselves, may prove inadequate on a heavily loaded network. By contrast, a connection-oriented network has powerful traffic-management and QoS capabilities. MPLS imposes a connection-oriented framework on an IP-based internet and thus provides the foundation for sophisticated and reliable QoS traffic contracts.

Traffic Engineering

MPLS makes it easy to commit network resources in such a way as to balance the load in the face of a given demand and to commit to differential levels of support to meet various user traffic requirements. The ability to dynamically define routes, plan resource commitments on the basis of known demand, and optimize network utilization is referred to as *traffic engineering*.

With the basic IP mechanism, there is a primitive form of automated traffic engineering. Specifically, routing protocols such as OSPF enable routers to dynamically change the route to a given destination on a packet-by-packet basis to try to balance load. But such dynamic routing reacts in a very simple manner to congestion and does not provide a way to support QoS. All traffic between two endpoints follows the same route, which may be changed when congestion occurs. MPLS, on the other hand, is aware of not just individual packets, but flows of packets in which each flow has certain QoS requirements and a predictable traffic demand. With MPLS, it is possible to set up routes on the basis of these individual flows, with two different flows between the same endpoints perhaps following different routers. Further, when congestion threatens, MPLS paths can be rerouted intelligently. That is, instead of simply changing the route on a packet-by-packet basis, with MPLS, the routes are changed on a flow-by-flow basis, taking advantage of the known traffic demands of each flow. Effective use of traffic engineering can substantially increase usable network capacity.

VPN Support

MPLS provides an efficient mechanism for supporting VPNs. With a VPN, the traffic of a given enterprise or group passes transparently through an internet in a way that effectively segregates that traffic from other packets on the internet, proving performance guarantees and security.

Multiprotocol Support

MPLS, which can be used on many networking technologies, is an enhancement to the way a connectionless IP-based internet is operated, requiring an upgrade to IP routers to support the MPLS features. MPLS-enabled routers can coexist with ordinary IP routers, facilitating the introduction of evolution to MPLS schemes. MPLS is also designed to work in ATM and Frame Relay networks. Again, MPLS-enabled ATM switches and MPLS-enabled Frame Relay switches can be configured to coexist with ordinary switches. Furthermore, MPLS can be used in a pure IP-based internet, a pure ATM network, a pure Frame Relay network, or an internet that includes two or even all three technologies. This universal nature of MPLS should appeal to users who currently have mixed network technologies and seek ways to optimize resources and expand QoS support.

For the remainder of this discussion, we focus on the use of MPLS in IP-based internets, with brief comments about formatting issues for ATM and Frame Relay networks.

MPLS Operation

An MPLS network or internet consists of a set of nodes, called *Label Switched Routers* (LSRs), that are capable of switching and routing packets on the basis of a label which has been appended to each packet. Labels define a flow of packets between two endpoints or, in the case of multicast, between a source endpoint and a multicast group of destination endpoints. For each distinct flow, called a *Forwarding Equivalence Class* (FEC), a specific path through the network of LSRs is defined. Thus, MPLS is a connection-oriented technology. Associated with each FEC is a traffic characterization that defines the QoS requirements for that flow. The LSRs do not need to examine or process the IP header, but rather simply forward each packet based on its label value. Therefore, the forwarding process is simpler than with an IP router.

Figure 1: MPLS Operation

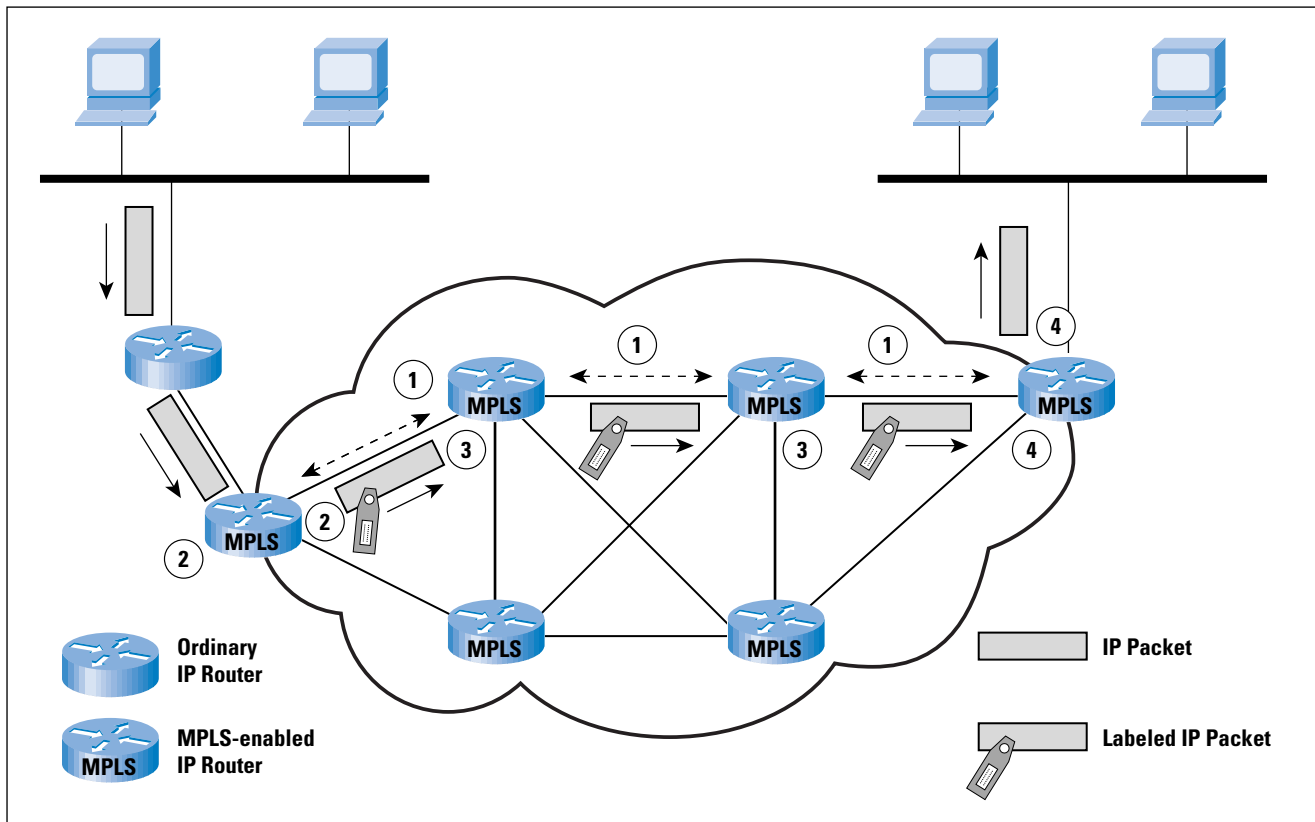


Figure 1, based on one in^[4], depicts the operation of MPLS within a domain of MPLS-enabled routers. The following are key elements of the operation.

1. Prior to the routing and delivery of packets in a given FEC, a path through the network, known as a *Label Switched Path* (LSP), must be defined and the QoS parameters along that path must be established. The QoS parameters determine (1) how many resources to commit to the path, and (2) what queuing and discarding policy to establish at each LSR for packets in this FEC. To accomplish these tasks, two protocols are used to exchange the necessary information among routers:
 - (a) An interior routing protocol, such as OSPF, is used to exchange reachability and routing information.
 - (b) Labels must be assigned to the packets for a particular FEC. Because the use of globally unique labels would impose a management burden and limit the number of usable labels, labels have local significance only, as discussed subsequently. A network operator can specify explicit routes manually and assign the appropriate label values. Alternatively, a protocol is used to determine the route and establish label values between adjacent LSRs. Either of two protocols can be used for this purpose: the *Label Distribution Protocol* (LDP) or an enhanced version of RSVP.

2. A packet enters an MPLS domain through an ingress edge LSR where it is processed to determine which network-layer services it requires, defining its QoS. The LSR assigns this packet to a particular FEC, and therefore a particular LSP, appends the appropriate label to the packet, and forwards the packet. If no LSP yet exists for this FEC, the edge LSR must cooperate with the other LSRs in defining a new LSP.
3. Within the MPLS domain, as each LSR receives a labeled packet, it:
 - (a) Removes the incoming label and attaches the appropriate outgoing label to the packet.
 - (b) Forwards the packet to the next LSR along the LSP.
4. The egress edge LSR strips the label, reads the IP packet header, and forwards the packet to its final destination.

Several key features of MLSP operation can be noted at this point:

1. An MPLS domain consists of a contiguous, or connected, set of MPLS-enabled routers. Traffic can enter or exit the domain from an endpoint on a directly connected network, as shown in the upper-right corner of Figure 1. Traffic may also arrive from an ordinary router that connects to a portion of the internet not using MPLS, as shown in the upper-left corner of Figure 1.
2. The FEC for a packet can be determined by one or more of a number of parameters, as specified by the network manager. Among the possible parameters:
 - Source or destination IP addresses or IP network addresses
 - Source or destination port numbers
 - IP protocol ID
 - Differentiated services codepoint
 - IPv6 flow label
3. Forwarding is achieved by doing a simple lookup in a predefined table that maps label values to next-hop addresses. There is no need to examine or process the IP header or to make a routing decision based on destination IP address.
4. A particular *Per-Hop Behavior* (PHB) can be defined at an LSR for a given FEC. The PHB defines the queuing priority of the packets for this FEC and the discard policy.
5. Packets sent between the same endpoints may belong to different FECs. Thus, they will be labeled differently, will experience different PHB at each LSR, and may follow different paths through the network.

Figure 2: MPLS Packet Forwarding

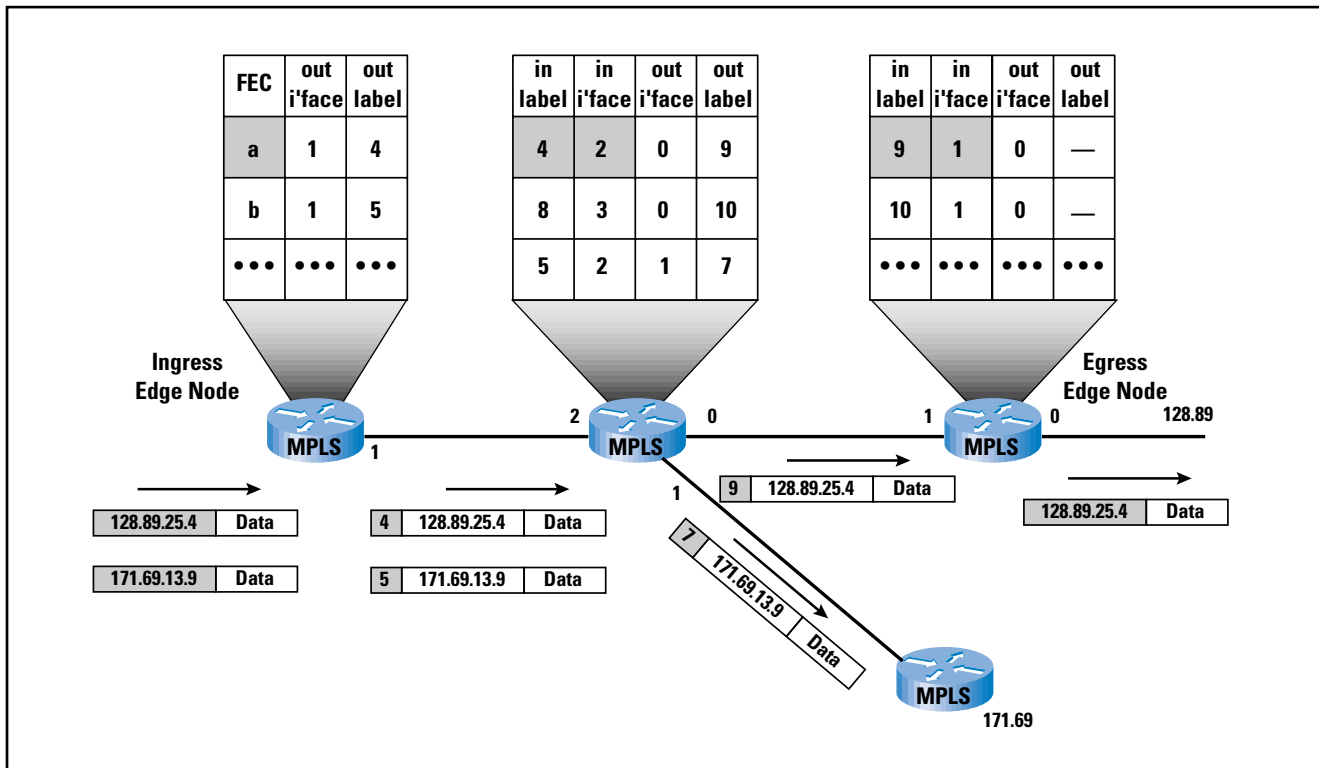


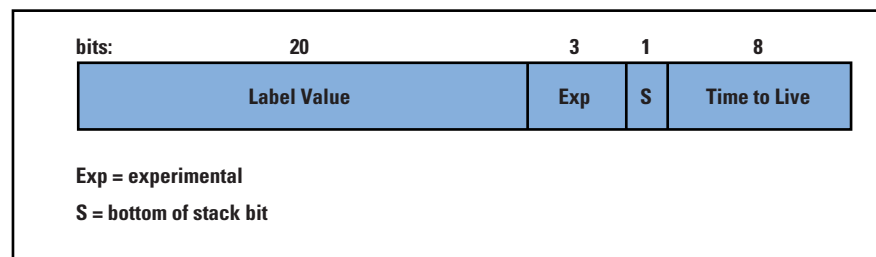
Figure 2 shows the label-handling and label-forwarding operation in more detail. Each LSR maintains a forwarding table for each LSP passing through the LSR. When a labeled packet arrives, the LSR indexes the forwarding table to determine the next hop. For scalability, as was mentioned, labels have local significance only. Thus, the LSR removes the incoming label from the packet and attaches the matching outgoing label before forwarding the packet. The ingress-edge LSR determines the FEC for each incoming unlabeled packet and, on the basis of the FEC, assigns the packet to a particular LSP, attaches the corresponding label, and forwards the packet.

Label Stacking

One of the most powerful features of MPLS is *label stacking*. A labeled packet may carry many labels, organized as a last-in-first-out stack. Processing is always based on the top label. At any LSR, a label may be added to the stack (push operation) or removed from the stack (pop operation). Label stacking allows the aggregation of LSPs into a single LSP for a portion of the route through a network, creating a *tunnel*. At the beginning of the tunnel, an LSR assigns the same label to packets from a number of LSPs by pushing the label onto the stack of each packet. At the end of the tunnel, another LSR pops the top element from the label stack, revealing the inner label. This is similar to ATM, which has one level of stacking (virtual channels inside virtual paths), but MPLS supports unlimited stacking.

Label stacking provides considerable flexibility. An enterprise could establish MPLS-enabled networks at various sites and establish numerous LSPs at each site. The enterprise could then use label stacking to aggregate multiple flows of its own traffic before handing it to an access provider. The access provider could aggregate traffic from multiple enterprises before handing it to a larger service provider. Service providers could aggregate many LSPs into a relatively small number of tunnels between points of presence. Fewer tunnels means smaller tables, making it easier for a provider to scale the network core.

Figure 3: MPLS Label Format



Label Format and Placement

An MPLS label is a 32-bit field consisting of the following elements (Figure 3):

- *Label value*: locally significant 20-bit label
- *Exp*: 3 bits reserved for experimental use; for example, these bits could communicate DS information or PHB guidance
- *S*: set to one for the oldest entry in the stack, and zero for all other entries
- *Time To Live* (TTL): 8 bits used to encode a hop count, or time to live, value

Time-to-Live Processing

A key field in the IP packet header is the TTL field (IPv4), or Hop Limit (IPv6). In an ordinary IP-based internet, this field is decremented at each router and the packet is dropped if the count falls to zero. This is done to avoid looping or having the packet remain too long in the internet because of faulty routing. Because an LSR does not examine the IP header, the TTL field is included in the label so that the TTL function is still supported. The rules for processing the TTL field in the label are as follows:

1. When an IP packet arrives at an ingress edge LSR of an MPLS domain, a single label stack entry is added to the packet. The TTL value of this label stack entry is set to the value of the IP TTL value. If the IP TTL field needs to be decremented, as part of the IP processing, it is assumed that this has already been done.

When an MPLS packet arrives at an internal LSR of an MPLS domain, the TTL value in the top label stack entry is decremented.

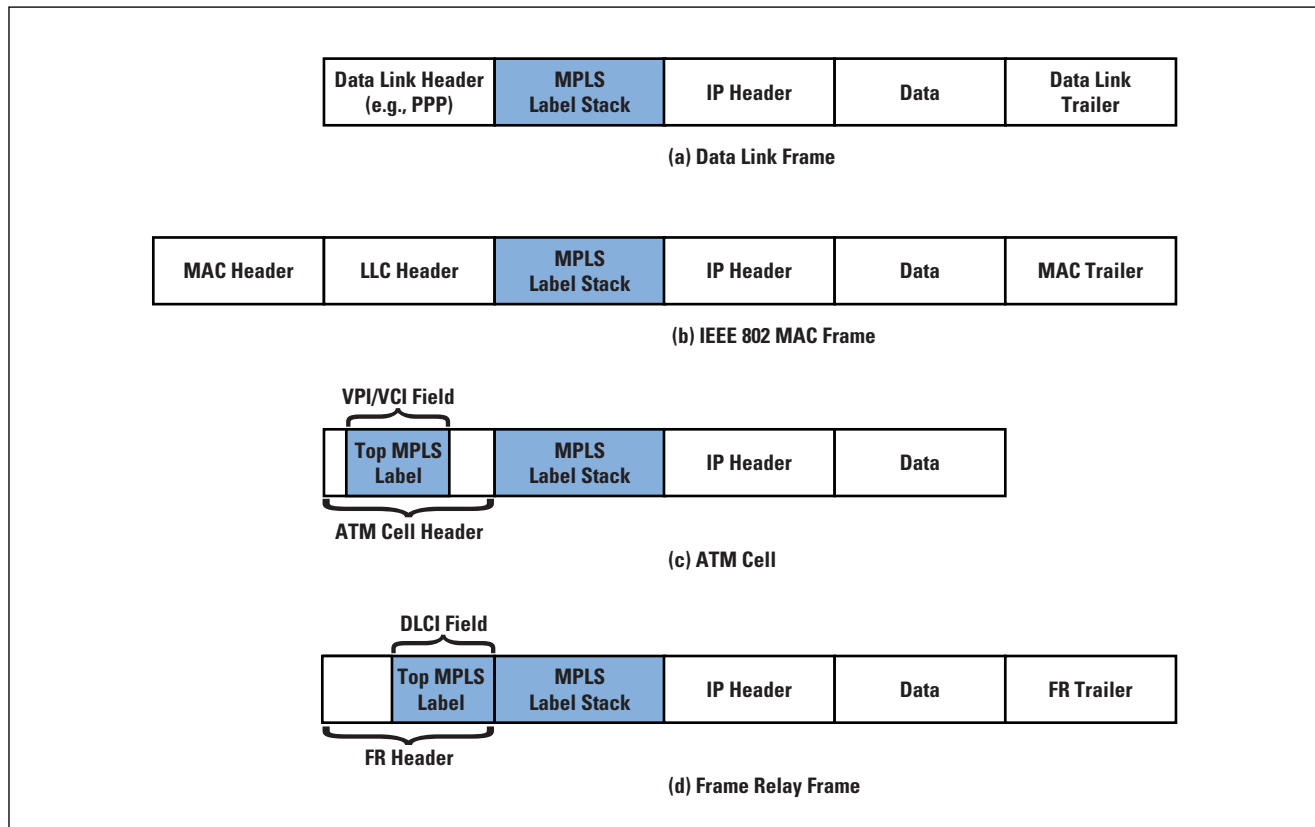
Then:

- (a) If this value is zero, the MPLS packet is not forwarded. Depending on the label value in the label stack entry, the packet may be simply discarded, or it may be passed to the appropriate “ordinary” network layer for error processing (for example, for the generation of an *Internet Control Message Protocol* [ICMP] error message).
 - (b) If this value is positive, it is placed in the TTL field of the top label stack entry for the outgoing MPLS packet, and the packet is forwarded. The outgoing TTL value is a function solely of the incoming TTL value, and is independent of whether any labels are pushed or popped before forwarding. There is no significance to the value of the TTL field in any label stack entry that is not at the top of the stack.
2. When an MPLS packet arrives at an egress edge LSR of an MPLS domain, the TTL value in the single label stack entry is decremented and the label is popped, resulting in an empty label stack. Then:
- (a) If this value is zero, the IP packet is not forwarded. Depending on the label value in the label stack entry, the packet may be simply discarded, or it may be passed to the appropriate “ordinary” network layer for error processing.
 - (b) If this value is positive, it is placed in the TTL field of the IP header, and the IP packet is forwarded using ordinary IP routing. Note that the IP header checksum must be modified prior to forwarding.

Label Stack

The label stack entries appear after the data link layer headers, but before any network layer headers. The top of the label stack appears earliest in the packet (closest to the network layer header), and the bottom appears latest (closest to the data link header). The network layer packet immediately follows the label stack entry that has the *S* bit set. In a data link frame, such as for the *Point-to-Point Protocol* (PPP), the label stack appears between the IP header and the data link header (Figure 4a). For an IEEE 802 frame, the label stack appears between the IP header and the *Logical Link Control* (LLC) header (Figure 4b).

Figure 4: Position of MPLS Label



If MPLS is used over a connection-oriented network service, a slightly different approach may be taken, as shown in Figure 4c and d. For ATM cells, the label value in the topmost label is placed in the *Virtual Path/Channel Identifier* (VPI/VCI) field in the ATM cell header. The entire top label remains at the top of the label stack, which is inserted between the cell header and the IP header. Placing the label value in the ATM cell header facilitates switching by an ATM switch, which would, as usual, need to look only at the cell header. Similarly, the topmost label value can be placed in the *Data Link Connection Identifier* (DLCI) field of a Frame Relay header. Note that in both these cases, the TTL field is not visible to the switch and so is not decremented. The reader should consult the MPLS specifications for the details of the way this situation is handled.

FECs, LSPs, and Labels

To understand MPLS, it is necessary to understand the operational relationship among FECs, LSPs, and labels. The specifications covering all the ramifications of this relationship are lengthy. In the remainder of this section, we provide a summary.

The essence of MPLS functionality is that traffic is grouped into FECs. The traffic in an FEC transits an MPLS domain along an LSP. Individual packets in an FEC are uniquely identified as being part of a given FEC by means of a *locally significant label*.

At each LSR, each labeled packet is forwarded on the basis of its label value, with the LSR replacing the incoming label value with an outgoing label value.

The overall scheme described in the previous paragraph imposes numerous requirements. Specifically:

1. Traffic must be assigned to a particular FEC.
2. A routing protocol is needed to determine the topology and current conditions in the domain so that a particular LSP can be assigned to an FEC. The routing protocol must be able to gather and use information to support the QoS requirements of the FEC.
3. Individual LSRs must become aware of the LSP for a given FEC, must assign an incoming label to the LSP, and must communicate that label to any other LSR that may send it packets for this FEC.

The first requirement is outside the scope of the MPLS specifications. The assignment needs to be done either by manual configuration, by means of some signaling protocol, or by an analysis of incoming packets at ingress LSRs. Before looking at the other two requirements, let us consider the topology of LSPs. We can classify these in the following manner:

- *Unique ingress and egress LSR*: In this case a single path through the MPLS domain is needed.
- *Unique egress LSR, multiple ingress LSRs*: If traffic assigned to a single FEC can arise from different sources that enter the network at different ingress LSRs, then this situation occurs. An example is an enterprise intranet at a single location but with access to an MPLS domain through multiple MPLS ingress LSRs. This situation would call for multiple paths through the MPLS domain, probably sharing a final few hops.
- *Multiple egress LSRs for unicast traffic*: RFC 3031 states that most commonly, a packet is assigned to a FEC based (completely or partially) on its network layer destination address. If not, then it is possible that the FEC would require paths to multiple distinct egress LSRs. However, more likely, there would be a cluster of destination networks, all of which are reached via the same MPLS egress LSR.
- *Multicast*: RFC 3031 lists multicast as a subject for further study.

Route Selection

Route selection refers to the selection of an LSP for a particular FEC. The MPLS architecture supports two options: hop-by-hop routing and explicit routing.

With *hop-by-hop routing*, each LSR independently chooses the next hop for each FEC. The RFC implies that this option makes use of an ordinary routing protocol, such as OSPF.

This option provides some of the advantages of MPLS, including rapid switching by labels, the ability to use label stacking, and differential treatment of packets from different FECs following the same route. However, because of the limited use of performance metrics in typical routing protocols, hop-by-hop routing does not readily support traffic engineering or policy routing (defining routes based on some policy related to QoS, security, or some other consideration).

With *explicit routing*, a single LSR, usually the ingress or egress LSR, specifies some or all of the LSRs in the LSP for a given FEC. For strict explicit routing, an LSR specifies all of the LSRs on an LSP. For loose explicit routing, only some of the LSRs are specified. Explicit routing provides all the benefits of MPLS, including the ability to do traffic engineering and policy routing.

Explicit routes can be selected by configuration, that is, set up ahead of time, or dynamically. Dynamic explicit routing would provide the best scope for traffic engineering. For dynamic explicit routing, the LSR setting up the LSP would need information about the topology of the MPLS domain as well as QoS-related information about that domain. An MPLS traffic engineering specification^[2] suggests that the QoS-related information falls into two categories:

- A set of attributes associated with an FEC or a collection of similar FECs that collectively specify their behavioral characteristics
- A set of attributes associated with resources (nodes, links) that constrain the placement of LSPs through them

A routing algorithm that accounts for the traffic requirements of various flows and the resources available along various hops and through various nodes is referred to as a *constraint-based routing algorithm*. In essence, a network that uses a constraint-based routing algorithm is aware of current utilization, existing capacity, and committed services at all times. Traditional routing algorithms, such as OSPF and the *Border Gateway Protocol* (BGP), do not employ a sufficient array of cost metrics in their algorithms to qualify as constraint-based.

Furthermore, for any given route calculation, only a single cost metric (for instance, number of hops, delay) can be used. For MPLS, it is necessary either to augment an existing routing protocol or to deploy a new one. For example, an enhanced version of OSPF has been defined^[1] that provides at least some of the support required for MPLS. Examples of metrics that would be useful to constraint-based routing include the following:

- Maximum link data rate
- Current capacity reservation
- Packet loss ratio
- Link propagation delay

Label Distribution

Route selection consists of defining an LSP for an FEC. A separate function is the actual setting up of the LSP. For this purpose, each LSR on the LSP must:

1. Assign a label to the LSP to be used to recognize incoming packets that belong to the corresponding FEC.
2. Inform all potential upstream nodes (nodes that will send packets for this FEC to this LSR) of the label assigned by this LSR to this FEC, so that these nodes can properly label packets to be sent to this LSR.
3. Learn the next hop for this LSP and learn the label that the downstream node (LSR that is the next hop) has assigned to this FEC. This process will enable this LSR to map an incoming label to an outgoing label.

The first item in the preceding list is a local function. Items 2 and 3 must be done either by manual configuration or by using some sort of label distribution protocol. Thus, the essence of a label distribution protocol is that it enables one LSR to inform others of the label/FEC bindings it has made. In addition, a label distribution protocol enables two LSRs to learn each other's MPLS capabilities. The MPLS architecture does not assume a single label distribution protocol but allows for multiple such protocols. Specifically, RFC 3031 refers to a new label distribution protocol and to enhancements to existing protocols, such as RSVP and BGP, to serve the purpose.

The relationship between label distribution and route selection is complex. It is best to look at in the context of the two types of route selection.

With hop-by-hop route selection, no specific attention is paid to traffic engineering or policy routing concerns, as we have seen. In such a case, an ordinary routing protocol such as OSPF is used to determine the next hop by each LSR. A relatively straightforward label distribution protocol can operate using the routing protocol to design routes.

With explicit route selection, a more sophisticated routing algorithm must be implemented, one that does not employ a single metric to design a route. In this case, a label distribution protocol could make use of a separate route selection protocol, such as an enhanced OSPF, or incorporate a routing algorithm into a more complex label distribution protocol.

References

The two most important defining documents for MPLS are [5] and [6]. Reference [3] provides a thorough treatment of MPLS; [8] covers not only MPLS but other Internet QoS concepts; it includes an excellent chapter on MPLS traffic engineering. Reference [7] includes a concise overview of the MPLS architecture and describes the various proprietary efforts that preceded MPLS.

- [1] Apostolopoulos, G., et al., “QoS Routing Mechanisms and OSPF Extensions,” RFC 2676, August 1999.
- [2] Awduche, D., et al. “Requirements for Traffic Engineering over MPLS,” RFC 2702, September 1999.
- [3] Black, U., *MPLS and Label Switching Networks*, ISBN 0130158232, Prentice Hall, 2001.
- [4] Redford, R., “Enabling Business IP Services with Multiprotocol Label Switching,” Cisco White Paper, July 2000 (www.cisco.com).
- [5] Rosen, E., et al. “Multiprotocol Label Switching Architecture,” RFC 3031, January 2001.
- [6] Rosen, E., et al. “MPLS Label Stack Encoding,” RFC 3032, January 2001.
- [7] Viswanathan, A., et al., “Evolution of Multiprotocol Label Switching,” *IEEE Communications Magazine*, May 1998.
- [8] Wang, Z., *Internet QoS: Architectures and Mechanisms for Quality of Service*, ISBN 1558606084, Morgan Kaufmann, 2001.

Useful Web Sites

- *MPLS Forum*: An industry forum to promote MPLS:
<http://www.mplsforum.org/>
- *MPLS Resource Center*: Clearinghouse for information on MPLS:
<http://www.mplsrc.com/>
- *MPLS Working Group*: Chartered by IETF to develop standards related to MPLS. The Web site includes all relevant RFCs and Internet Drafts:
<http://www.ietf.org/html.charters/mpls-charter.html>

WILLIAM STALLINGS is a consultant, lecturer, and author of over a dozen books on data communications and computer networking. He also maintains a computer science resource site for CS students and professionals at WilliamStallings.com/StudentSupport.html. He has a PhD in computer science from M.I.T. His latest book is *Wireless Communications and Networks* (Prentice Hall, 2001). His home in cyberspace is WilliamStallings.com and he can be reached at ws@shore.net

A Unique, Authoritative Root for the DNS

by M. Stuart Lynn, ICANN

The following *Internet Coordination Policy* (ICP) is being posted for the information of the Internet community by the *Internet Corporation for Assigned Names and Numbers* (ICANN) and is a statement of policy currently followed in administering the authoritative root of the Domain Name System. Comments on this article are welcome and should be directed to comments@icann.org

Abstract

This article reaffirms ICANN's commitment to a single, authoritative public root for the Internet *Domain Name System* (DNS) and to the management of that unique root in the public interest according to policies developed through community processes. This commitment is founded on the technical and other advice of the community and is embodied in existing ICANN policy.

The DNS is intended to provide a convenient means of referring to sites available on the Internet. By offering users an easy-to-use and reliable means of unambiguously referring to Web sites, e-mail servers, and the Internet's many other services, the DNS has helped the Internet achieve its promise as a global communications medium for commerce, research, education, and cultural and other expressive activities.

The DNS is a globally distributed database of domain name (and other) information. One of its core design goals is that it reliably provides the same answers to the same queries from any source on the public Internet, thereby supporting predictable routing of Internet communications. Achievement of that design goal requires a globally unique public name space derived from a single, globally unique DNS *root*.

Although the Internet allows a high degree of decentralized activities, coordination of the assignment function by a single authority is necessary where unique parameter values are technically required. Because of the uniqueness requirement, the content and operation of the DNS root must be coordinated by a central entity.

Where central coordination is necessary, it should be performed by an organization dedicated to serving the public interest and that acts according to policies developed through processes that are developed through the participation of affected stakeholders. Traditionally, the responsibility for performing the central coordinating functions of the global Internet for the public good, including management of the unique public DNS root, has been carried out by the *Internet Assigned Numbers Authority* (IANA)^[12]. ICANN's core mission is to continue the work of the IANA in a more formalized and globally representative framework, to ensure the views of all the Internet's stakeholders are taken into account in carrying out this public trust.

Over the past several years, some private organizations have established DNS roots as alternates to the authoritative root. Some uses of these alternate roots do not jeopardize the stability of the DNS. For example, some are purely private roots operating inside institutions and are carefully insulated from the DNS. Others are purely experimental in the best traditions of the Internet and are carefully managed so as not to interfere with the operation of the DNS. These both operate within community-established norms.

Frequently, however, these alternate roots have been established to support top-level or pseudo-top-level domain name registries that are operated for profit. Yet other alternate roots have been established by certain individuals to protest the policies developed by the broader community processes for management of the authoritative root, or to express their disinterest in participating in those processes. These alternate roots have not been launched through any ICANN consensus processes, so they have not been entered into the authoritative root managed by the IANA or ICANN.

These alternate roots typically substitute insular concerns in place of the community-based processes that govern the management of the authoritative root. Their operators decide to include particular top-level domains in these alternate roots that have not been subjected to the tests of community support and conformance with consensus processes—coordinated by ICANN—that would allow their inclusion in the authoritative root. These decisions of the alternate-root operators have been made without any apparent regard for the fundamental public-interest concern of Internet stability. The widespread use of active domain names in these alternate roots could in fact impair the uniqueness of the authoritative name-resolution mechanism and hence the stability of the DNS.

ICANN's mandate to preserve stability of the DNS requires that it avoid encouraging the proliferation of these alternate roots that could cause conflicts and instability. This means that ICANN continues to adhere to community-based processes in its decisions regarding the content of the authoritative root. Within its current policy framework, ICANN can give no preference to those who choose to work outside of these processes and outside of the policies engendered by this public trust.

None of this precludes experimentation done in a manner that does not threaten the stability of name resolution in the authoritative DNS. Responsible experimentation is essential to the vitality of the Internet. Nor does it preclude the ultimate introduction of new architectures that may ultimately obviate the need for a unique, authoritative root. But the translation of experiments into production and the introduction of new architectures require community-based approaches, and are not compatible with individual efforts to gain proprietary advantage.

The Technical Need for a Single Authoritative Root

The DNS was originally deployed in the mid-1980s^[13] as an improved means of mapping easy-to-remember names (i.e., **example.com**) to the IP addresses (i.e., **128.9.176.32**) by which packets are routed on the Internet. It is a distributed database that holds this mapping information (as well as various other types of technical information regarding computers on the Internet) in *resource records*. The DNS provides these resource records in response to queries it receives from programs called *resolvers* on individual computers throughout the Internet. The resolvers translate domain names into the corresponding IP addresses.

From the inception of the DNS, its most fundamental design goal has been to provide the same answers to the same queries issued from any place on the Internet. As stated in RFC 1034, the basic specification of the DNS's "Concepts and Facilities,"^[16] "The primary (design) goal is a consistent name space which will be used for referring to resources." And as reiterated in RFC 2535, "Domain Name System Security Extensions,"^[15] "It is part of the design philosophy of the DNS that the data in it is public and that the DNS gives the same answers to all inquirers."

The DNS is hierarchical. By design, the hierarchy begins with a group of *root nameservers* (often called simply *root servers*), which are specially-designated computers operated under common coordination that provide information about which other computers are authoritative regarding the top-level domains in the DNS naming structure. These set of root servers house the *authoritative root*. Thus, a resolver seeking information concerning a domain name such as **www.example.com** obtains one of the root servers' resource records about **.com**, which tells the resolver which computers have authoritative information about names within the **.com** top-level domain. The resolver then queries one of those authoritative **.com** nameservers about **example.com**, to locate the nameservers for **example.com**. A query is then made to one of those nameservers obtain the IP address of the computer designated by the name **www.example.com**.

The principal advantage of this hierarchical structure is that it allows different parts of the naming database to be maintained by different entities. According to the DNS's design, each domain was intended to be administered by a single entity.^[19]

When the DNS was deployed in the mid-1980s, a set of root nameservers was designated and several top-level domains were established. These root nameservers (there are now 13 of them distributed around the world) are intended to provide authoritative information about which nameservers hold the naming information for each of the top-level domains. Since the authoritative root nameservers operate at the top of the hierarchy, resolvers find them by referring to IP addresses pre-stored at local computers throughout the Internet.

Over the past several years, some groups have established alternate root nameservers on the public Internet that distribute different information than the information distributed by the authoritative root nameservers. These groups then seek to persuade ISPs and Internet users to replace the pre-stored IP addresses of the authoritative root nameservers with those of their alternate servers. For a variety of reasons, these alternate roots have not to date achieved a significant level of usage on the public Internet.

Fortunately, the rare usage of alternate roots has thus far limited their practical effect on the Internet. If these alternate roots were to become prevalent, however, they would have the potential for seriously disrupting the reliable functioning of the DNS. Some of the consequences include:

- *Providing the Wrong Location:* The presence of alternate public DNS roots can result in different answers being given to the same DNS query issued from different computers on the Internet, depending on whether the inquiring computer is programmed to access the authoritative root or a particular one of the alternate roots (or more precisely a domain-name resolver associated with one or the other of these). The fundamental DNS design goal of providing consistent answers to DNS queries is therefore frustrated.^[1]
- *Reaching the Wrong Computer:* The main consequence of such inconsistent data is that the same domain name can identify different computers depending on where the name is used. Put another way, *Uniform Resource Locators* (URLs) are no longer uniform. Thus, typing in a Web site address at two different computers configured to reference different roots can result in reaching different Web sites—a particularly disturbing possibility if, for example, money is to change hands or privacy or security concerns are violated. Similarly, the same piece of e-mail sent to the same address from the two computers can be directed to different recipients. The return of inconsistent DNS data defeats the globally consistent resolution of domain names that is vital to the Internet achieving its promise as a universal communications and applications medium for commerce, research, education, cultural exchange, expressive activities, and other uses.
- *Consequences Unpredictable to Most Users:* The set of DNS answers that will be received (from the authoritative root or one of the several alternate roots) is not predictable by most end users. Most users on the Internet employ a local DNS resolver that is configured by another person. Few users are likely to appreciate the significance of the resolver's DNS configuration; even fewer are likely to have detailed knowledge of that configuration. As the number of users on the Internet has grown, the proportion of users knowledgeable about technical concepts such as DNS resolvers and root servers has diminished. Yet these non-technical users are precisely those for whom the Internet in general—and the DNS in particular—hold the greatest potential benefits.

- *Intermediate Hosts Add to Confusion:* Moreover, some Internet services depend on the actions of DNS resolvers employed by intermediate hosts. Alternate roots introduce the possibility that the DNS answer obtained by the intermediate host alters the character of the service in an unexpected way. A similar phenomenon can occur where one user sends another a reference to a URL, such as an e-mail reply address or a link on a Web site. If the recipient of an e-mail or the visitor to the Web site is using a computer that employs a different DNS root than intended by the sender of the e-mail or the designer of the Web site, unexpected results are likely to occur. For example, the e-mail could end up with the wrong person.
- *Cache Poisoning:* Alternate roots also introduce the possibility of misdirected Internet activities due to the phenomenon known as cache poisoning. For performance reasons, the DNS design calls for resource records to be passed around among the nameservers on the Internet, so that a resolver can obtain quicker access to a local copy of the resource record. Because the DNS assumes a single-root system, resource records are not marked to distinguish them according to the root from which they emanate. Thus, the presence of alternate roots introduces the possibility that Internet activities by those intending to use the authoritative root could be misdirected by a stray resource record emanating from an alternate root. Indeed, some malicious hacking attacks have been based on this principle, prompting the *Internet Engineering Task Force* (IETF) to propose a series of not-yet-fully-implemented improvements known as *DNS-Security* or *DNSSEC*.

(It should be noted that the original design of the DNS provided a way to operate alternate roots in a way that does not imperil stability. See “Experimentation” below for details.)

These potentially destructive effects of alternate roots have long been accepted by the vast majority of Internet engineers. Despite this broad-based recognition, some have sought to justify the alternate roots by downplaying these effects. In response, and to document what it referred to as “some of the problems inherent in a family of recurring technically naive proposals,” in May 2000 the *Internet Architecture Board* (IAB)^[14] issued RFC 2826, entitled “IAB Technical Comment on the Unique DNS Root.” The IAB summarized its comments (in relevant part) as follows:

“Summary: To remain a global network, the Internet requires the existence of a globally unique public name space. The DNS name space is a hierarchical name space derived from a single, globally unique root. This is a technical constraint inherent in the design of the DNS. Therefore it is not technically feasible for there to be more than one root in the public DNS. That one root must be supported by a set of coordinated root servers administered by a unique naming authority.

“Put simply, deploying multiple public DNS roots would raise a very strong possibility that users of different ISPs who click on the same link on a Web page could end up at different destinations, against the will of the Web page designers.”

For some concrete examples of potential failures and instabilities that would likely result from alternate roots prevalently used on the public Internet, see the draft “Alt-Roots, Alt-TLDs.”^[17]

In the face of the destabilizing consequences of alternate roots, as articulated by the IAB and others, ICANN’s prime directive of preserving the stability of the Internet and DNS requires an unwavering commitment to promote the continued prevalence of a single authoritative root for the public DNS. Any other course of action by ICANN would be irresponsible.

The Public Trust in Coordinated Assignment Functions

The Internet’s proper operation requires assignment of unique values to various identifiers for different computers or services on the Internet. To be effective, these assigned values must be made broadly available and their significance must be respected by the many people responsible for the Internet’s operation. For example, every computer on the public Internet is assigned a unique IP address; this address is made known to routers throughout the Internet to cause TCP/IP packets with that destination address to be routed to the intended computer. Without common agreement to respect the assignment, the Internet would not reliably route communications to their intended destinations.

Beginnings to 1998: Central Coordination as a Public Trust

From the very beginnings of the Internet, the technical community has recognized the need for central coordination of the unique assignment of the values of identifiers. The IANA, now operated by ICANN was created to fill this need; it now makes assignments of unique values for approximately 120 different identifier types. This responsibility has always been understood to be a public trust, and the IANA long ago adopted the motto: “Dedicated to preserving the central coordinating functions of the global Internet for the public good.”

The most commonly known of the Internet’s uniquely assigned identifiers, of course, are domain names. From the time the DNS was deployed, the Internet community made the IANA “responsible for the overall coordination and management of the Domain Name System (DNS), and especially the delegation of portions of the name space called top-level domains.”^[18] As in its other assignment responsibilities, the IANA’s role is to act in the public interest, neutrally, and without proprietary motives.

Competition as a Value Guiding the Internet's Technical Management

In the Internet's early years, with limited exceptions day-to-day registration activities for domain names were done by a single company (first SRI International and later Network Solutions) under the IANA's guidance.

By the mid-1990s, however, the growth and increasing commercialization of the Internet led the U.S. Government's Green^[2] and White^[3] Papers to note the emergence of "widespread dissatisfaction about the absence of competition in domain name registration." This dissatisfaction prompted the Green and White Papers to include the promotion of competition in registration services as one of the four values (stability; competition; private, bottom-up coordination; and representation) that should guide the Internet's technical management. Both documents made clear that, of these four values, preservation of stability was to be paramount.

Building on the IANA model of a non-profit entity carrying the public trust to perform the vital central coordination functions, the U.S. Government reconciled the need to ensure Internet stability with the desire to introduce competitive domain-name registration services as follows:

"In keeping with these principles, we divide the name and number functions into two groups, those that can be moved to a competitive system and those that should be coordinated. We then suggest the creation of a representative, not-for-profit corporation to manage the coordinated functions according to widely accepted objective criteria. We then suggest the steps necessary to move to competitive markets in those areas that can be market driven."^[4]

This dichotomy recognizes that the Internet is, after all, a network (albeit a network of networks), and networks require coordination among their participants to operate in a stable and efficient manner. It also reflects the phenomenal success of the Internet's tradition of cooperatively developed open and non-proprietary standards. Those standards have provided an environment of highly interoperable systems that has allowed competition and innovation to flourish.

ICANN Assumes the Public Trust

After public comment on the Green Paper, the United States Government issued the White Paper, which laid out the basic charter on which ICANN was founded and continues to operate. The White Paper re-emphasized the prime directive of stability and, to that end, the need to avoid creation of alternate roots:

"The introduction of a new management system should not disrupt current operations or create competing root systems. During the transition and thereafter, the stability of the Internet should be the first priority of any DNS management system."^[5]

The United States Government then invited the Internet community to form a not-for-profit corporation to perform the “coordinated functions” that should be handled as a matter of public trust, rather than according to a competitive regime that would not be conducive to stability. Among the “coordinated functions” were management of the root-server system and decisions to introduce new TLDs:

“Similarly, coordination of the root server network is necessary if the whole system is to work smoothly. While day-to-day operational tasks, such as the actual operation and maintenance of the Internet root servers, can be dispersed, *overall policy guidance and control of the TLDs and the Internet root server system should be vested in a single organization* that is representative of Internet users around the globe.

“Further, changes made in the administration or the number of gTLDs contained in the authoritative root system will have considerable impact on Internet users throughout the world. In order to promote continuity and reasonable predictability in functions related to the root zone, the *development of policies for the addition, allocation, and management of gTLDs and the establishment of domain name registries and domain name registrars to host gTLDs should be coordinated.*”^[6]

In response to this invitation for the formation of a non-profit, Internet-community-based organization, ICANN was established in 1998. ICANN was subsequently selected by the United States Government from among several proposals submitted precisely because it was open, consensus-based, and rooted in the Internet community. The establishment of ICANN had followed extensive dialogs among different constituencies of the Internet community to ensure that ICANN could be responsive to the needs of these various constituencies.

ICANN, among its other responsibilities, now acts as the coordinator for operation of the authoritative root-server system and the policy forum for decisions about the policies governing what TLDs are to be included in the authoritative DNS root.^[7]

In linking the formation of ICANN to the global Internet community, the White Paper established a public trust that required that the DNS be administered in the public interest as the unique-rooted,^[8] authoritative database for domain names that provides a stable addressing system for use by the global Internet community. This commitment to a unique and authoritative root is a key part of the broader public trust—to carry out the Internet’s central coordination functions for the public good—that is ICANN’s reason for existence.

The Public Trust and the Introduction of New TLDs

It is essential that the centrally coordinated functions be performed in the public interest, not out of proprietary or otherwise self-interested motives. For this reason, ICANN was founded as a not-for-profit public-benefit organization, accountable to the Internet community. Longstanding Internet principles also require that the policies guiding the coordinated functions be established openly based on community deliberation and input. For these reasons ICANN's structure is representative of the geographic and functional diversity of the Internet, and relies to the extent possible on private-sector, bottom-up methods.

As the White Paper emphasized, the decisions about the introduction of new TLDs are appropriately done within this open, non-proprietary, and broadly representative framework, rather than by individuals or entities not accountable to the community and that ordinarily act for their own proprietary motives:

“As Internet names increasingly have commercial value, the decision to add new top-level domains cannot be made on an ad hoc basis by entities or individuals that are not formally accountable to the Internet community.”^[9]

Within the framework of its commitment to a unique root system and to the stability of the Internet, last year ICANN launched a process for carefully introducing several new generic TLDs to the DNS. This introduction was fashioned as a proof of concept of the technical and business feasibility of introducing more TLDs into the DNS. Proceeding with an initial proof of concept was in response to the advice of ICANN's *Protocol Supporting Organization* (PSO) and its *Domain Name Supporting Organization* (DNSO) to proceed cautiously and in an orderly fashion. The PSO and the DNSO represent the consensus views of the technical and the user/business/other institutional communities, respectively. Generic TLDs had not been introduced for many years, and there were and still are serious questions as to what the effect of introducing new TLDs will be on the stability and reliability of the DNS; and many questions about what should be the appropriate contractual and business context.

In response to an issued RFP, forty-seven institutions and groups submitted proposals for the establishment of new TLDs. They chose to work within the community-based ICANN process, even though they knew that only a “limited number” of TLDs would be selected—at least in the first round. In fact, seven were selected, and, following a methodology which allowed for considerable community input, contracts have or will shortly be signed with these initial seven. ICANN looks forward to the successful introduction of these new TLDs and will work with the community to monitor their performance so that a community decision can be made on moving forward with the introduction of more TLDs, should this be the conclusion of the proof of concept.

Outside the Process

Some private organizations have established DNS roots as alternates to the authoritative root. Some uses of these alternate roots do not jeopardize the stability of the DNS. For example, many are purely private roots operating inside institutions and are carefully insulated from the DNS. Others are purely experimental in the best traditions of the Internet and are carefully managed so as not to interfere with the operation of the DNS. These both operate within community-established norms.

Frequently, however, these alternate roots have been established to support top-level or pseudo-top-level domain name registries that are operated for profit. Yet other alternate roots have been established by certain individuals to protest the policies developed by the broader community processes for management of the authoritative root, or to express their disinterest in participating in those processes. These alternate roots have not been launched through any ICANN consensus processes, so they have not been entered into the authoritative root managed by the IANA or ICANN.

These alternate roots typically substitute insular concerns in place of the community-based processes that govern the management of the authoritative root. Their operators decide to include particular top-level domains in these alternate roots that have not been subjected to the tests of community support and conformance with consensus processes—coordinated by ICANN—that would allow their inclusion in the authoritative root. These decisions of the alternate root operators have been made with no apparent regard for the fundamental public-interest concern of Internet stability. The widespread introduction of active domain names into these alternate roots could in fact impair the uniqueness of the authoritative name resolution mechanism and hence the stability of the DNS.

In fact, some of the operators of these alternate roots state that stability is not an important attribute for the DNS. This thesis, for reasons already stated, is at fundamental variance with ICANN policy as embodied in its founding documents. Some of these operators and their supporters assert that their very presence in the marketplace gives them preferential right to TLDs to be authorized in the future by ICANN. They work under the philosophy that if they get there first with something that looks like a TLD and invite many registrants to participate, then ICANN will be required by their very presence and force of numbers to recognize in perpetuity these pseudo TLDs, inhibiting new TLDs with the same top-level name from being launched through the community's processes.

No current policy allows ICANN to grant such preferential rights. To do so would effectively yield ICANN's mandate to introduce new TLDs in an orderly manner in the public interest to those who would simply grab all the TLD names that seem to have any marketplace value, thus

circumventing the community-based processes that ICANN is required to follow. For ICANN to yield its mandate would be a violation of the public trust under which ICANN was created and under which it must operate. Were it to grant such preferential rights, ICANN would abandon this public trust, rooted in the community, to those who only act for their own benefit. Indeed, granting preferential rights could jeopardize the stability of the DNS, violating ICANN's fundamental mandate.

Alternate roots inherently endanger DNS stability—that is, they create the real risk of name resolvers being unable to determine to which numeric address a given name should point. This violates the fundamental design of the DNS and impairs the Internet's utility as a ubiquitous global communications medium. Some of these alternate systems also employ special technologies that—ingenious as they may be—may conflict with future generations of community-established Internet standards. Indeed, can there be any guarantee that these proprietary technologies can or will be adapted to future changes in Internet standards?

Experimentation

Experimentation has always been an essential component of the Internet's vitality. Working within the system does not preclude experimentation, including experimentation with alternate DNS roots. But these activities must be done responsibly, in a manner that does not disrupt the ongoing activities of others and that is managed according to experimental protocols.

DNS experiments should be encouraged. Experiments, however, almost by definition have certain characteristics to avoid harm: (a) they are clearly labeled as experiments, (b) it is well understood that these experiments may end without establishing any prior claims on future directions, (c) they are appropriately coordinated within a community-based framework (such as the IETF), and (d) the experimenters commit to adapt to consensus-based standards when they emerge through the ICANN and other community-based processes. This is very different from launching commercial enterprises that lull users into a sense of permanence without any sense of the foregoing obligations or contingencies.

Moreover, it is essential that experimental operations involving alternate DNS roots be conducted in a controlled manner, so that they do not adversely affect those who have not consented to participate in them. Given the design of the DNS, and particularly the intermediate-host and cache poisoning issues described earlier, special care must be taken to insulate the DNS from the alternate roots' effects. For example, alternate roots are commonly operated by large organizations within their private networks without harmful effects, since care is taken to prevent the flow of the alternate resource records onto the public Internet.

It should be noted that the original design of the DNS provides a facility for future extensions that accommodates the possibility of safely deploying multiple roots on the public Internet for experimental and other purposes. As noted in RFC 1034, the DNS includes a “class” tag on each resource record, which allows resource records of different classes to be distinguished even though they are commingled on the public Internet. For resource records within the authoritative root-server system, this class tag is set to “IN”; other values have been standardized for particular uses, including 255 possible values designated for “private use” that are particularly suited to experimentation.^[10]

As described in a recent proposal within the IETF,^[11] this “class” facility allows an alternate DNS namespace to be operated from different root servers in a manner that does not interfere with the stable operation of the existing authoritative root-server system. To take advantage of this facility, it should be noted, requires the use of client or applications software developed for the alternate namespace (presumably deployed after responsible testing), rather than the existing software that has been developed to interoperate with the authoritative root. Those who operate alternate roots for global commercial purposes, however, have not followed this course.

In an ever-evolving Internet, ultimately there may be better architectures for getting the job done where the need for a single, authoritative root will not be an issue. But that is not the case today. And the transition to such an architecture, should it emerge, would require community-based approaches. In the interim, responsible experimentation should be encouraged, but it should not be done in a manner that affects those who do not consent after being informed of the character of the experiment.

Conclusion

The success of the Internet and the guarantee of Internet stability rest on the cooperative activities of thousands, even millions, of people and institutions collaborating worldwide towards a common end. This extraordinary—even unprecedented—community effort has served to impel the incredible growth of the Internet. Many of these people and institutions compete intensely among themselves yet agree to do so within a common framework for the overall public good. Their collective efforts provide a policy framework for technical and entrepreneurial innovation, and the advancement of economic, social, and educational goals.

Most members of the global community and most institutions with which they are associated recognize that it is in their best long-term interests to work within these community-based processes, even if that means foregoing short-term advantages to particular individuals or groups. The over-arching principles outlined in this document override exclusive and narrowly focused self-interest.

Community-based policy development is not perfect. It may proceed slower than some would wish. The introduction of new TLDs has proceeded at deliberate speeds. Impatience in the context of Internet timescales is perfectly understandable. The outcome of orderly processes based on the wishes of the community, however, is assurance that the Internet will continue to function in a stable and holistic manner that benefits the global community, and not become captured by the self-interests of the few. That, in the minds of most, is a price worth paying.

ICANN—in deference to its public trust—will continue to collaborate with these citizens of the Internet community to advance the notions of a unique root system as a prerequisite to Internet stability, and to ensure that community-based policies take precedence. ICANN encourages responsible experimentation designed to further advance the Internet as a useful, stable, and accessible medium for the public good.

References

- [1] Ironically, to avoid name conflicts in a multi-root system, a single-root system would need to be created—adding a higher level to the hierarchy.
- [2] “Improvement of Technical Management of Internet Names and Addresses,” (Green Paper), 63 *Federal Register* 8825, 8827 (20 February, 1998).
- [3] “Management of Internet Names and Addresses,” (White Paper), 63 *Federal Register* 31741, 31742 (10 June, 1998).
- [4] Green Paper, 63 *Federal Register* at 8827.
- [5] White Paper, 63 *Federal Register* at 31749. The Green and White Papers both made additional references to the need for a single authoritative root system. For example, in response to comments received from the Green Paper, the White Paper notes:

“In the absence of an authoritative root system, the potential for name collisions among competing sources for the same domain name could undermine the smooth functioning and stability of the Internet.”
- [6] White Paper, 63 *Federal Register* at 31749 (emphasis added).
- [7] ICANN’s corporate charter emphasizes its role in overseeing operation of the unique DNS root:

“... the Corporation shall ... pursue the charitable and public purposes ... of promoting the global public interest in the operational stability of the Internet by ... (iv) overseeing operation of the authoritative Internet DNS root server system ...”

ICANN Articles of Incorporation, para. 3. The phrase “the authoritative Internet DNS root server system” is decidedly in the *singular*.

See: <http://www.icann.org/general/articles.htm>

- [8] The Memorandum of Understanding between the United States Government and ICANN that governs the transfer of responsibilities from the U.S. Department of Commerce to ICANN also makes reference to the authoritative root in the singular, not in the plural:

“In the DNS Project, the parties will jointly design, develop, and test the mechanisms, methods, and procedures to carry out the following DNS management functions: ...

“b. Oversight of the operation of the authoritative root server system;

“c. Oversight of the policy for determining the circumstances under which new top level domains would be added to the root system ... ”

See also: www.icann.org/general/icann-mou-25nov98.htm

- [9] White Paper, 63 *Federal Register* at 31742.
- [10] Eastlake, D., Brunner-Williams, E., Manning, B., “Domain Name System (DNS) IANA Considerations,” section 3.2, RFC 2929, September, 2000.
- [11] Klensin, J., “Internationalizing the DNS—A New Class,” Internet Draft, work in progress, December, 2000.
- [12] Internet Assigned Numbers Authority (IANA). See www.iana.org
- [13] Postel, J., “Domain Name System Implementation Schedule—Revised,” RFC 921, October 1984.
- [14] Internet Architecture Board (IAB). See <http://www.iab.org>
- [15] Eastlake, D., “Domain Name System Security Extensions,” RFC 2535, March 1999.
- [16] Mockapetris, P., “Domain Names—Concepts and Facilities,” RFC 1034, November 1987.
- [17] <http://www.icann.org/stockholm/draft-crispin-alt-roots-tlds-00.txt>
- [18] Postel, J., “Domain Name System Structure and Delegation,” RFC 1591, March 1994.
- [19] Postel, J., and Reynolds, J., “Domain Requirements,” RFC 920, October 1984.

Dr. M. STUART LYNN is President & CEO of The Internet Corporation for Assigned Names and Numbers (ICANN). Dr. Lynn has had a distinguished career in computing and information technology that dates back almost four decades. His most recent position until his retirement in 1999 was as Associate Vice President for Information Resources and Communications for the University of California Office of the President where he served as chief information officer for the combined University of California system. Dr. Lynn also served as President and Chairman of the Board of the Corporation for Education Network Initiatives in California (CENIC). Dr. Lynn has also held positions at Cornell University, UC Berkeley, Rice University, Baylor College of Medicine, IBM and Chevron. Over the course of his career, he has been active in several professional organizations including the Association for Computing Machinery (ACM) and the American Federation of Information Processing Societies. In 1994, he was elected a Fellow of the ACM. In addition, he has served on numerous boards of directors, advisory committees and as a consultant to academia, government and industry. Dr. Lynn holds a M.A. and Ph.D. in Mathematics from the University of California at Los Angeles and a B.A. and M.A. in Mathematics from Oxford University.

E-mail: lynn@icann.org

Book Review

Web Protocols and Practice

Web Protocols and Practice: HTTP/1.1, Networking Protocols, Caching, and Traffic Measurement, by Balachander Krishnamurty and Jennifer Rexford, ISBN 0-201-71088-9, Addison-Wesley, 2001.

If you want to know something about the underlying workings of the Web, you can find it somewhere out there on the Web itself. But, as we all know, it is not always easy to find the page you want, and particularly not if you are in a hurry and don't want to have to wade through documentation hierarchies or download PDF files. In these cases a real book is unbeatable, if one is available. Sadly, for information about the lower reaches of Web protocols there has been no single useful printed reference source available.

Organisation

This book fills that gap. It provides a detailed look at all the low level protocol issues as well as many other things; the book's subtitle sums it up admirably. The first section provides a brief history of the Web and its development which introduces all the important terminology and, most importantly, also says what the book is *not* about: nothing on XML (hurrah!), HTML, scripting languages, administration of Web servers, or specific products.

Section two moves on to more technical matters looking at Web clients, proxies and servers. The client chapter has a particularly useful section on spiders with an excellent table showing the names and calling hosts of the commonest spider programs. The information about proxies and servers is also of high quality and provide a solid grounding in how they interact with each other and the potential problems that can arise.

The third section looks at the protocols involved when using the Web. Starting with a concise run through TCP and the use of the DNS, the authors then glance at FTP, SMTP and NNTP, before going to a detailed examination of HTTP/1.1. In my personal experience, information on HTTP/1.1 has always been particularly inaccessible, both from the point of view of discoverability and readability, and this chapter explained several things that I had been puzzled about, especially about cache control which is rather a black art. (Also featured is a comprehensive table of HTTP return codes to which I shall turn quite often.) To finish this section of the book, there is a chapter on how HTTP interacts with TCP—a whole area that I had never really thought about before and which is much more complex that I would have thought it to be.

Next is a short section devoted to measuring and characterizing Web traffic. This a hugely contentious area and the discussion is well balanced and sensible. Following this the authors look in more detail at caching and at multimedia streaming, and manage to cover the latter topic without going into much unnecessary details about the actual bits that get sent whilst still giving a good coverage of the important material.

To round off the book, there are three chapters devoted to research topics, looking again at caching, measurement and protocol issues. Much of the material here is not directly of relevance to someone who is dealing with Web protocols on a daily basis, but there is still much here that will be of interest as the authors draw attention to places where improvements can be expected and how these might be realised.

Excellent Book

As you might expect, there is also a comprehensive bibliography and index. All in all an excellent book that is well researched, well written, and clearly set out without the excess of white space that is so common in computing books today. The price is perhaps rather high (I certainly could not recommend this as a textbook to my students—they simply could not afford it), but for people working in the industry it would be a worthwhile purchase and I think that they would soon find it an indispensable source of reference.

—Lindsay Marshall, *University of Newcastle upon Tyne*

Lindsay.Marshall@ncl.ac.uk

Summary of Acronyms

DNS: *Domain Name System*

FTP: *File Transfer Protocol*

HTTP: *HyperText Transfer Protocol*

NNTP: *Network News Transfer Protocol*

PDF: *Portable Document Format*

SMTP: *Simple Mail Transfer Protocol*

TCP: *Transmission Control Protocol*

XML: *Extensible Markup Language*

Would You Like to Review a Book for IPJ?

We receive numerous books on computer networking from all the major publishers. If you've got a specific book you are interested in reviewing, please contact us and we will make sure a copy is mailed to you. The book is yours to keep if you send us a review. We accept reviews of new titles, as well as some of the "networking classics." Contact us at **ipj@cisco.com** for more information.

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

Next ICANN Meeting, Marina del Rey, November 13–15, 2001

Many members of the *Internet Corporation for Assigned Names and Numbers* (ICANN) community wrote in response to a call for input as to whether the events of September 11 would affect their plans to travel to Los Angeles in November to attend the scheduled ICANN meetings. Almost without exception the respondents emphatically encouraged ICANN to hold its meetings and stated unequivocally that they planned to attend unless the international situation deteriorated to where travel was not practical.

Given this response and given the need to address emerging priorities, ICANN is planning to proceed with its November meeting, subject to any further serious change in the international situation that would affect travel conditions. However, as discussed below, the format of the meeting will differ significantly from what had previously been announced.

The events of September 11 have caused institutions worldwide to rethink their priorities and plans. As an international institution, ICANN is not immune. Although those events raise logistical and other concerns for holding meetings, they also underscore the need to address Internet stability issues, and security as a key component of stability. ICANN is not responsible for the overall security of the Internet. However, given ICANN's global responsibilities for the stability of the Internet's naming and addressing systems and under the new circumstances facing the international community, it would be irresponsible for ICANN not to conduct an in depth assessment of the robustness and security of these systems, and to take steps, if necessary, to strengthen the Internet in these regards. These are urgent matters and of worldwide importance.

The Internet is global in reach, as are the threats of terrorism. The events of September 11 offered a stark and tragic reminder of the incalculable importance of a reliable and secure naming and addressing system to support emergency response, personal and other communications, and information sharing. E-mail, instant messaging, and the Web, for example, all played essential roles.

Accordingly, the November ICANN meetings will focus on stability and security of the Internet's naming and addressing systems and of their operational implementation globally. This will be the overriding imperative for the meeting. As such, this will be a very different kind of meeting than previous ICANN meetings and will not follow the usual format.

At this meeting, ICANN will be seeking to promote discussion throughout the community on how to reassess areas of potential threats that could affect services within the scope of ICANN's responsibilities, how to improve readiness to meet these threats, and what additional policies or other actions should be considered and implemented to facilitate such improvements.

Clearly not all these questions will be answered in one meeting, but ICANN must now devote its energies as members of the global Internet community towards obtaining answers. Every constituency and supporting organization will be asked to report on its efforts to ensure the stability of the Internet's naming and addressing systems and what additional steps it proposes to take to improve that stability and security among its member organizations. Agenda items will be assessed for inclusion by what they contribute to the overall focus of the meeting.

Although a precise schedule has not yet been mapped out, these meetings will last three days from November 13 through 15, inclusive. Constituencies and supporting organizations will be asked to meet during this time to focus on the topic of the meeting. There will be a Board meeting at the end of the meeting to address essential business. The Board agenda will concentrate on topics where time is of the essence.

The focus of the meetings may well delay progress on some of the worthy and important initiatives that are currently underway. The effects of such delays have to be measured against the importance of ensuring the stability and security of the Internet itself. This will require patience on the part of those who may experience delays in matters of importance to them so that the ICANN community can bear down on the issue at hand.

This is only a preliminary announcement to enable attendees to firm up their travel plans. Details of the meeting will be announced as soon as possible. Please visit the ICANN Web site (<http://www.icann.org>) for further updates.

Van Jacobson Receives 2001 ACM SIGCOMM Award

Van Jacobson, the man widely credited with saving the Internet from an otherwise inevitable congestion collapse in the late 1980s, has been named the 2001 recipient of the ACM SIGCOMM Award. Jacobson is chief scientist at networking startup Packet Design, LLC.

The award is given annually by the *Association for Computing Machinery's Special Interest Group in Data Communications* (ACM SIGCOMM) to a recipient with a long and distinguished history of contributing to the field of data communications. Jacobson began his career in data communications developing control systems for the Department of Energy in the 1970s. He is best known for redesigning the TCP/IP protocol's flow-control algorithms to better handle congestion, preventing the Internet's collapse from traffic congestion in 1988–89. He is also widely recognized for his work on network synchronization effects, scalable multimedia protocols and applications, IP operations tools (for example *traceroute* and *pathchar*) and high-performance TCP implementations.

Prior to joining Packet Design as a member of the founding team, Jacobson was chief scientist at Cisco Systems, and before that had been group leader for Lawrence Berkeley Laboratory's Network Research Group.

The SIGCOMM Award has been presented every year since 1989. Prior recipients include Paul Baran, Vinton G. Cerf, David Farber and Leonard Kleinrock. ACM SIGCOMM is the world's largest professional society devoted to data communications. For more information, see: <http://www.acm.org/sigcomm/>

Useful Links

The following is a list of Web addresses that we hope you will find relevant to the material typically published in *The Internet Protocol Journal*. In the near future we will make these and other links available on our Web site: <http://www.cisco.com/ipj>

If you have suggestions for other pointers to include, please drop us a line at ipj@cisco.com

- The *Internet Engineering Task Force* (IETF). The primary standards-setting body for Internet technologies. <http://www.ietf.org>
- *Internet-Drafts* are working documents of the IETF, its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts. Internet-Drafts are not an archival document series. These documents should not be cited or quoted in any formal document. Unrevised documents placed in the Internet-Drafts directories have a maximum life of six months. After that time, they must be updated, or they will be deleted. Some Internet-Drafts become RFCs (see below). <http://www.ietf.org/ID.html>
- The *Request For Comments* (RFC) document series. The RFCs form a series of notes, started in 1969, about the Internet (originally the ARPANET). The notes discuss many aspects of computer communication, focusing on networking protocols, procedures, programs, and concepts but also including meeting notes, opinion, and sometimes humor. The specification documents of the Internet protocol suite, as defined by IETF and its steering group the IESG, are published as RFCs. Thus, the RFC publication process plays an important role in the Internet standards process. <http://www.rfc-editor.org/>
- The *Internet Society* (ISOC) is a non-profit, non-governmental, international, professional membership organization. <http://www.isoc.org>
- The *Internet Corporation for Assigned Names and Numbers* (ICANN) "... is the non-profit corporation that was formed to assume responsibility for the IP address space allocation, protocol parameter assignment, domain name system management, and root server system management functions previously performed under U.S. Government contract by IANA and other entities." <http://www.icann.org>

- The *North American Network Operators' Group* (NANOG) “...provides a forum for the exchange of technical information, and promotes discussion of implementation issues that require community cooperation. Coordination among network service providers helps ensure the stability of overall service to network users.”
<http://www.nanog.org>
- The *Regional Internet Registries* (RIRs) provide IP address block assignments for Internet Service Providers and others. Currently, there are three active RIRs:
 - The *Asia Pacific Network Information Centre* (APNIC):
<http://www.apnic.net>
 - *RIPE Network Coordination Centre*—the RIR responsible for Europe and Northern Africa: <http://www.ripe.net>
 - *American Registry for Internet Numbers* (ARIN)—the RIR responsible for the Americas and Sub-Saharan Africa:
<http://www.arin.net>

Two more RIRs are in the process of formation: *AfriNIC* for Africa and *LACNIC* for Central- and Latin America.
- The *World Wide Web Consortium* (W3C) “... develops interoperable technologies (specifications, guidelines, software, and tools) to lead the Web to its full potential as a forum for information, commerce, communication, and collective understanding.”
<http://www.w3.org/>
- The *International Telecommunication Union* (ITU) “... is an international organization within which governments and the private sector coordinate global telecom networks and services.”
<http://www.itu.int>
- The *International Organization for Standardization* (ISO) “... is a worldwide federation of national standards bodies from some 140 countries, one from each country. The mission of ISO is to promote the development of standardization and related activities in the world with a view to facilitating the international exchange of goods and services, and to developing cooperation in the spheres of intellectual, scientific, technological and economic activity. ISO's work results in international agreements which are published as International Standards.” <http://iso.org>

This is by no means intended to be a complete list of organizations that are related to Internet development in one way or another, but this list should give you a good starting point.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Technology
WorldCom, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 2001 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-10/5
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSR STD U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol Journal

December 2001

Volume 4, Number 4

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Scaling Inter-Domain Routing.....	2
Regional Internet Registries ...	17
Book Reviews	30
Letters to the Editor.....	34
Fragments	38
Call for Papers	39

FROM THE EDITOR

In a previous article entitled “Analyzing the Internet BGP Routing Table,” Geoff Huston examined many issues relating to the operation of today’s Internet. In this issue he goes a step further and suggests ways in which the fundamental routing architecture could be changed to solve problems related to routing-table growth. The article is called “Scaling Inter-Domain Routing—A View Forward.”

The IP address space is administered by three entities, namely APNIC, ARIN and RIPE NCC. Collectively referred to as the *Regional Internet Registries* (RIRs), these organizations are responsible for address allocation to their member organizations (typically national registries or large Internet Service Providers). Since the IPv4 address space is a limited resource, this allocation has to be done with care, while accounting for the needs of the address space consumers. We asked the RIRs for an overview of the work they perform. What we received was a joint effort that not only describes the RIR structure, but also gives some historical background on the evolution of IP addressing and routing.

We were pleased to receive a couple of Letters to the Editor recently, both in response to articles in our previous issue. This kind of feedback is most welcome and we encourage you to send your comments and suggestions to ipj@cisco.com

We’d like to remind you that all back issues of *The Internet Protocol Journal* can be downloaded from www.cisco.com/ipj. Click on “IPJ Issues” and you will be taken to the appropriate section.

By the time you read this, our online subscription system should be operational. You will find it at our Web site: www.cisco.com/ipj. Please let us know if you encounter any difficulties by sending e-mail to ipj@cisco.com

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Scaling Inter-Domain Routing—A View Forward

by Geoff Huston, Telstra

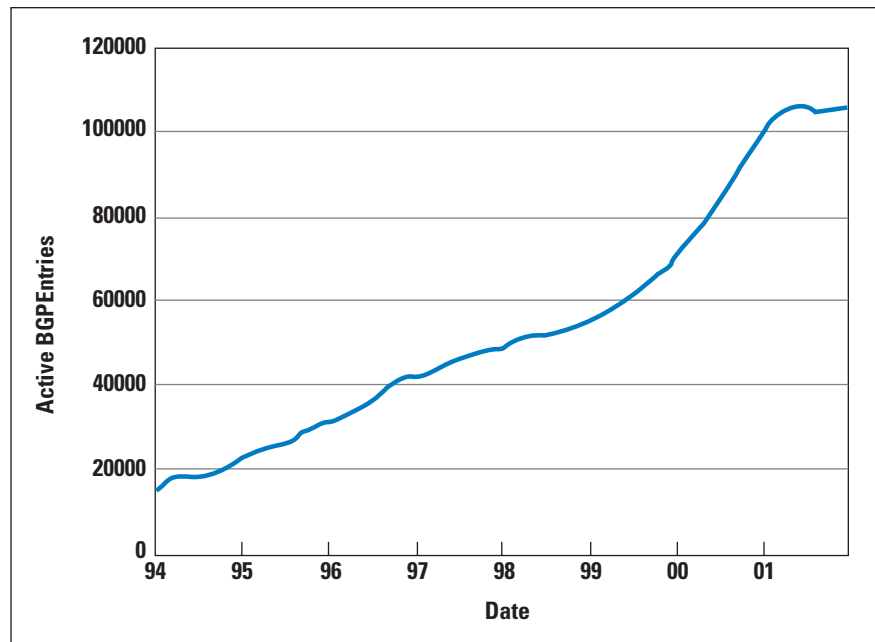
In the previous IPJ article, “Analyzing the Internet BGP Routing Table,” (Vol. 4, No. 1, March 2001) we looked at the characteristics of the growth of the routing table in recent years. The motivation for this work is to observe aspects of the Internet routing table in order to understand the evolving structure of the Internet and thereby attempt to predict some future requirements for routing technology for the Internet.

The conclusions drawn in the previous article included the observation that multihomed small networks appeared to be a major contributor to growth of the Internet routing system. It also observed that there was a trend toward a denser mesh of inter-Autonomous System connectivity within the Internet. At the same time there has been an increase of various forms of policy-based constraints imposed upon this connectivity mesh, probably associated with a desire to undertake various forms of inter-domain traffic engineering through manipulation of the flow of routing information.

Taken together, these observations indicate that numerous strong growth pressures are being exerted simultaneously on the inter-domain routing space. Not only is the network itself growing in size, but also the internal interconnectivity of the network is becoming more densely meshed. The routing systems that are used to maintain a description of the network connectivity are being confronted with having to manipulate smaller route objects that describe finer levels of network detail. This is coupled with lengthening lists of qualifying attributes that are associated with each route object. The question naturally arises as to whether the *Border Gateway Protocol* (BGP) and the platforms used to support BGP in the Internet today can continue to scale at a pace that matches the growth in demands that are being placed upon it.

The encouraging news is that there appears to be no immediate cause for concern regarding the capability of BGP to continue to support the load of routing the Internet. The processor and memory capacity in current router platforms is easily capable of supporting the load associated with various forms of operational deployment models, and the protocol itself is not in imminent danger of causing network failure through any internal limitation within the protocol itself. Also, numerous network operators have exercised a higher level of care as to how advertisements are passed into the Internet domain space and, as a result, the growth rates for the routing table over 2001 shows a significant slowdown over the rates of the previous two years (Figure 1).

Figure 1: BGP Table
Size 1994–2001



However, the observed trends in inter-domain routing of an increasingly detailed and highly qualified view of a more densely interconnected and still-growing network provide adequate grounds to examine the longer-term routing requirements. It is useful, therefore, to pose the question as to whether we can continue to make incremental changes to the BGP protocol and routing platforms, or whether the pace of growth will, at some point in time, mandate the adoption of a routing architecture that is better attuned to the evolving requirements of the Internet.

This article does not describe the operation of an existing protocol, nor does it describe any current operational practice. Instead it examines those aspects of inter-domain routing that are essential to today's Internet, and the approaches that may be of value when considering the evolution of the Internet inter-domain routing architecture. With this approach, the article illustrates one of the initial phases in any technology development effort—that of an examination of various requirements that could or should be addressed by the technology.

Attributes of an Inter-Domain Routing Architecture

Let's start by looking at those aspects of the inter-domain routing environment that could be considered a base set of attributes for any inter-domain routing protocol.

Accuracy

For a routing system to be of any value, it should accurately reflect the forwarding state of the network. Every routing point is required to have a consistent view of the routing system in order to avoid forwarding loops and black holes (points where there is no relevant forwarding information and the packet must be discarded). Local changes in underlying physical network, or changes in the policy configuration of the network at any point, should cause the routing system to compute a new distributed routing state that accurately reflects the changes.

This requirement for accuracy and consistency is not, strictly speaking, a requirement that every node in a routing system has global knowledge, nor a requirement that all nodes have precisely the same scope of information. In other words, a routing system that detects and avoids routing loops and inconsistent black holes does not necessarily need to use routing systems that rely on uniform distribution of global knowledge frameworks.

Scalability

Scalability can be expressed in many ways, including the number of routing entries, or prefixes, carried within the protocol, the number of discrete routing entities within the inter-domain routing space, the number of discrete connectivity policies associated with these routing entries, and the number of protocols supported by the protocol. Scalability also needs to encompass the dynamic nature of the network, including the number of routing updates per unit of time, time to converge to a coherent view of the connectivity of the network following changes, and the time taken for updates to routing information to be incorporated into the network forwarding state. In expressing this ongoing requirement for scalability in the routing architecture, there is an assumption that we will continue to see an Internet that is composed of a large number of providers, and that these providers will continue to increase the density of their interconnection.

The growth trends in the inter-domain routing space do not appear to have well-defined upper limits, so placing bounds on various aspects of the routing environment is impractical. The only practical way to describe this attribute is that it is essential to use a routing architecture that is scalable to a level well beyond the metrics of today's Internet.

In the absence of specific upper bounds to quantify this family of requirements, the best we conclude here is that at present we are working in an inter-domain environment that manipulates some 10^5 distinct routing entries, and at any single point of interconnection there may be of the order of 10^6 routing protocol elements being passed between routing domains. Experience in scaling transmission systems for the Internet indicates that an improvement of a single order of magnitude in the capacity of a technology has a relatively short useful lifetime. It would, therefore, be reasonable to consider that a useful attribute is to be able to operate in an environment that is between two to three orders of magnitude larger than today's system.

Policy Expressiveness

Routing protocols perform two basic tasks: first, determining if there is at least one viable path between one point in the network and another, and secondly, where there is more than one such path, determining the "best" such path to use. In the case of interior routing protocols, "best" is determined by the use of administratively assigned per-link metrics, and a "best" path is one that minimizes the sum of these link metrics.

In the case of the inter-domain routing protocols, no such uniformly interpreted metric exists, and “best” is expressed as a preference using network paths that yield an optimal price and performance outcome for each domain.

The underlying issue here is that the inter-domain routing system must straddle a collection of heterogeneous networks, and each network has a unique set of objectives and constraints that reflect the ingress, egress, and transit routing policies of a network. Ingress routing policies reflect how a network learns information, and which learned routes have precedence when selecting a routing entry from a set of equivalent routes. In a unicast environment, exercising control over how routes are learned by a domain has a direct influence over which paths are taken by traffic leaving the domain. Egress policies reflect how a domain announces routes to its adjacent neighbors. A domain may, for example, wish to announce a preferential route to a particular neighbor, or indicate a preference that the route not be forwarded beyond the adjacent neighbor. In a unicast environment, egress routing policies have a bearing on which paths are used for traffic to reach the domain. Transit routing policies control how the routes learned from an adjacent domain are advertised to other adjacent domains. If a domain is a transit provider for another domain, then a typical scenario for the transit provider would be to announce all learned routes to all other connected domains. For a multi-homed transit customer, routes learned from one transit provider would normally not be announced to any other transit provider.

This requirement for policy expressiveness implies that the inter-domain routing protocol should be able to attach various attributes to protocol objects, allowing a domain to communicate its preferences relating to handling of the route object to remote domains.

Robust Predictable Operational Characteristics

A routing system should operate in such a way that it achieves predictable outcomes. The inference here is that under identical initial conditions a routing system should always converge to the same routing state, and that with knowledge of the rules of operation of the protocol and the characteristics of the initial environment, an observer can predict what this state will be. Predictability also implies stability of the routing environment, such that a routing state should remain constant for as long as the environment itself remains constant.

The routing protocol should operate in a way that tends to damp propagation of dynamic changes to the routing system rather than amplify such changes. This implies that minor variations in the state of the network should not cause large-scale instability across the entire network while a new stable routing state is reached. Instead, routing changes should be propagated only as far as necessary to reach a new stable state, so that the global requirement for stability implies some degree of locality in the behavior of the system.

The routing system should have robust convergence properties. A change in the physical configuration or policy environment in any part of the network causes a distributed computation of the routing state. Convergence implies that this distributed computation reaches a conclusion at some point. The requirement for a robust convergence property implies that the distributed computation should always halt, that the halting point be reached quickly, and the system should avoid generating transitory incorrect intermediate routing states. The interpretation of “quickly” in this context is variable. Currently, this value for BGP convergence time is of the order of tens to hundreds of seconds. In order to support increasingly time-critical applications, there appears to be an emerging requirement to reduce the median convergence time for the inter-domain routing protocol to a small number of seconds.

Efficiency

The routing system should be efficient, in that the amount of network resources, in terms of bandwidth and processing capacity of the network switching elements, should not be disproportionately large. This is an area of trade-off in that the greater the amount of information passed within the routing system and the greater the frequency of such information exchanges, the greater the level of expectation that the routing system can continuously maintain an accurate view of the connectivity of the network, but at a cost of higher overhead. It is necessary to pass enough information across the system to allow each routing element to have a sufficiently accurate view of the network, yet ensure that the total routing overhead is low.

Evolving Requirements of Inter-Domain Routing

Layered on top of the base set of routing requirements listed above are a second set of requirements that can be seen as reflecting current directions in the deployed Internet, and are not necessarily well integrated into the existing routing architecture.

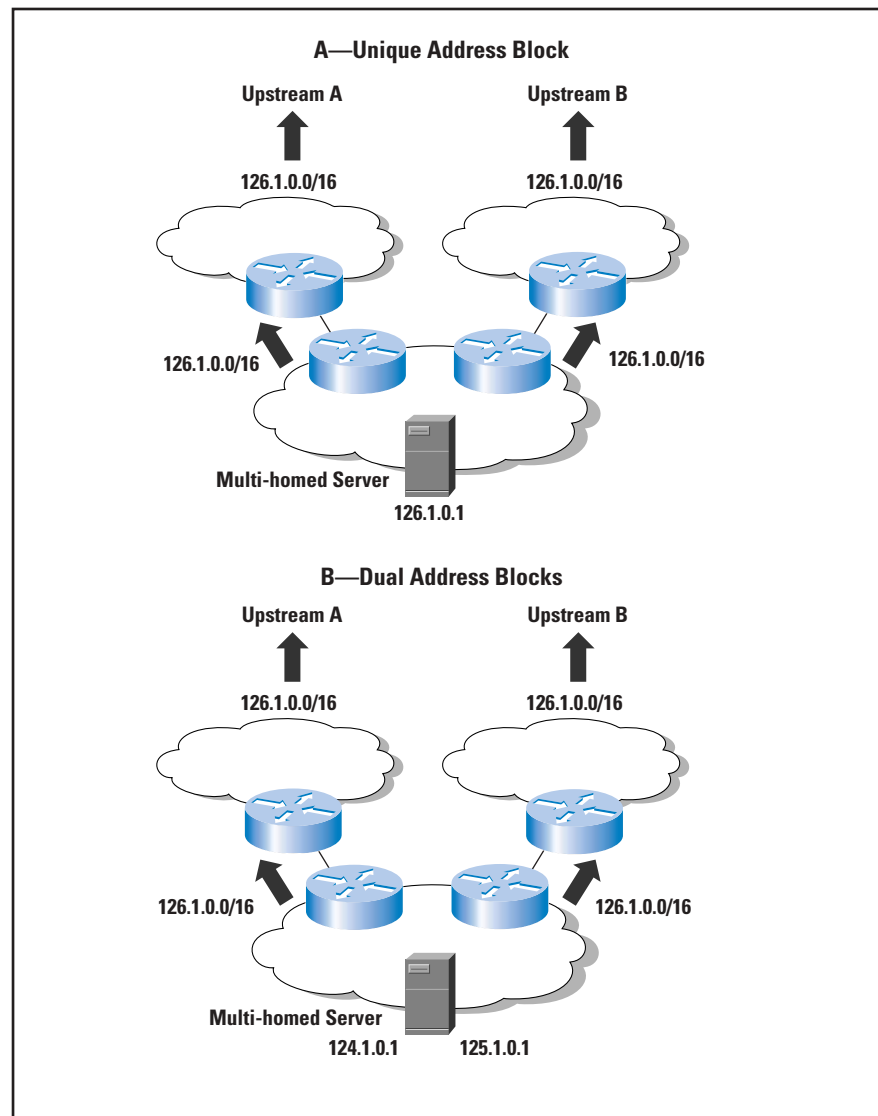
Multi-Homing of Edge Networks

Multi-homing refers to the practice of using more than one upstream transit provider. The common motivation for such a configuration is that if service from one transit provider fails, the customer can use the other provider as a means of service restoration. It may also allow some form of traffic balancing across multiple services. With careful use of route policies, the customer can direct traffic to each provider to minimize delay and loss, achieving some improved application performance.

The issue presented by multi-homing is that the multi-homed network is now not wholly contained within a service hierarchy of any particular provider. This implies that routing information describing reachability to the multi-homed customer cannot readily be aggregated into any single provider’s routing advertisements, and the usual outcome is that the multi-homed customer must independently announce its reachability to each transit provider, who in turn must propagate this information across the routing system.

The evolving requirement here is one that must be able to integrate the demands of an increasing use of multi-homing into the overall network design. Two basic forms of approach can be used here—one is to use a single address block across the customer network and announce this block to all transit providers as an unaggregatable routing advertisement into the inter-domain routing system, and the other is to use multiple address blocks drawn from each provider's address block, and use either host-based software or some form of dynamic address translation within the network in order to use a source address drawn from a particular provider's block for each network transaction (Figure 2). The second approach is not widely used, and for the immediate future the requirement for multi-homing is normally addressed by using unique address blocks for the multi-homed network that are not part of any provider's aggregated address blocks. The consequence of this is that widespread use of multi-homing as a means of service resiliency will continue to have an impact on the inter-domain routing system.

Figure 2: Routing Approaches to Multi-Homing



Inter-Domain Traffic Engineering

In an increasingly densely interconnected network, selecting and using just one path between two points is not an optimal outcome of a routing architecture. Of more importance is the ability to identify a larger set of viable paths between these points and distribute the associated traffic flows in such a way that each individual transaction uses a single path, but the total set of flows is distributed across the set of paths.

To achieve this outcome, more information must be placed into the routing system, allowing a route originator to describe the policy-based preferences of which sets of paths should be preferred for traffic destined to the route originator, allowing a transit service operator to add information regarding current preferences associated with using particular transit paths, and allowing the traffic originator the ability to use local traffic egress policies to reach the destination. These traffic engineering-related preferences are not necessarily represented by static values of routing attributes. One of the requirements of traffic engineering is to allow the network to dynamically respond to shifting traffic load patterns, and this implies that there is a component of dynamic information update that is associated with such traffic engineering-related aspects of the routing system.

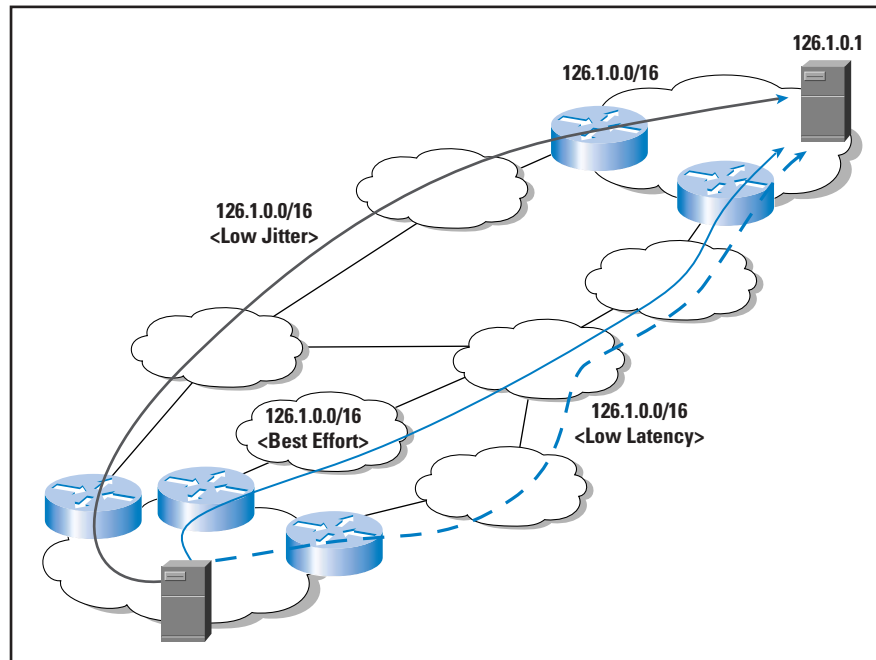
At an abstract level, this greater volume of routing information is needed in order to address the dual role of the routing system as both an inter-domain connectivity maintenance protocol and as a traffic-engineering tool.

Inter-Domain Quality of Service

Quality of Service (QoS) is a term that encompasses a wide variety of mechanisms. In the case of routing, the term is used to describe the process of modifying the normal routing response of associating a single forwarding action with a destination address prefix in such a way that there may be numerous forwarding decisions for a particular address prefix. Each forwarding decision is associated with a particular service response, so that a “best-effort” path to a particular destination address may differ from a “low-latency” path, which in turn may differ from a “high-bandwidth” path, and so on.

As with inter-domain traffic engineering, this requirement is one which would be expected to place greater volumes of information into the routing domain. At an abstract level this requirement can be seen as the association of a service quality attribute with an address prefix, and passing the paired entity into the routing domain as a single routing object. The inference is that multiple quality attributes associated with a path to a particular prefix would require the routing system to independently manipulate multiple route objects, because it would be reasonable to anticipate that the routing system would select different paths to reach the same address prefix if different QoS service attributes were used as a path qualifier (Figure 3).

Figure 3: Inter-Domain
Routing with QoS



Approaches to Inter-Domain Routing

Let's now take this set of requirements and attempt to match them to various approaches to routing protocols.

Routing is a distributed computation wherein each element of the computation set must reach an outcome that is consistent with all other computations undertaken by other members of the set. There are two major approaches to this form of distributed computation, namely *serial* or *parallel* computation. Serial computation involves each element of the set undertaking a local computation and then passing the outcomes of this computation to its adjacent elements. This approach is used in various forms of distance-vector routing protocols where each routing node computes a local set of selected paths, and then propagates the set of reachable prefixes and the associated path metric to its neighbors. Parallel computation involves rapid flooding of the current state of connectivity within the set to all elements, and all set elements simultaneously compute forwarding decisions using the same base connectivity data. This approach is used in various forms of link-state routing protocols, where the protocol uses a flooding technique to rapidly propagate updated link-status information and then relies on each routing node to perform a local path selection computation for each reachable address prefix. Is one of these approaches substantially better suited than the other to the inter-domain routing environment?

Open or Closed Routing Policies

One of the key issues behind consideration of this topic is that of the role of *local policy*. Using a distance-vector protocol, a routing domain gathers selected path information from its neighbors, applies local policy to this information, and then distributes this updated information in the form of selected paths to its neighbor domains.

In this model the nature of the local policy applied to the routing information is not necessarily visible to the domain neighbors, and the process of converting received route advertisements into advertised route advertisements uses a local policy process whose policy rules are not visible externally. This scenario can be described as *policy opaque*. The side effect of such an environment is that a third party cannot remotely compute which routes a network may accept and which may be readvertised to each neighbor.

In link-state protocols, a routing domain effectively broadcasts its local domain adjacencies, and the policies it has with respect to these adjacencies, to all nodes within the link-state domain. Every node can perform an identical computation upon this set of adjacencies and associated policies in order to compute the local inter-domain forwarding table. The essential attribute of this environment is that the routing node has to announce its routing policies in order to allow a remote node to compute which routes will be accepted from which neighbor, and which routes will be advertised to each neighbor and what, if any, attributes are placed on the advertisement. Within an interior routing domain the local policies are in effect metrics of each link, and these policies can be announced within the routing domain without any consequent impact.

In the exterior routing domain it is not the case that interconnection policies between networks are always fully transparent. Various permutations of supplier/customer relationships and peering relationships have associated policy qualifications that are not publicly announced for business competitive reasons. The current diversity of interconnection arrangements appears to be predicated on policy opacity, and to mandate a change to a model of open interconnection policies may be contrary to operational business imperatives. An inter-domain routing tool should be able to support models of interconnection where the policy associated with the interconnection is not visible to any third party. If the architectural choice is a constrained one between distance vector and link state, then this consideration would appear to favor the continued use of a distance-vector approach to inter-domain routing. This choice, in turn, has implications on the convergence properties and stability of the inter-domain routing environment. If there is a broader spectrum of choice, the considerations of policy opacity would still apply.

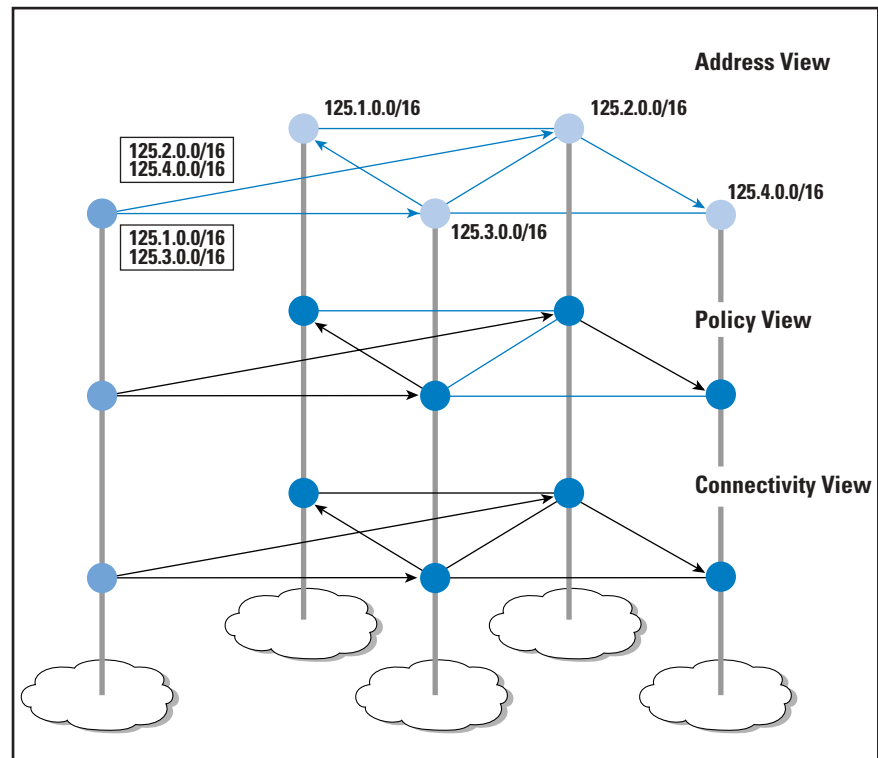
Separation of Functions

The inter-domain routing function undertakes many roles simultaneously. First, it maintains the current view of inter-domain connectivity. Any changes in the adjacency of a domain are reflected in a distributed update computation that determines if the adjacency change implies a change in path selection and in address reachability. Secondly, it maintains the set of currently reachable address prefixes. And finally, the protocol binds the first two functions together by associating each prefix with a path through the inter-domain space.

This association uses a policy framework to allow each domain to select a path that optimizes local policy constraints within the bounds of existing constraints applied by other domains. This policy may be related to traffic-engineering objectives, QoS requirements, local cost optimization, or related operational or business objectives.

An alternative approach to inter-domain routing is to separate the functions of connectivity maintenance, address reachability, and policy negotiation. As an example of this approach, a connectivity protocol can be used to identify all viable paths between a source and a destination domain. A policy negotiation protocol can be used to ensure that there are a consistent sequence of per-domain forwarding decisions that will pass traffic from the source domain to the destination domain. An address reachability protocol can be used to associate a collection of address prefixes with each destination domains. This framework is illustrated in Figure 4.

Figure 4: A Multi-Tiered Approach to Inter-Domain Routing



Address Prefixes and Autonomous System Numbers

One observation about the current inter-domain routing system is that it uses a view of the network based on computing the optimal path to each address prefix. This view is translated into an inter-domain routing protocol that uses the address prefix as the basic protocol element and attaches various attributes to each address prefix as they are passed through the network

As of late 2001, the routing system had some 100,000 distinct address prefixes and 11,500 origin domains. This implies that each origin domain is responsible for an average of 8 to 9 address prefixes. If each domain advertised its prefixes with a consistent policy, then each address prefix would be advertised with identical attributes. If the routing protocol were to be inverted such that the routing domain identifier, or *Autonomous System* number, were the basic routing object and the set of prefixes and associated common set of route attributes were attributes of the Autonomous System object, then the number of routing objects would be reduced by the same factor of between 8 and 9.

The motivation in this form of approach is that seeking clear hierarchical structure in the address space as deployed is no longer feasible, and that no further scaling advantage can be obtained by various forms of address aggregation within the routing system. This approach replaces this address-based hierarchy with a two-level hierarchy of routing domains. Within a routing domain, routing is undertaken using the address prefix. Between routing domains, routing is undertaken using domain identifiers and associated sets of domain attributes.

Although this approach appears to offer some advantage in creating a routing domain, one-tenth of the size of the address prefix-based routing domain, it is interesting to note that since late 1996 the average number of address prefixes per Autonomous System has fallen from 25 to the current value of 9. In other words, the number of distinct routing domains is growing at a faster rate than the number of routed address prefixes. While the adoption of a domain-based routing protocol offers some short-term advantages in scaling, the longer-term prospects are not so attractive, given these relative growth rates.

Routing Hierarchies of Information

The scaling properties of an inter-domain routing protocol are related on the ability of the protocol to remove certain specific items of information from the routing domain at the point where it ceases to have any differentiating impact. For example, it is important for a routing protocol to carry information that a particular domain has multiple adjacencies and that there are a number of policies associated with each adjacency, and propagate this information to all local domains. At a suitably distant point in the network, the forwarding decision remains the same regardless of the set of local adjacencies, and propagation of the detail of the local environment to points where the information ceases to have any distinguishing outcome is unproductive.

From this perspective, scaling the routing system is not a case of determining what information can be added into the routing domain, but instead it's a case of determining how much information can be removed from the routing domain, and how quickly.

One way of removing information is through the use of *hierarchies*. Within a hierarchical structure, a set of objects with similar properties are aggregated into a single object with a set of common properties. One way to perform such aggregation is by increasing the amount of information contained in each aggregate route object. For example, if single route objects are to be used that encompass a set of address prefixes and a collection of Autonomous Systems, then it would be necessary to define additional attributes within the route object to further qualify the policies associated with the object in terms of specific prefixes, specific Autonomous Systems, and specific policy semantics that may be considered as policy exceptions to the overall aggregate. This approach would allow aggregation of routing information to occur at any point in the network, allowing the aggregator to create a compound object with a common set of attributes, and a set of additional attributes that apply to a particular subset of the aggregate.

Another approach to using hierarchies to reduce the number of route objects is to reduce the scope of advertisement of each routing object, allowing the object to be removed and proxy aggregated into some larger object when the logical scope of the object is reached. This approach would entail the addition of route attributes that could be used to define the circumstances where a specific route object would be subsumed by an aggregate route object without impacting the policy objectives associated with the original set of advertisements. This approach places control of aggregation with the route object originator, allowing the originator to specify the extent to which a specific route object should be propagated before being subsumed into an aggregate object.

It is not entirely clear that the approach of exploiting hierarchies in an address space is the most appropriate response to scaling pressures. Viewed from a more general perspective, scaling of the routing system requires the systematic removal of information from the routing domain. The way this is achieved is by attempting to align the structure of deployment with some structural property of the syntax of the protocol elements that are being used as routing objects. Information can then be eliminated through systematic aggregation of the routing objects at locations within the routing space that correspond to those points in the topology of the network where topology aggregation is occurring. The maintenance of this tight coupling of the structure of the deployed network to the structure of the identifier space is the highest cost of this approach. Alterations to the topology of the network through the relocation or reconfiguration of networks requires renumbering of the protocol element if hierarchical aggregation is to be maintained. If the address space is the basis of routing, as at present, then this becomes a large-scale exercise of renumbering networks that in turn implies an often prohibitively disruptive and expensive exercise of renumbering collections of host systems and associated services.

One view of this is that the connectivity properties of the Internet are already sufficiently meshed that there is no readily identifiable hierarchical structure, and that this trend is becoming more pronounced, not less. In that case, the most appropriate course of action may be to reexamine the routing domain and select some other attribute as the basis of the routing computation that does not have the same population, complexity, and growth characteristics as address prefixes, and base the routing computation on this attribute. One such alternative approach is to consider Autonomous System numbers as routing “atoms” where the routing system converges to select an Autonomous System path to a destination Autonomous System, and then uses this information to add the associated set of prefixes originated by this Autonomous System, and next-hop forwarding decision to reach this Autonomous System into the local forwarding table.

Extend or Replace BGP

A final consideration is to consider whether these requirements can best be met by an approach of a set of upward-compatible extensions to BGP, or by a replacement to BGP.

The rationale for extending BGP would be to increase the number of commonly supported transitive route attributes, and, potentially, allow a richer syntax for attribute definition which in turn would allow the protocol to use a richer set of semantic definitions in order to express more complex routing policies.

This direction may sound like a step backward, in that it proposes an increase in the complexity of the route objects carried by the protocol and potentially increases the amount of local processing capability required to generate and receive routing updates. However, this can be offset by potential benefits that are realizable through the greater expressive capability for the policy attributes associated with route objects. It can allow a route originator an ability to specify the scope of propagation of the route object, rather than assuming that propagation will be global. The attributes can also describe intended service outcomes in terms of policy and traffic engineering. It may also be necessary to allow BGP sessions to negotiate additional functionality intended to improve the convergence behavior of the protocol. Whether such changes can produce a scalable and useful outcome in terms of inter-domain routing remains, at this stage, an open question.

An alternative approach is that of a replacement protocol. Use of a parallel-processing approach to the distributed computation of routing, such as that used in the link-state protocols, can offer the benefits of faster convergence times and avoidance of unstable transient routing states. On the other hand, link-state protocols present issues relating to policy opaqueness, as described above. Another major issue with such an approach is the need to address the efficiency of inter-domain link-state flooding.

The inter-domain space would need some further levels of imposed structure similar to intra-domain areas in order to ensure that individual link updates are rapidly propagated across the relevant subset of the network. The use of such an area structure may well imply the need for an additional set of operator relationships, such as mutual transit. Such inter-domain relationships may prove challenging to adapt to existing operator practices.

Another approach could be based on the adoption of a multi-layer approach of separate protocols for separate functions, as described above. A base inter-domain connectivity protocol could potentially be based on a variant of a link-state protocol, using the rapid convergence properties of such protocols to maintain a coherent view of the current state of connectivity within the network. The overlay of a policy protocol would be intended as a signaling mechanism to allow each domain to make local forwarding decisions that are consistent with those adopted by adjacent domains, thereby maintaining a collection of coherent inter-domain paths from source to destination. Traffic engineering can also be envisaged as an overlay mechanism, allowing a source to make a forwarding decision that selects a path to the destination where the characteristics of the path optimize the desired service outcomes.

Directions for Further Activity

Although short-term actions based on providing various incentives for network operators to remove redundant or inefficiently grouped entries from the BGP routing table may exist, such actions are short-term palliative measures, and will not provide long-term answers to the need for a scalable inter-domain routing protocol. One approach to the longer-term requirements may be to preserve many of the attributes of the current BGP protocol, while refining other aspects of the protocol to improve its scaling and convergence properties. A minimal set of alterations could retain the Autonomous System concept to allow for administrative boundaries of information summarization, as well as retaining the approach of associating each prefix advertisement with an originating Autonomous System. The concept of policy opaqueness would also be retained in such an approach, implying that each Autonomous System accepts a set of route advertisements, applies local policy constraints, and readvertises those advertisements permitted by the local policy constraints. It could be feasible to consider alterations to the distance-vector path-selection algorithm, particularly as it relates to intermediate states during processing of a route withdrawal. It is also feasible to consider the use of compound route attributes, allowing a route object to include an aggregate route, and numerous specifics of the aggregate route, and attach attributes that may apply to the aggregate or a specific address prefix. Such route attributes could be used to support multi-homing and inter-domain traffic-engineering mechanisms. The overall intent of this approach is to address the major requirements in the inter-domain routing space without using an increasing set of globally propagated specific route objects.

Another approach is to consider the feasibility of decoupling the requirements of inter-domain connectivity management with the applications of policy constraints and the issues of sender- and receiver-managed traffic-engineering requirements. Such an approach may use a link-state protocol as a means of maintaining a consistent view of the topology of inter-domain network, and then use some form of overlay protocol to negotiate policy requirements of each Autonomous System, and use a further overlay to support inter-domain traffic-engineering requirements. The underlying assumption of such an approach is that if the functional role of inter-domain routing is divided into distinct components, each component will have superior scaling and convergence properties which in turn will result in superior properties for the entire routing system. Obviously, this assumption requires some testing.

Research topics with potential longer-term application include the approach of drawing a distinction between the identity of a network, its location relative to other networks, and maintenance of a feasible path set between a source and destination network that satisfies various policy and traffic-engineering constraints. Again the intent of such an approach would be to divide the current routing function into numerous distinct scalable components rather than using a single monolithic routing protocol.

Further Reading

- [0] Huston, G., "Analyzing the Internet BGP Routing Table," *The Internet Protocol Journal*, Vol. 4, No. 1, March 2001.
www.cisco.com/warp/public/759/ipj_4-1/ipj_4-1_bgp.html
- [1] Huitema, C., *Routing in the Internet, 2nd Edition*, ISBN 0130226475, Prentice Hall, January 2000. *A good introduction to the general topic of IP routing.*
- [2] Rekhter, Y., and Li T., "A Border Gateway Protocol 4 (BGP-4)," RFC 1771, March 1995. *The base specification of BGP 4. This document is currently being updated by the IETF. The state of this work in progress as of November 2001 is documented as an Internet Draft,*
[draft-ietf-idr-bgp4-15.txt](#)
- [3] Elwyn Davies et al., "Future Domain Routing Requirements," work in progress, July 2001. *This work is currently documented as an Internet Draft, **[draft-davies-fdr-reqs-01.txt](#)**. It contains a review of an earlier effort in enumerating routing requirements ("Goals and Functional Requirements for Inter-Autonomous System Routing," RFC 1126, October 1989), as well as a commentary on a proposed set of current routing requirements.*

GEOFF HUSTON holds a B.Sc. and M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. He is currently the Chief Scientist in the Internet area for Telstra, a member of the Internet Architecture Board, and is the Secretary of the APNIC Executive Committee. He is author of *The ISP Survival Guide*, ISBN 0-471-31499-4, *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, and coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. All three books are published by John Wiley & Sons. E-mail: **gih@telstra.net**

Development of the Regional Internet Registry System

by Daniel Karrenberg, RIPE-NCC; Gerard Ross, APNIC; Paul Wilson, APNIC; Leslie Nobile, ARIN

The current system of managing Internet address space involves *Regional Internet Registries* (RIRs), which together share a global responsibility delegated to them by the *Internet Assigned Numbers Authority* (IANA). This regime is now well established, but it has evolved over ten years from a much simpler, centralized system. Internet number spaces were originally managed by a single individual “authority,” namely the late Jon Postel, co-inventor of some of the most important technical features of today’s Internet.

It is important to understand that the evolution of the RIR system was not simply the result of Internet growth and the natural need to refine and decentralize a growing administrative task. On the contrary, it arose from, and closely tracked, the technical evolution of the Internet Protocol, in particular the development of today’s IP addressing and routing architecture.

In a relatively short time, the Regional Internet Registry system has evolved into a stable, robust environment for Internet address management. It is maintained today through self-regulatory practices that are well established elsewhere in the Internet and other industries, and it maintains its legitimacy and relevance by firmly adhering to open, transparent, participatory decision-making processes.

Before the RIRs:

IP Address Architecture

An important feature of the Internet Protocol (IP) is the ability to transparently use a wide variety of underlying network architectures to transport IP packets. This is achieved by encapsulating IP packets in whatever packet or frame structure the underlying network uses. Routers connecting different networks forward IP traffic by decapsulating incoming IP packets and then re-encapsulating them as appropriate for the next network to carry them.

To achieve this task with full transparency, the IP needed an addressing structure, which developed as a two-level hierarchy in both addressing and routing. One part of the address, the *network* part, identifies the particular network a host is connected to, while the other part, the *local* part, identifies the particular end system on that network.

Internet routing, then, has to deal only with the network part of the address, routing the packet to a router directly connected to the destination network. The local part is not used at all in Internet routing itself; rather it is used to determine the intended address within the addressing structure of the destination network.

The method by which the local part of an IP address is translated to a local network address depends on the architecture of the destination network—static tables, simple conversions, or special-purpose protocols are used as appropriate.

The original Internet addresses comprised 32 bits, the first 8 bits providing the network part and the remaining 24 bits the local part. These addresses were used for many years. However, in June 1978, in Internet Engineering Note (IEN) 46 “A proposal for addressing and routing in the internet,” Clark and Cohen observed:

“The current internet header has space to name 256 networks. The assumption, at least for the time being, is that any network entering the internet will be assigned one of these numbers. While it is not likely that a great number of large nets, such as the ARPANET, will join the internet, the trend toward local area networking suggests that a very large number of small networks can be expected in the internet in the not too distant future. We should thus begin to prepare for the day when there are more than 256 networks participating in the internet.”

Classful Addressing

As predicted, it soon became necessary to adapt the address architecture to allow more networks to be connected. By the time the Internet Protocol itself was comprehensively specified (in RFC 790, published in 1981, edited by Jon Postel), the IP address could be segmented in numerous ways to provide three classes of network address.

In Class A, the high-order bit is zero, the next 7 bits are the network, and the last 24 bits are the local address. In Class B, the high-order 2 bits are one-zero, the next 14 bits are the network, and the last 16 bits are the local address. In Class C, the high-order 3 bits are one-one-zero, the next 21 bits are the network, and the last 8 bits are the local address.

This so-called “classful” architecture served the Internet for the next 12 years, during which time it grew from a small U.S.-based research network to a global academic network showing the first signs of commercial development.

Early Registration Models

In the 1980s, the American *National Science Foundation’s* (NSF’s) high-speed network, NSFNET, was connected to the ARPANET, a U.S. *Defense Advanced Research Projects Agency* (ARPA, now DARPA) wide-area network, which essentially formed the infrastructure that we now know as the Internet.

From these early days of the Internet, the task of assigning addresses was a necessary administrative duty, to ensure simply that no two networks would attempt to use the same network address in the Internet.

At first, the elementary task of maintaining a list of assigned network addresses was carried out voluntarily by Jon Postel, using (according to legend) a paper notebook.

As the Internet grew, and particularly as classful addressing was established, the administrative task grew accordingly. The IANA was established, and within it the Internet Registry (IR). But as the task of the IR outgrew Postel's notebook, it was passed to SRI International in Menlo Park, California, under a NSF contract, and was called the *Defense Data Network (DDN) Network Information Center (NIC)*.

During this time, under the classful address architecture, networks were allocated liberally and to any organization that fulfilled the simple request requirements. However, with the accelerating growth of the Internet during the late 1980s, two problems loomed: the rapid depletion of address space, due to the crude classful divisions; and the uncontrolled growth of the Internet routing table, due to unaggregated routing information.

Conservation vs. Aggregation

The problems of “three sizes fit all” highlight the basic dilemma of address space assignment: conservation versus aggregation. On the one hand, one wants to conserve the address space by assigning as little as possible; on the other hand, one wants to ease routing-table pressures by aggregating as many addresses as possible in one routing-table entry.

This can be illustrated by looking at a typical networking setup of the time. Within organizations having a single Internet connection, buildings, departments, or campuses would have their own local networks. Often the use of multiple networks was dictated by distance limitations inherent in the emerging local-area networking technologies, such as Ethernet.

These networks typically had to accommodate more than the 254 hosts addressable by a Class C address, but would rarely exceed 1000 hosts. Using pure classful addressing, one could either subdivide networks artificially to remain below the 254 host limit, or use a Class B address for each local network, possibly wasting more than 60,000 addresses in each. Whereas the latter solution is obviously wasteful in terms of address space, the former is obviously cumbersome. Less obviously, the former also puts an additional burden on the Internet routing system, because each of these networks would require a separate route propagated throughout the whole Internet.

This basic dilemma persists to this day. Assigning address space generously tends to reduce the routing-table size, but wastes address space. Assigning conservatively will waste less, but cause more stress for the routing system.

Subnetting

In order to address some of the problems of classful addressing, the technique of *subnetting* was invented. Described in RFC 791 in 1984, subnetting provided another level of addressing hierarchy by inserting a *subnet* part into the IP address between the network and local parts. Global routing remained the same using the *network* part of the address (Class A, B, or C) until traffic reached a router on the network identified by the network part of the address. This router, configured for subnetting, would interpret a statically configured number of bits from the local part of the address (the subnet part) to route the packet further among a set of similarly configured routers. When the packet reached a router connected to the destination subnet, the remaining bits of the local part would be used to determine the local address of the destination as usual. So, in the previous example, the organization could have used a Class B address with 6-bit subnetting, a setup that would allow for 62 networks of 1022 hosts each.

Subnetting nicely solved the routing-table problem, because now only one global routing-table entry was needed for the organization. It also helped address space conservation somewhat because it provided an obvious alternative to using many sparsely populated Class B networks.

Because the boundary between the subnet part and the local part of an address could not be determined from the address itself, this local knowledge needed to be configured into the routers. At first this was done by static configuration. Later, interior routing protocols carried that information. Refer to RFC 791 for numerous historically interesting case studies.

Supernetting

Within seven years, however, it was becoming clear that subnetting was no longer sufficient to keep up with Internet growth. RFC 1338 stated the problem:

“As the Internet has evolved and grown ... in recent years, it has become painfully evident that it is soon to face several serious scaling problems. These include:

1. Exhaustion of the Class-B network address space. One fundamental cause of this problem is the lack of a network class of a size that is appropriate for a midsized organization; Class C, with a maximum of 254 host addresses, is too small while Class B, which allows up to 65534 addresses, is too large to be widely allocated.
2. Growth of routing tables in Internet routers beyond the ability of current software (and people) to effectively manage.
3. Eventual exhaustion of the 32-bit IP address space.

It has become clear that the first two of these problems are likely to become critical within the next one to three years.”

The solution proposed was to extend the subnetting technique beyond the local organization, into the Internet itself. In other words, RFC 1338 proposed abolishing classful addressing, and replacing it with *supernetting*. The proposal was summarized as follows:

“The proposed solution is to hierarchically allocate future IP address assignment, by delegating control of segments of the IP address space to the various network service providers.”

CIDR

In 1993, the supernetting technique was published as a standards track RFC under the name *Classless Inter-Domain Routing* (CIDR), by which it is known and used today. Two main ingredients were necessary to make CIDR work: routing system changes and new address allocation and assignment procedures.

Under CIDR, routers could no longer determine the network part of an address from the address itself. This information now needed to be conveyed by Internet routing protocols. Fortunately, there was only one such protocol in widespread use at the time, and it was quickly extended by the major router vendor of the time. According to legend, the necessary extensions of the *Border Gateway Protocol* (BGP)-3 to BGP-4 were designed on a napkin, with all implementors of significant routing software present. The changes were implemented in a matter of days, but only much later described by the Internet standards track RFC 1654.

CIDR also required that forwarding decisions of routers be changed slightly. The network part of an address, now more generally called the *prefix*, can be of any length. This means that a router can have multiple valid routes covering a specific 32-bit destination address. Routers need to use the most specific of these routes—the *longest prefix*—when forwarding packets.

In addition to technical changes, the success of CIDR also relied on the development of administrative procedures to allocate and assign address space in such a way that routes could be aggregated as much as possible. Because the Internet was evolving toward the current state of arbitrarily interconnected networks of *Internet Service Providers* (ISPs), it was obvious that ISPs should play a role in address space distribution. In the new technique, ISPs would now, as much as possible, assign address space to their customers in contiguous blocks, which could be aggregated into single routes to the rest of the Internet.

Emergence of the RIRs:

Internationalization

While the engineering-driven need for topological address space assignment was becoming clear, there was also an emerging recognition that the administrative mechanisms of address space distribution needed further development. A central system just would not scale for numerous reasons, including:

- Sheer volume
- Distance from the address space consumers
- Lack of an appropriate global funding structure
- Lack of local community support

The need to change administrative procedures was formally recognized by August 1990, when the Internet Activities Board published a message it had sent to the U.S. Federal Networking Council, stating “it is timely to consider further delegation of assignment and registration authority on an international basis” (RFC 1174).

The increasing cultural diversity of the Internet also posed administrative challenges for the central IR. In October 1992, the *Internet Engineering Task Force* (IETF) published RFC 1366, which described the “growth of the Internet and its increasing globalization” and set out the basis for an evolution of the registry process, based on a regionally distributed registry model. This document stressed the need for a single registry to exist in each geographical region of the world (which would be of “continental dimensions”). Registries would be “unbiased and widely recognized by network providers and subscribers” within their region. Each registry would be charged with allocating remaining address space in a manner “compatible with potential address aggregation techniques” (or CIDR).

RIPE NCC

While in the United States the Government continued to support and fund registry functions, this was not the case in other parts of the world. In Europe, IP network operators cooperating in *Réseaux IP Européens* (RIPE) realized the need for professional coordination and registration functions. Establishment of the *RIPE Network Coordination Centre* (NCC) was proposed in the same month that RFC 1174 was published. The RIPE NCC was to “function as a ‘Delegated Registry’ for IP numbers in Europe, as anticipated and defined in RFC 1174” (RIPE-19).

Although consensus among IP network operators was quickly established, it took almost two years of organizing and fund-raising before the first RIR was fully operational in May 1992. The RIPE NCC was organized as a highly independent part of RARE, the organization of European research networks. It was to be funded by contributions from those networks, as well as a small number of emerging commercial networks. The RIPE NCC published its first regional address distribution policy in July 1992 (RIPE-65).

During the following months, European regional policies were refined and, for the first time, global guidelines were published as RFCs (RFC 1366, RFC 1466).

The RIPE NCC is presently organized as a membership association, performing the essential coordination and administration activities required by the RIPE community. Located in Amsterdam, Netherlands, the RIPE NCC service region incorporates 109 countries covering Europe, the Middle East, Central Asia, and African countries located north of the equator. The RIPE NCC currently consists of more than 2700 members. At the time of publication, RIPE NCC is performing the secretariat function for the *Address Supporting Organization* (ASO) of The *Internet Corporation for Assigned Names and Numbers* (ICANN). More information about RIPE NCC is available at <http://www.ripe.net>

APNIC

Asia Pacific Network Information Centre (APNIC), the second RIR, was established in Tokyo in 1993, as a pilot project of APCCIRN (Asia Pacific Coordination Committee for Intercontinental Research Networks, now *Asia Pacific Networking Group* [APNG]).

The project was an intended as a trial model for servicing the Internet addressing needs of national *Network Information Centres* (NICs) and other networks throughout the region.

After a successful ten-month trial period, APNIC was established as a permanent organization to serve the Asia Pacific region (which includes 62 economies from Central and South Asia to the Islands of Oceania and the Western Pacific).

Originally, APNIC relied on the support of networking organizations and national NICs. However, in 1996, APNIC implemented a tiered membership structure.

APNIC relocated to Brisbane, Australia, in mid-1998. It currently services approximately 700 member organizations, across 39 economies of the region. Within the APNIC membership, there are also five *National Internet Registries* (NIRs), in Japan, China, Taiwan, Korea, and Indonesia. The NIRs perform analogous functions to APNIC at a national level and together represent the interests of more than 500 additional organizations.

In 2000, APNIC hosted the secretariat functions of the ASO in its inaugural year. More information about APNIC is available at:

<http://www.apnic.net>

ARIN

In 1991, the contract to perform the IR function was awarded to Network Solutions, Inc. in Herndon, Virginia. This included the transition of services including IP address registration, domain name registration and support, *Autonomous System Number* (AS) registration, user registration, online information services, help-desk operations, and RFC and Internet-Draft archive and distribution services (RFC 1261).

With explosive Internet growth in the early 1990s, the U.S. Government and the NSF decided that network support for the commercial Internet should be separated from the U.S. Department of Defense. The NSF originated a project named InterNIC under a cooperative agreement with *Network Solutions, Inc.* (NSI) in 1993 to provide registration and allocation of domain names and IP address numbers for Internet users.

Over time, after lengthy consultation with the IANA, the IETF, RIPE NCC, APNIC, the NSF, and the *Federal Networking Council* (FNC), a further consensus was reached in the general Internet community to separate the management of domain names from the management of IP numbers. This consensus was based on the recognition that the stability of the Internet relies on the careful management of IP address space.

Following the examples of RIPE NCC and APNIC, it was recommended that management of IP address space then administered by the InterNIC should be under the control of, and administered by, those that use it, including ISPs, end-user organizations, corporate entities, universities, and individuals.

As a result, ARIN (*American Registry for Internet Numbers*) was established in December 1997, as an independent, nonprofit corporation, with a membership structure open to all interested entities or individuals.

ARIN is located in Chantilly, Virginia, United States. Its service region incorporates 70 countries, covering North America, South America, the Caribbean, and African countries located south of the equator. ARIN currently consists of more than 1500 members. Within the ARIN region, there are two national delegated registries, located in Mexico and Brazil.

Until now, ARIN has carried the responsibility for maintaining registration of resources allocated before the inception of the RIRs. However, a major project is now under way to transfer these legacy records to the relevant RIRs. More information about ARIN is available at:

<http://www.arin.net>

Emerging RIRs

The existing RIRs currently serve countries outside their core regions to provide global coverage; however, new RIRs are expected to emerge, necessitating changes to the existing service regions. Because the regions are defined on continental dimensions, the number of new RIRs will be low.

Currently, two groups have made significant progress in seeking to establish new RIRs. *AfriNIC* (for the Africa region) and *LACNIC* (for Latin America and the Caribbean) have each conducted public meetings, published documentation, and participated in the activities of the

existing RIRs. In recognition of the regional support they have so far obtained, each organization has been granted observer status at ICANN ASO meetings. The existing RIRs have also sought to provide as much assistance and support as possible to these emerging organizations.

More information about AfriNIC is available at;
<http://www.afrinic.org/>

More information about LACNIC is available at:
<http://lacnic.org/>

The RIR System:

Goals of the RIRs

RFC 2050, published in November 1996, represented a collaboration of the global Internet addressing community to describe a set of goals and guidelines for the RIRs. Although IANA was to retain ultimate responsibility for the entire address pool, RFC 2050 recognizes that RIRs operate under the consensus of their respective regional Internet community. This document, along with a history of RIR coordination, has helped to form the basis for a set of consistent global policies.

The three primary goals of the RIR system follow:

- *Conservation*: to ensure efficient use of a finite resource and to avoid service instabilities due to market distortions (such as stockpiling or other forms of manipulation);
- *Aggregation (routability)*: to assist in maintenance of Internet routing tables at a manageable size, by supporting CIDR techniques to ensure continued operational stability of the Internet;
- *Registration*: to provide a public registry documenting address space allocations and assignments, necessary to ensure uniqueness and provide information for Internet troubleshooting at all levels.

The Open Policy Framework

It was always recognized that these goals would often be in conflict with each other and with the interests of individuals and organizations. It was also recognized that legitimate regional interests could justify varying approaches in balancing these conflicts. Therefore, within the global framework, each regional community has always developed its own specific policies and procedures.

However, whereas the specific approaches may differ across the RIRs, all operate on a basic principle of open, transparent, consensus-based decision-making, following self-regulatory practices that exist elsewhere in the Internet and other industries. Furthermore, the RIRs all maintain not-for-profit cost-recovery systems and organizational structures that seek to be inclusive of all interested stakeholders.

The activities and services of each of the RIRs are defined, performed, discussed, and evaluated in open forums, whose participants are ultimately responsible for decision-making.

To facilitate broad participation, open policy meetings are hosted by RIRs regularly in each of the regions. Ongoing discussions are carried out on the public mailing lists of each RIR, which are open to both the RIR constituents and the broader community. The RIRs also participate actively in other Internet conferences and organizations and, importantly, each RIR has a strong tradition of participating in the public activities of the others.

A current example of the coordinated efforts of the RIRs is the Provisional IPv6 Assignment and Allocation Policy Document, a joint effort of the RIRs with the assistance of the IETF, The *Internet Architecture Board* (IAB), and the *Internet Engineering Steering Group* (IESG) to describe the allocation and assignment policies for the first release of IPv6 address numbers.

Also, the RIRs recently published the RIR Comparative Policy Overview, which is available at: <http://www.ripe.net/ripenc/mem-services/registration/rir-comp-matrix-rev.html>

These documents help illustrate that the well-established combination of bottom-up decision-making and global cooperation of the RIRs has created a stable, robust environment for Internet address management.

RIR Functions

The primary function of each RIR is to ensure the fair distribution and responsible management of IP addresses and the related numeric resources that are required for the stable and reliable operation of the Internet. In particular, the resources allocated, assigned, and registered by RIRs are Internet address numbers (IPv4 and IPv6) and AS numbers. RIRs are also responsible for maintaining the reverse delegation registrations of the parent blocks within their respective ranges.

Complementing their registry function, the RIRs have an important role in educating and informing their communities. The activities carried out by the individual RIRs vary, but include open policy meetings, training courses, seminars, outreach activities, statistical reporting, and research.

Additionally, a crucial role for the RIRs is to represent the interests of their communities by participating in global forums and providing support to other organizations involved in Internet addressing issues.

RIRs and The Global Internet Community:

Formation of ICANN and the ASO

The global Internet governance landscape began to undergo radical changes in mid-1998, with the publication of a U.S. Government white paper outlining the formation of a “not-for-profit corporation formed by private sector Internet stakeholders to administer policy for the Internet name and address system.” ICANN was formed later that year.

At the heart of the ICANN structure are “supporting organizations” that are formed to “assist, review and develop recommendations on Internet policy and structure” within specialized areas. In October 1999, the existing RIRs and ICANN jointly signed a *Memorandum of Understanding* (MoU) to establish the principles for forming and operating the *Address Supporting Organization* (ASO). It is intended that new RIRs will sign the MoU as they emerge.

Under the ASO MoU, the policy forums within each of the RIR regions continue to be responsible for development of regional IP address policy. In addition, each signatory RIR is responsible for electing three members to the ICANN *Address Council*.

The purpose of the Address Council, as described in the MoU, is to review and develop recommendations on issues related to IP address space, using the open processes that exist in the three regions; and to advise the ICANN Board on these matters. In addition, the Address Council is responsible for the appointment of three ICANN Directors to the ICANN Board.

RIR–ASO Coordination

Since the formation of the ASO, the RIRs have played an integral part in facilitating its activities. By joint agreement, the RIRs will share the ASO secretariat duties, including the hosting of the ASO Web site, on a revolving basis. APNIC provided these services in the ASO’s first year of operation, and RIPE NCC is currently performing this role.

The ASO Address Council holds monthly telephone conferences, which are attended by representatives of the RIRs (and emerging RIRs on a listener basis). In accordance with the MoU, the ASO also holds regular open meetings in conjunction with the open policy meetings of the RIRs.

RIRs and Industry Development

As noted previously, the RIRs maintain high levels of participation in the conferences and activities of other organizations. Similarly, they invite the participation of interested parties in their own activities.

The RIRs are active in many areas of new technology implementation (such as *General Packet Radio Service* [GPRS] and *Universal Telecommunications System* [UMTS] mobile telephony, IPv6, and cable and *Digital Subscriber Line* [xDSL]-based Internet services).

The established regional processes have proved both flexible and open enough to incorporate such new developments into policy formation. Industry representatives frequently join policy discussions, present at plenary sessions, and participate in working groups.

The RIRs pursue relationships with industry bodies, particularly those with representative and developmental functions, to facilitate industry convergence on open standards and policy processes.

Many diverse parties have legitimate interests in the allocation and registration of IP addresses, and the RIRs remain committed to participating with these parties to achieve a consensus among the Internet community on IP address allocation issues.

The Future of RIRs

In Internet time it can be easy to forget that eight years is actually not long. Since it was first proposed in 1990, the RIR system has evolved rapidly, enjoyed strong community support, and has been relatively free of the political wrangling that has characterized the registration systems of other Internet resources. Without doubt, this position is largely due to the early determination to provide accessible, open forums for the interested stakeholders in the various regions.

New technologies, such as GPRS, broadband services, and IPv6 may raise operational and policy challenges to the RIRs, yet at the same time they bring opportunities for increased global cooperation, in a context where distinct regional concerns are represented more effectively than ever before.

It is hoped that the emergence of new RIRs will only serve to expand and enhance the inclusive nature of RIR activities.

References

- [1] Clark, D., and Cohen, D., "A Proposal for Addressing and Routing in the Internet," IEN 46, June 1978.
- [2] Postel, J., "Assigned Numbers," RFC 790, September 1981.
- [3] Information Sciences Institute, "Internet Protocol, DARPA Internet Program, Protocol Specification," RFC 791, September 1981.
- [4] Cerf, V., "IAB Recommended Policy on Distributing Internet Identifier Assignment and IAB Recommended Policy Change to Internet 'Connected' Status," RFC 1174, August 1990.
- [5] Williamson, S., and Nobile, L., "Transition of NIC Services," RFC 1261, September 1991.
- [6] Fuller, V., Li, T., Yu, J., and Varadhan, K., "Supernetting: An Address Assignment and Aggregation Strategy," RFC 1338, June 1992.
- [7] Gerich, E., "Guidelines for Management of IP Address Space," RFC 1366, October 1992.
- [8] Gerich, E., "Guidelines for Management of IP Address Space," RFC 1466, May 1993.
- [9] Rekhter, Y., and Li, T., "A Border Gateway Protocol 4 (BGP-4)," RFC 1654, July 1994.
- [10] Hubbard, K., Kesters, M., Conrad, D., Karrenberg, D., and Postel, J., "Internet Registry IP Guidelines," RFC 2050, November 1996.
- [11] Blokzijl, R., Devillers, Y., Karrenberg, D., and Volk, R., "RIPE Network Coordination Center," RIPE-19, September 1990.
- [12] Terpstra, M., "RIPE NCC Internet Numbers Registration Procedures," RIPE-65, July 1992.

DANIEL KARRENBORG has helped to build the European Internet since the early 1980s. As one of the founding members of the German UNIX Users Group, he has been involved in the setting up of EUnet, a pan-European cooperative network providing electronic mail and news to businesses and academic institutions all over Europe. While at CWI in Amsterdam, Karrenberg helped to expand this network and convert it to a fully IP-based service. During this time he created a whois database of operational contacts, which was the nucleus of the current RIPE database. Karrenberg is one of the founders of RIPE, the IP coordination body for Europe and surrounding areas. In 1992 he was asked to set up the RIPE NCC, the first regional Internet registry providing IP numbers to thousands of Internet service providers in more than 90 countries. Karrenberg led the RIPE NCC until 1999, when it had an international staff of 59 with more than 20 nationalities; he currently helps to develop new RIPE NCC services. Recently his contributions have been recognized by the Internet Society with its *Jon Postel Service Award*. Karrenberg's current interests include measurements of Internet performance and routing as well as security within the Internet infrastructure. In general he likes building new and interesting things. Mr. Karrenberg holds an MSc in computer science from Dortmund University. E-mail: **Daniel.Karrenberg@ripe.net**

GERARD ROSS holds a BA and LLB from University of Queensland and a Grad.Dip. (Communication) from Queensland Institute of Technology. He was employed as the technical writer at APNIC in 1998 and has been involved in the development and drafting of several major policy documents both in the APNIC region and as part of coordinated global RIR activities. He was the ASO webmaster in its inaugural year. He is currently the APNIC Documentation Manager. E-mail: **gerard@apnic.net**

PAUL WILSON has been Director-General of APNIC since August 1998. Previously, he was a founding staff member and subsequently Chief Executive Officer at Pegasus Networks, the first private ISP in Australia. Over an eight-year period he worked as a consultant to the United Nations and other international agencies on Internet projects in many countries. Since 1994, he has worked with the International Development Research Centre (IDRC) on its Pan-Asia Networking (PAN) Programme, supporting projects in Mongolia, Vietnam, Cambodia, Maldives, Nepal, Bhutan, PNG, and China. He continues to serve as a member of the PAN Research and Development Grants Committee. E-mail: **pwilson@apnic.net**

LESLIE NOBILE received her B.A. from the American University in Washington, D.C. She has over 15 years of experience in the Internet field, and has been involved with the Internet Registry system since 1991. Prior to that, she held various technical management positions while working under a U.S. Government contract that supported the engineering and implementation of the Defense Data Network, a high-speed data network that evolved from the ARPANET. Her experience with the Registry system began in 1991 working as one of the Operations managers who transitioned the Internet Network Information Center (NIC) from SRI to Network Solutions, Inc. She remained a registration services manager with the DDN/DoD NIC until August 2000, when she became Director of Registration Services at the American Registry for Internet Numbers (ARIN). She has been a contributing author to RFCs, Internet Society (ISOC) articles, and various other industry publications and has been actively involved in the global coordination of Internet addressing policy. Her e-mail address is **leslie@arin.net**

Book Reviews

Web Caching *Web Caching* by Duane Wessels, ISBN 1-56592-536-X, O'Reilly, June 2001.

It's always a pleasure to read a technical book written by someone who has not just studied the topic, but has been so involved that he has spent years living and breathing the subject. Such books do more than just describe the technology, because they are invariably able to add a dimension of deeper insight and interest, and in so doing, bring the topic to life for the reader. Duane Wessel's experiences in the Harvest project, and then as self-confessed "Chief Procrastinator" in the *Squid* Web cache project, certainly place him in the category of an author who has lived the topic. The outcome is a well-researched and very readable book on the topic of Web caching.

Web Caching

Web caching has been an integral part of the architecture of the World Wide Web since its inception, and is now a broad topic encompassing a range of approaches, a range of technologies, and a range of deployment issues for the end consumer, the content publisher, and the service provider intermediaries. The book starts with a clear introduction that outlines the elements of the architecture of the Web, and describes the terminology used within the book. This section also provides a basic introduction to the operation of the *Hypertext Transfer Protocol* (HTTP). This section also describes the various forms of Web caches that are in use today.

The way in which a cache interprets the directives at the header of a delivered Web object is described in some detail. I learned something unexpected here, in that a Web object that includes a directive of the form "Cache-control: no-cache" is defined in RFC 2616 as allowing a cache to store a copy of the object and use it, subject to revalidation, for subsequent requests. It seems that if you really want the object not to be stored in a cache, then "no-store" is what you are after, because "no-cache" allows the object to be cached! As well as describing the definition of the cache control directives, this section provides a clear explanation of how document ageing is defined, and when a cache server determines that a cached object should be checked against the original to ensure that the cached copy remains a faithful reproduction.

Caching has its champions and its detractors, and the book attempts to present both perspectives in a balanced fashion. On the positive side, caching is seen as an effective way to improve the performance of the delivery of Web-based services, and to relieve network and server load. The claim is made here that a large busy cache can achieve a hit ratio of some 70 percent. Don't get too enthusiastic, however, because a more common achieved ratio is somewhere between 30 and 40 percent.

On the negative side is the ever-present issue of accuracy of the cache, the inability for a content provider to track content access, and the issue of integrity of the cache in the face of service attacks that are directed to the cached copy of the content.

The Politics of Caching

This section of the book intrigued me, because it is certainly rare to see a technical book address the various social implications of the technology. The study includes the issues of privacy, request blocking, copyright control, content integrity, cache busting, and the modifications to the trust model in the presence of cache intermediaries. The book exposes the tension between the content provider, the user, and the service provider. The content provider would generally like to exercise some control over tracking who is accessing the content and how each client uses the content and how they navigate through the Web site. The user is interested in efficiency of content delivery, and also has to place a high level of trust in the integrity of the content-delivery system. The service provider is also interested in rapid delivery of content, as well as managing network load. Third parties, such as regulatory or law-enforcement bodies, may be interested in ensuring that the content originator is unambiguously traceable, and that various regulations with respect to content are enforced by content originators and service providers.

Practical Advice

From this overview, the book moves onto more practical topics, and first describes how to configure browsers to take advantage of caches. It also covers how various proxy auto-configurators work. The topic that has generated some attention is that of *interception caching*, where a user's Web-browser commands are intercepted by a provider cache without the direct knowledge of the user of the user's browser. The techniques of implementing such interception caches are described, including a description of the operation of the *Web Cache Coordination Protocol* (WCCP), policy routing, and firewall interception. Interception caching, or transparent caching, is a topic that has generated its fair share of controversy in the past, and the book does take the time to clearly describe the issues associated with this caching approach.

The other topic covered under the general topic of practical advice is advice to server operators and content providers on how to make servers and content work in a predictable fashion with caches, describing which HTTP reply headers affect cacheability. This section provides advice on how to build a cache-friendly Web site, and motivates this with reasons why a content provider would want to ensure that content is readily cacheable. This includes some practical advice on how a content provider can still receive hit counts and site navigation information while still allowing the content of a site to be cached.

Fun with Caches—Cache Hierarchies and Clusters

Although caches can operate in a standalone configuration, it is possible to interconnect caches so that a cache will refer to another cache in the event of a cache miss, rather than directly refer to the origin server. I gather that the author is not overly keen on such an approach, given that the arguments against such configurations consume five times as much space as the arguments in favor! The alternative to a strict hierarchy is a set of cooperating peer caches, together with an intercache protocol to allow a cache to efficiently query its peers for an object. The book describes the *Internet Cache Protocol* (ICP), the *Cache Array Routing Protocol* (CARP), which is pointed out to be an algorithm, not a protocol, despite its name, the *Hypertext Caching Protocol* (HTCP), and *Cache Digests*. The scenarios where each approach would be preferred is a helpful addition to this section. Cache clusters are also described; if I have a criticism of the book, it is that this section is too terse—I was looking for more details of cache-balancing and content-distribution techniques.

Cache Operation

The final section of the book looks at the tasks associated with designing, benchmarking, and operating cache servers. How much disk space is enough for a cache? How much memory? Where should the caches be placed in the network? What aspects of the cache operation should you monitor? And if you are considering purchasing caches, what aspects of the cache should you carefully examine?

Conclusion

This is not a book about how to build a cache, although if you are considering doing that it's a good place to start your research. Nor is it a book about every detail on how to operate a cache. But if you are operating a cache, it will be useful. Although it's not a book about how to operate a Web server, if you are operating a Web server, then caches will attempt to store your content, and this book will help you configure your server to interoperate predictably with caches.

The Web is a large part of today's Internet, and Web caches can make the Web faster, more efficient, and more resilient. If you want to understand how caches work and understand how you can use caches to improve the user's experience rather than making things worse, then this book is essential reading.

—Geoff Huston
gih@telstra.net

IPSec *IPSec: The New Security Standard for the Internet, Intranets, and Virtual Private Networks*, by Naganand Doraswamy and Dan Harkins, ISBN 0-13-011898-2, 1999, Prentice Hall PTR Web Infrastructure series. <http://www.phptr.com>

We all know that Internet security is a major concern. Evolving technologies such as *Virtual Private Networks* (VPNs) are making it easier to deploy secure networks at low costs. VPN technology is based upon encryption techniques that make use of different algorithms. Most of these algorithms are specified in the form of *Requests for Comments* (RFCs). Though RFCs provide the minute details, they are not exactly lively reading. This is where the *IP Security* (IPSec) book comes in handy. The authors have done their best to explain IPSec technology in layman's language, although one encounters a lot of technical jargon in this book.

Organization

The book is divided into three parts. Part I gives a history of cryptography and techniques and cryptographic tools, and overviews of TCP/IP and IPSec. Authentication methods such as *Public Key Infrastructure* (PSI), RSA, and DSA are discussed. Key exchange methods such as Diffie-Hellman and RSA Key Exchange are discussed, along with their advantages and disadvantages. IPSec architecture is explored in the IP Security Overview section, which describes the security services provided by IPSec, how packets are constructed and processed, and the interaction of IPSec processing with policy. IPSec protocols—*Authentication Header* (AH) and *Encapsulation Security Payload* (ESP)—are the basic ingredients of the IPSec stack to provide security. Both AH and ESP can be operated in either the transport mode or tunnel mode. Part II offers a detailed analysis of IPSec, the different modes, IPSec implementation, the ESP, AH, and the *Internet Key Exchange* (IKE). The authors do a good job of describing the IPSec road map, which defines how various components within IPSec interact with each other. Detailed packet formats of different IPSec formats are discussed in Chapter 4. ESP, AH, and IKE are discussed in depth in Chapters 5 through 7. Part III deals with most of the deployment issues concerned with IPSec, as well as policy definition, policy management, implementation architecture, and end-to-end security are discussed in this section. Chapter 11 discusses the future of IPSec and what it means to the world of security. Though IPSec may be thought of as a totally secure method of communication, it has its conflicts when it comes to *Network Address Translation* (NAT), multicasting, and key management in a multicast environment.

Prerequisites

Although the authors have done a good job delivering the IPSec concept, understanding this text requires more than basic computer and communication concepts. One should understand hacking and different types of Internet attacks. OSI layer details and packet-level understanding of every layer within the OSI model is a must.

—Manohar Chandrashekar, WorldCom Inc

mchandra@wcom.com

Letters to the Editor

ICANN Mr. Jacobsen,

I very much enjoy the *Internet Protocol Journal* and put it at the top of my reading stack as soon as it is received. In particular, I enjoy the standards and high technical detail and view it as a safe place from overt commercial advertisement and politics.

That is why I was disappointed by the article from Mr. Lynn. My opinion of ICANN is that it is undemocratic in any tradition, uninterested in experimentation, and uninterested in outside views. I took offense at his continued use of the phrase “public trust” and interpreted the article as propaganda. Further, I found the technical content of the article to be zero.

On the other hand, William Stallings article on MPLS was exactly the kind of article I’ve come to enjoy. I wasn’t familiar with MPLS and the article helped me understand the concepts, vocabulary, and high-level issues. I hope that “MPLS” serves as a model of the articles in future IPJ issues.

I keep back issues of IPJ in a binder and continue to hope you uncover more articles like “The Social Life of Routers.” My copy of Mr. Krebs article has notes in all the margins—I was excited—but it was a twist on something that I thought I knew and he exposed a different design vocabulary by making an unexpected comparison.

I apologize for complaining about something that is a gift from Cisco; I do understand how crass that is. I hope that you will interpret my note in a complementary manner: I’ve come to respect the journal and found that it fits an unfilled niche in my reading.

—Brent D. Stewart, Global Knowledge
<brent@stewart.hickory.nc.us>

Brent,

I appreciate your feedback, as I am sure Mr. Lynn will if you send it to him. The article was, after all, published for public comment.

ICANN has unfortunately tended to polarize people and has become a forum in which a certain amount of politics is played out. I don’t think this is entirely ICANN (the board)’s fault. What was set up as an organization to take over the work of one man—the late Jon Postel, is seen by some as an opportunity for “Internet Governance” and “world-wide electronic democracy.”

Having watched the ICANN process since its beginnings in 1998, I would say that Mr. Lynn’s version of history is pretty much on target. When the IANA was in the hands of Jon Postel, it most certainly was a “public trust” (a limited resource to say the least), and if ICANN does not take that responsibility seriously, it certainly will have failed.

However, I do not think this is the case. Yes, ICANN is now a fairly large and slow moving machinery, and I would have liked to see more new domains deployed sooner, but to some extent the slowness is caused by the structure of Supporting Organizations as much as it is by the board itself. There is a lot to sort out, a lot to comment on, and *many* divergent views are indeed being expressed in all kinds of ICANN forums, including the public meetings. So, I cannot agree that ICANN is “uninterested in outside views.” A perfect democracy it is not, nor was it ever intended to be, and yes, some of the topics on the agenda such as the *Uniform Dispute Resolution Process* (UDRP) are indeed non-technical. But it is not as if ICANN had much choice in that particular matter. (Although some would argue that it could be moved outside the ICANN process.)

Being part of the ICANN process, through e-mail discussion, public meetings or through the Supporting Organizations is not difficult. Nor do I think that ICANN ignores any of the feedback it gets.

Back to the article. No, it was not particularly technical, but if you read IPJ’s Call for Papers you will see that it mentions “Legal, policy and regulatory topics...” Also, in the wake of September 11, I thought it was important to provide some background on the thinking of ICANN, and why they chose to refocus the most recent meeting on security etc. IPJ, by the way, also encourages the occasional “Opinion Piece,” although the article by Mr. Lynn was not intended as such. The issue of alternate roots is indeed a matter of debate, and while the IAB has already expressed its view, I appreciate that there might be other (valid) ones.

In any case, thank you for taking the time to write. I certainly don’t intend to steer IPJ away from topics such as MPLS and I hope that the occasional policy or even opinion piece won’t steer you away from IPJ.

—Ole Jacobsen, Editor and Publisher <ole@cisco.com>

MPLS Ole,

William Stallings otherwise-excellent article on MPLS in the *Internet Protocol Journal* Vol. 4, No. 3 had a serious error in it with respect to Virtual Private Networks (VPNs). He said that MPLS is an efficient mechanism for supporting VPNs and that MPLS provides security; neither is true.

As the rest of the article shows, MPLS provides a transport tunnel for IP packets, meaning that it helps create virtual networks. However, there is no privacy on those virtual networks, so it is inappropriate and probably dangerous to call MPLS tunnels virtual private networks.

To most Internet users, security means preventing snooping of sensitive traffic, preventing malicious changes to content, or both. MPLS does not provide either service. Instead of relying on insecure MPLS, users who want secure tunnels use systems that employ the IPsec protocol.

Many dozens of vendors supply IPsec systems appropriate for everything from tiny home offices to gigantic telco central switches, all with the same high security. Although the article showed that MPLS has many valuable features, IPJ readers should not fall into the trap of thinking that VPN support or security are MPLS features.

—Paul Hoffman, Director, VPN Consortium
<paul.hoffman@vpnc.org>

Ed: We presented this letter to a panel of experts, and here are some samples of the responses we received:

The term “VPN” has been used in many different contexts. I saw a group once call a VLAN a “VPN” as well. I honestly couldn’t say that they were incorrect. It may be appropriate to say that there are IPsec VPNs and that there are MPLS VPNs, but I have a problem calling one “right” and another “wrong” simply because of some perceived, implied definition of the security level that should be provided by a “VPN.” Most people support the notion that an MPLS VPN provides about as much “security” as a Frame Relay link. This amount of “security” in a VPN is acceptable to many people.

—Chris Lowick, Cisco Systems <clonvick@cisco.com>

We have different views on security, I’m sure. One view is that a secure private network: a) ensures that a third party cannot impose a condition on the network such that a customer’s traffic is directed to another customer b) ensures that a third party cannot inject traffic into a customer’s private network, c) a third party cannot alter customer traffic and d) a third party cannot discern that communications is taking place between two parts of a private network.

MPLS uses the same mechanisms as X.25, ATM and Frame, and has similar properties—the objectives above can be met with adequate confidence as long as the network is carefully configured and managed.

Edge to edge IPsec has a different set of security principles—the basic mode of operation is that such networks may be subject to attacks that redirect customer’s traffic to third party sites, and allow third parties to inject traffic into the VPN, and allow a third party to discern that communications is taking place within a private context. The essential attribute of edge to edge IPsec is that the encryption is intended to ensure that leakage can be identified: foreign injected traffic or altered traffic can be identified and rejected and leaking traffic cannot be decoded.

Both approaches have vulnerabilities and weaknesses. The first approach places trust in the integrity of the host platform. The second approach is prone to various forms of DOS attacks and traffic profiling.

But I would not concur with a view that labels the MPLS approach as inefficient or insecure, nor would I label X.25 networks, ATM or Frame as *intrinsically* inefficient and insecure. There are insecure operating practices and there are cautious operating practices.

IPSec networks have similar issues—relating particularly to the vulnerabilities of third party disruption and profiling eavesdropping.

So it's not that I believe that all MPLS networks are well designed and well operated—on the contrary! But as an architectural approach I am not able to agree with a comment that appears to condemn MPLS as *intrinsically* a poor choice for a VPN host technology.

So if the comment is that the article provides the impression that MPLS is such a robust technology that it creates secure private network applications such as VPNs, and appears to make this assertion so strongly that it gives the impression that this outcome occurs irrespective of MPLS network design and operating practices, and that this impression is ill-founded, then I would agree entirely with Mr. Hoffman. Secure networks, or at least robust networks, are a result of careful choice of technologies coupled with careful design and careful operation.

—Geoff Huston, Telstra <gih@telstra.net>

Ed.: We forwarded these comments to Mr. Hoffman, and he responded:

Geoff believes that it a network that does not prevent an active attacker from seeing or modifying traffic, and does not prevent a passive attacker from seeing packets, is secure and private; I do not. The fact that MPLS restricts the flow of traffic to a particular defined network is sufficient for him; it is not for me, given the fact that an attacker breaking into any node on that defined network can compromise the privacy and integrity of the traffic.

It is typical for ISPs to not want to do the work of actually securing the traffic they say they have put in a VPN by using IPsec. That work is not cheap, and takes more management than vanilla MPLS, but it is the only way to really secure the data. I am absolutely not saying that the IPsec community is without blame here: we have a tendency to ignore the valuable features of MPLS and have done almost nothing to make it easier to intelligently tunnel IPsec in MPLS (we also pretty much stonewalled the IPsec under L2TP work that is now finally standardized). But our lack of openness doesn't make MPLS a VPN technology.

—Paul Hoffman, Director, VPN Consortium
<paul.hoffman@vpnc.org>

Ed.: We would love to hear from you. Please send your letters to:
ipj@cisco.com

Fragments

ACM Assembles Security and Privacy Panel

Prompted by increased public concerns about personal privacy and the security of networked information systems, the *Association for Computing Machinery* (ACM) has announced the formation of a new *Advisory Committee on Security and Privacy* (ACSP). Led by Peter Neumann and Eugene H. Spafford, the ACSP brings together a dozen leaders and innovators in the field of privacy and information assurance to serve as a powerful resource for the ACM community and the public at large.

Comprising experts from research, industry, academia, and government, the diverse group represents a wide range of viewpoints. Commenting on the formation of the ACSP, Co-Chair Peter Neumann noted, “The ACSP will provide timely and accurate assessments of situations relating to information security that are otherwise clouded by confusion, uncertainty, and often, misinformation.”

Added ACSP Co-Chair Gene Spafford, “Until recently, computing professionals have been primarily concerned with making computers work consistently, cheaply, and effectively. Now it is critical that we also bring expertise to bear on how computers can be made to operate safely, keep information resources secure from attack, and protect privacy.”

The ACSP consists of 12 distinguished members with expertise in information security and assurance, privacy, cybercrime, and allied fields. The group will coordinate with other ACM Committees, including the *U.S. ACM Committee on Public Policy* (USACM) and ACM Law Committee, to provide objective advice to the computing community, the public at large, and to policy-makers. ACSP is expected to provide statements and testimony on information security and privacy issues, as well as undertaking studies of related topics. For more information about the ACSP, see the web site at:

<http://www.acm.org/usacm/ACSP/homepage.htm>

Members of the ACSP (affiliations provided for identification purposes only) are:

Steve Bellovin (AT&T Labs Research)
Matthew Blaze (AT&T Labs Research)
David Clark (MIT)
Dorothy Denning (Georgetown University)
Ed Felten (Princeton University)
David Farber (University of Pennsylvania)
Susan Landau (Sun Microsystems)
Robert Morris (Dartmouth College)
Peter Neumann (SRI International)
Fred Schneider (Cornell University)
Eugene H. Spafford (Purdue University CERIAS)
Willis Ware (RAND Corporation)

For more information, see ACM’s Web site at: **<http://www.acm.org>**

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Technology
WorldCom, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 2001 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-10/5
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSR STD
U.S. Postage
PAID
Cisco Systems, Inc.

The Internet Protocol *Journal*

March 2002

Volume 5, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
IEEE 802.11	2
Code Signing.....	14
Book Review.....	27
Call for Papers	30
Fragments	31

FROM THE EDITOR

Major Internet events such as the IETF meetings, the Regional Internet Registry meetings, APRICOT, SIGCOMM, and NetWorld+Interop to name a few, all provide Internet access for attendees. Commonly referred to as the “Terminal Room,” these facilities have evolved into complex high-speed networks with redundant paths, IPv6 routing, multicast, and more. In the last five years or so, these networks have also been providing wireless access using various flavors of the IEEE 802.11 standard. As I write this, I am sitting in the lobby of the Minneapolis Hilton Hotel, where the 53rd IETF meeting is being held. The lobby area and two floors of meeting rooms have IEEE 802.11 coverage, and a directional high-gain antenna provides access in the pub across the street. Wireless Internet computing is a reality, at least when you have a large gathering of engineers such as an IETF meeting. In our first article, Edgar Danielyan takes a closer look at this technology, its applications and evolution.

More and more software is being distributed via the Internet rather than through the use of conventional media such as CD ROMs or floppy disks. Downloading software via the Internet is very convenient, especially if you have reasonably high bandwidth. However, with this convenience comes a certain risk that you may be receiving a modified copy of the software, perhaps one that contains a virus. Code signing is a method wherein software is cryptographically signed and later verified. Eric Fleischman explains the details of code signing.

I should have known better than to announce the imminent availability of our online subscription system in the previous issue. We are working on it, but it isn't ready yet, so please continue to send your subscription requests and updates to: ipj@cisco.com

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

IEEE 802.11

by Edgar Danielyan

Introduced in 1997, the IEEE Standard 802.11 for wireless local-area networks has seen modifications and improvements in the past years and is promising a brighter wireless future, so yearned for by many of us. However, during its lifetime, the standard also has had a few setbacks, which are reminders that nothing is perfect in this world, much less in networking. This article provides a brief but comprehensive introduction to IEEE 802.11 wireless networking, its present and future, and highlights some of its security, performance, and safety aspects.

IEEE 802.11

The initial IEEE Standard 802.11 was published by the *Institute of Electrical and Electronics Engineers* (IEEE) in 1997. That standard is known as IEEE 802.11-1997 and is now updated by the current standard, IEEE 802.11-1999. The current standard has also been accepted as an American national standard by the *American National Standards Institute* (ANSI) and has been adopted by the *International Organization for Standardization* (ISO) as ISO/IEC 8802-11:1999. The completion of IEEE 802.11 in 1997 set in motion the development of standards-based wireless LAN networking. The 1997 standard specified a bandwidth of 2 Mbps, with fallback to 1 Mbps in hostile (noisy) environments with *Direct Sequence Spread Spectrum* (DSSS) modulation, and bandwidth of 1 Mbps with *Frequency Hopping Spread Spectrum* (FHSS) modulation, with possible 2-Mbps operation in friendly (noiseless) environments. Both methods operate in the unlicensed 2.4-GHz band. What is less known about IEEE 802.11 is that it also defines a baseband infrared medium, in addition to the DSSS and FHSS radio specifications, although its usefulness seems somewhat limited. There are also several task groups inside the 802.11 working group itself that work on substandards of 802.11:

- 802.11D: Additional Regulatory Domains
- 802.11E: Quality of Service (QoS)
- 802.11F: Inter-Access Point Protocol (IAPP)
- 802.11G: Higher data rates at 2.4 GHz
- 802.11H: Dynamic Channel Selection and Transmission Power Control
- 802.11i: Authentication and Security

The IEEE 802 group has an official Web site at www.ieee802.org, and IEEE 802.11 has an official Web site at www.ieee802.org/11/.

DSSS

Direct Sequence Spread Spectrum (DSSS) is one of the modulation techniques provided for by the IEEE 802.11 and the one chosen by the 802.11 Working Group for the widely used IEEE 802.11b devices. DSSS modulation is governed in the United States by FCC Regulation 15.247 and in Europe by ETSI Regulations 300-328. DSSS in IEEE 802.11 uses *Differential Binary Phase Shift Keying* (DBPSK) for 1 Mbps, and *Differential Quadrature Phase Shift Keying* (DQPSK) for 2 Mbps. The *Higher-Rate DSSS* (DSSS/HR) defined in IEEE 802.11b uses *Complementary Code Keying* (CCK) as its modulation scheme and provides 5.5- and 11-Mbps data rates. Because of their compatibility, all three modulation schemes can coexist using the rate-switching procedures defined in the IEEE 802.11. The *Orthogonal Frequency Division Multiplexing* (OFDM) used by the IEEE 802.11a is regulated in the United States by Title 47 Section 15.407 of the U.S. *Code of Federal Regulation* (CFR). IEEE 802.11a uses a system of 52 subcarriers modulated by BPSK or QPSK and 16-quadrature amplitude modulation. It also uses *forward error correction* (FEC) coding, also used by the Digital Video Broadcasting (DVB) standard with coding rates of 1/2, 2/3, and 3/4.

FHSS

Although specified by the original IEEE 802.11, *Frequency Hopping Spread Spectrum* (FHSS) modulation is not favored by vendors and, it seems, the 802.11 working group itself. DSSS has won the battle—very few vendors support 802.11/FHSS, and further developments with 802.11 use DSSS. Some have expressed ideas that frequency hopping in FHSS may contribute to the security of 802.11, but these are invalid expectations—the hopping codes used by FHSS are specified by the standard and are available to anyone, thus making the expectation of security through FHSS unreasonable.

Two supplements to the IEEE 802.11-1999, known as IEEE 802.11a and IEEE 802.11b, brought considerable changes and improvements to the IEEE 802.11-1999 standard.

IEEE 802.11a

IEEE 802.11a specifies a high-speed physical layer operating in the 5-GHz unlicensed band utilizing a complex coding technique known as OFDM. The data rates specified by IEEE 802.11a are 6, 9, 12, 18, 24, 36, 48, and 54 Mbps, with support for 6, 12, and 24 Mbps as a mandatory requirement. IEEE 802.11a is seen by some in the industry as the future of IEEE 802.11. Some products already implement the IEEE 802.11a, such as the chip from Atheros (www.atheros.com) and a PCMCIA/CardBus adapter from Card Access Inc (www.cardaccess-inc.com) based on it. However, 802.11a is not without disadvantages. The increased bandwidth of IEEE 802.11a results in a shorter operation range.

Additionally, because of the protocol overhead and interference/error correction, the real bandwidth may be considerably less than the nominal. New surveys and installation will also be required in many cases; the underlying infrastructure will also be more expensive because of the shorter operation range (about 1/3 of 802.11b) and higher density of *base stations* (also known as *access points*).

IEEE 802.11b

Probably the most widely implemented and used wireless LAN technology today, IEEE 802.11b specifies 5.5- and 11-Mbps data rates (in addition to the already specified 1 and 2 Mbps), but operates in the original 2.4-GHz band also using DSSS modulation. Most currently selling IEEE 802.11 products implement IEEE 802.1b. IEEE 802.11b-compliant devices can operate at 1, 2, 5.5, and 11 Mbps.

It is important to note that both incarnations of IEEE 802.11 use the same *Media Access Control* (MAC) protocol, *Carrier Sense Multiple Access with Collision Avoidance* (CSMA/CA); therefore, these modifications affect only the physical layer (PHY layer in IEEE parlance) of the standard. The 1/2- and 5.5/11-Mbps DSSS (IEEE 802.11b) networks can coexist, enabling a painless transition to IEEE 802.11b (High Rate) at 11 Mbps. Eleven to fourteen radio channels are available for use with IEEE 802.11b in the 2.4-GHz band, depending on the local legal and administrative restrictions.

Distance, Power, and Speed Issues

It is obvious that all three of these parameters of wireless systems are interconnected. However, as with other radio-based technologies, the external conditions (such as the line of sight in case of outdoor use) greatly affect the operation of IEEE 802.11 devices.

Antennae

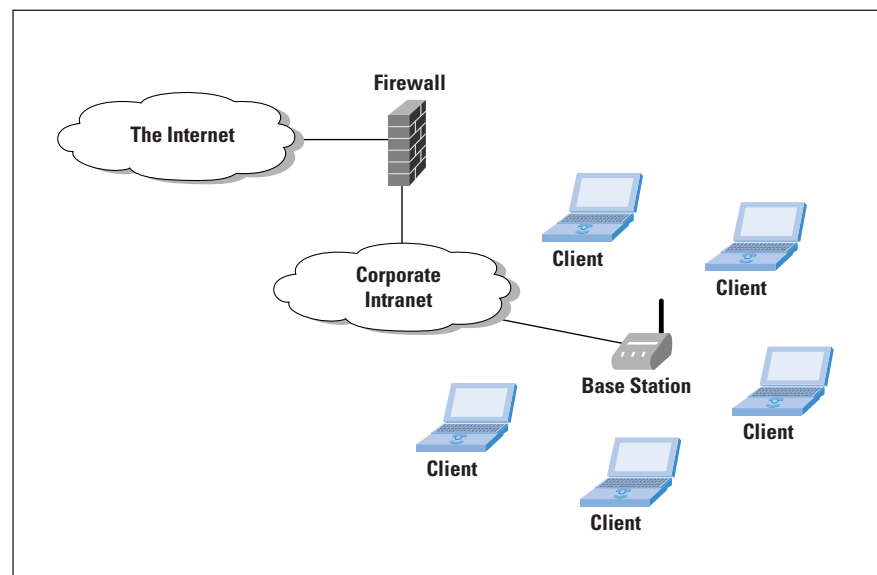
Antennae used with IEEE 802.11b devices may be grouped into two categories: *omnidirectional* and *point-to-point*. Obviously, omnidirectional antennae are the easiest to use, because they do not require positioning. Omnidirectional antennae are used in most base stations, as well as in most access cards. However, because of their nature, omnidirectional antennae do not work well over longer distances, unless used with external amplifiers; and these are not always legal or appropriate to use. Directional, or point-to-point antennae, on the other hand, require careful positioning and are used outdoors. Although the typical range for an omnidirectional antenna system is 150 ft (45m), configurations with high-gain directional antennae can work on distances up to 25 miles (about 40 km). In localities where amplifiers are allowed, the maximum distance may be considerably increased and is limited only by the line of sight.

Among other factors affecting the operational range of IEEE 802.11b devices are the base-station placement (when used in the infrastructure mode) and radio interference. As mentioned earlier, IEEE 802.11b devices will auto-configure for the highest possible speed and fall back to lower speeds when circumstances so require.

Performance Issues

Aside from obvious factors that affect performance (such as antennae, distance, radio interference) there are numerous other, more subtle issues. In the infrastructure mode, when all devices have to register with the base station(s), the load on the base station(s) increases with the number of clients and may reach a point when the performance reaches unacceptable lows. For example, Apple's AirPort Base Station (Version 2) can support up to 50 simultaneous clients. However, the actual performance of the whole system also depends on the kind of traffic. In particular, isochronous traffic (time-sensitive traffic, such as some types of video, audio, and telemetry), as well as multicast traffic, are particularly taxing for IEEE 802.11 networks and are better kept off the wireless LAN. However, several groups are currently working on extensions to 802.11 to provide for such kinds of traffic in a future version of the standard.

Figure 1: Typical IEEE 802.11 Configuration in Infrastructure Mode



IEEE 802.11 Base Stations and Clients

All IEEE 802.11 devices can be grouped into one of two groups: base stations or clients. Base stations can function as clients; however, not all clients can function as base stations. The reason for this is that base stations are required to provide certain network services to clients (association, distribution, integration, reassociation, and so on) that not all client hardware, firmware, or software can or intended to provide.

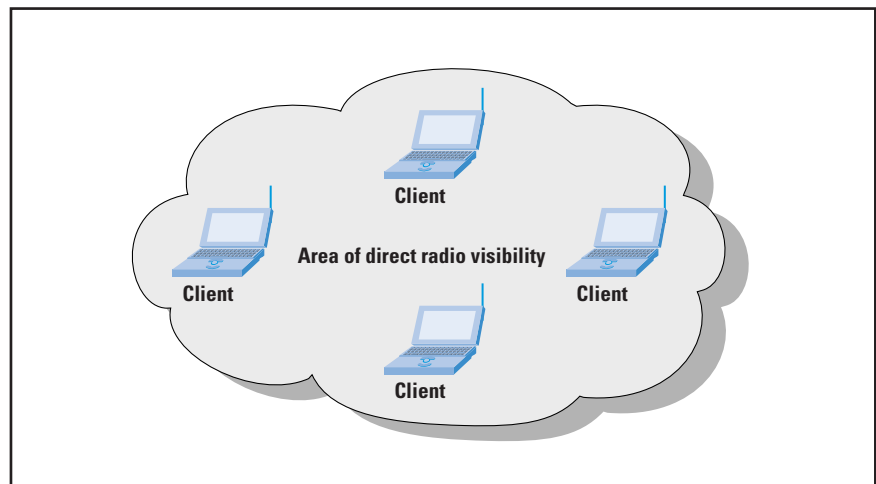
These considerations apply when the infrastructure mode of IEEE 802.11 is deployed. In *ad hoc* networks, where there are no base stations, all clients communicate directly with each other, reminiscent of a traditional shared Ethernet network, with all nodes sharing equal rights and responsibilities. As noted earlier, 11 to 14 radio channels are available, but separate networks may coexist on the same frequency (using different network IDs (*Service Set Identifiers* [SSIDs])), albeit with performance penalties.

The workings of 802.11 devices also differ in the infrastructure and *ad hoc* modes. In the infrastructure mode (Figure 1), clients associate (and optionally authenticate) themselves with a base station, and the presence of the base station is necessary for the operation of the network.

Complex 802.11 networks may be built using the infrastructure mode, with numerous base stations providing coverage over relatively large physical areas, and clients may roam within this roaming domain, which theoretically may extend from a single building to the entire campus or town. The *Spanning-Tree Protocol* (STP) is usually used in these cases to provide loop-free bridging in this wireless LAN.

In the *ad hoc* mode (Figure 2), base stations are not used and are not necessary, because all nodes of the wireless LAN have direct reachability (that is, they “see” each other). This mode is usually used in circumstances where all devices are in close proximity to each other (such as a floor or office) and when omnidirectional antennae are used.

Figure 2: IEEE 802.11
ad hoc Network



IEEE 802.11 Roaming and Mobility

IEEE 802.11 provides for roaming and mobility of 802.11 client devices and allows clients to roam among multiple 802.11 base stations that may be operating on the same or different frequencies (channels). This is achieved through the use of *beacon frames*, which are used to synchronize 802.11 devices and, in the infrastructure mode, to associate with a base station.

There are two ways to scan for existing 802.11 networks: active and passive scanning. In active scanning mode, the 802.11 device sends out “probe” frames, soliciting “I am here” responses from existing 802.11 devices. In the passive mode, the devices just listen for beacon frames, which are periodically transmitted by the active devices. In addition, the IEEE 802.11 Task Group F is working on the IAPP, which is to provide better and interoperable mobility and roaming mechanisms.

Security of IEEE 802.11

Up to this point IEEE 802.11 could be considered an absolute success; however, security of IEEE 802.11 is not quite on par with other aspects of the standard. Although an entire chapter (Chapter 8) of the standard is dedicated to authentication and privacy, it is now the common consensus that designers of IEEE 802.11 did not excel in this area. Two reports widely covered in the media, “Your 802.11 Wireless Network Has No Clothes”^[7], and “Intercepting Mobile Communications: The Insecurity of 802.11”^[6], shed light on the apparent shortcomings of the standard, or to be more exact, on its “vulnerability by design.” They demonstrated that although the designers were well aware of the need to plan for authentication and privacy, the actual implementation was not an excellent one. The WEP algorithm, used to provide authentication and privacy in 802.11 wireless networks, is the problem.

WEP

Before discussing the security weaknesses discovered in IEEE 802.11, we quote the aim of the *Wired Equivalent Privacy* (WEP) algorithm as specified in the IEEE 802.11 standard document:

“Eavesdropping is a familiar problem to users of other types of wireless technology. IEEE 802.11 specifies a wired LAN equivalent data confidentiality algorithm. Wired equivalent privacy is defined as protecting authorized users of a wireless LAN from casual eavesdropping. This service is intended to provide functionality for the wireless LAN equivalent to that provided by the physical security attributes inherent to a wired medium.”

As you see, the aim of WEP is to provide a level of privacy equivalent to that of a wired LAN. The wording of standard is very important here: the developers of the standard did not intend to provide a level of security superior to or higher than that of a regular wired LAN, such as Ethernet. The very name of the algorithm, “Wireless Equivalent Privacy,” signifies the actual intention of the developers. However, as the practice has shown, the level of security roughly equivalent to the level of security provided by wired LANs is not sufficient—and it is the assumption that “it is OK if wireless LANs are as secure as wired LANs” that is wrong. Other problems, such as the choice of *Cyclic Redundancy Check 32* (CRC-32) instead of *Message Digest Algorithm 5* (MD5) or some other secure hash algorithm, just worsen the problem.

How WEP Works

Let's now look at the workings of WEP. WEP uses a secret key shared between 802.11 nodes to encrypt 802.11 frames (Layer 2). It also uses a checksum (CRC-32) to provide data integrity. The checksum itself is also encrypted using the shared secret key. The decryption is the reverse of the encryption process: the frame is decrypted using the key and the CRC-32 checksum is computed and checked. The cipher used in WEP is RC4, a stream cipher designed by Ron Rivest, and believed to be cryptographically strong. The key is 40 or more bits long (up to 128 bits in some implementations). However, the *Initialization Vector* that is used during the encryption process is only 24 bits long. It is difficult to understand why the designers chose such a small number—more about this later. WEP does not provide any key management—the standard itself does not specify how the shared secret key should be managed and distributed. This leaves one of the most vulnerable parts of any cryptographic system—*key distribution*—open for misuse.

The Borisov Goldberg Wagner Attacks (February 2001)

In their paper entitled “Intercepting Mobile Communications: The Insecurity of 802.11,” Nikita Borisov, Ian Goldberg, and David Wagner describe the vulnerabilities present in WEP and attacks against it. In the introduction to their paper, they state:

“Unfortunately, WEP falls short of accomplishing its security goals. Despite employing the well-known and believed-secure RC4 cipher, WEP contains several major security flaws. The flaws give rise to a number of attacks, both passive and active, that allow eavesdropping on, and tampering with, wireless transmissions.”

They go on to say that WEP fails to achieve all three of its security goals, namely confidentiality, access control, and data integrity.

As has been noted earlier, WEP uses the RC4 stream cipher with a 24-bit Initialization Vector for encryption. Borisov, Goldberg, and Wagner show that the poor design of WEP makes the system vulnerable in many areas, and one of the weakest parts of WEP is the 24-bit Initialization Vector, which may result in keystream reuse. Keystream reuse in turn permits successful cryptanalysis attacks against the ciphertext. However, what is surprising is that:

“The WEP protocol contains vulnerabilities despite the designers’ apparent knowledge of the dangers of keystream reuse attacks.”

Another not less important but equally poorly designed aspect of WEP is the use of CRC-32. It is known that CRCs are not cryptographically strong and are not intended to be used in place of message digest or hash functions such as MD5 or the *Secure Hash Algorithm* (SHA). Because of the nature of CRC, it fails to provide the required integrity protection.

Some in the industry suggest that MD5 or SHA would introduce performance penalties if used—and indeed they would—one cannot disagree. But let’s not forget that CRC-32 was intended as a security measure—which it isn’t—yes, it is fast, but it is also insecure. Presumably, a slower but really secure solution is better than an inadequate though fast solution.

The Arbaugh Shankar Wau Attack (April 2001)

In the paper “Your 802.11 Wireless Network Has No Clothes,”^[7] authors present their research of the authentication flaws in the IEEE 802.11 and demonstrate a simple eavesdropping attack against IEEE 802.11 authentication. This work is partially based on the knowledge obtained by Borisov, Goldberg, and Wagner in the paper described previously. The attack described in this work is possible even with WEP enabled; however, in that case it will also require application of attack(s) against WEP presented by Borisov et al. The authors also note that a good key management architecture would increase the security of the system; however, in their opinion only a comprehensive redesign of the standard would provide a good long-term solution to these issues.

The Fluhrer Mantin Shamir Attack (August 2001)

Scott Fluhrer, Itsik Mantin, and Adi Shamir describe a passive ciphertext-only attack against the key scheduling algorithm of RC4 as used in WEP^[11]. They identify a large number of weak keys, in which knowledge of a small number of key bits suffices to determine many state and output bits with nonnegligible probability. They also show that the first byte generated by the RC4 leaks information about individual key bytes. This paper in particular shows how to reconstruct the secret key in WEP by analyzing enough WEP-encrypted packets. The authors have not tried to do this in practice—others did that.

The Stubblefield Ioannidis Rubin Implementation of Fluhrer Mantin Shamir Attack (August 2001)

In an AT&T Laboratories report published on August 21, 2001^[14], Adam Stubblefield, John Ioannidis, and Aviel Rubin describe a real-world successful implementation of the Fluhrer Mantin Shamir attack using a \$100 Linksys card on a Linux machine. They report that it took less than a week from ordering the card to recovering the WEP key on a production network. This practical work has shown that no expensive hardware or software is necessary in order to break WEP. They summarize that it is the poor implementation of reasonable secure technologies (such as RC4) that is responsible for WEP weaknesses.

WECA's Response

The *Wireless Ethernet Compatibility Alliance* (WECA) is the organization responsible for certifying compliance with the IEEE 802.11 standards. It also awards the WiFi (*Wireless Fidelity*) industry mark to the products that have passed IEEE 802.11 compliance testing.

In response to the Berkeley paper, WECA has published an official statement, clarifying its understanding of the situation. The main line of this statement is that poor security is better than no security, as well as that WEP was not intended to be a panacea for all security needs. The statement correctly notes that the biggest security threat is the failure to use available protection methods, including WEP.

IEEE 802.11 Chair's Response

In response to the research made at UC Berkeley and the University of Maryland, the Chair of the IEEE 802.11 Working Group, Stuart Kerry, has published a Chair's response intended to clarify some of the issues around the security of IEEE 802.11. He denied allegations made in the media that the security weaknesses of WEP are due to the closed standardization process. In fact, because WEP is a part of IEEE 802.11, it was developed through an open process, like other IEEE standards. The IEEE 802.11 Working Group itself is open to all interested parties to participate. He also rejects the viewpoint that frequency-hopping wireless networks would be less vulnerable to security attacks. It is evident that this is not true because both hopping codes and timing are unencrypted and are available to the attacker. Reminding us that the goal of WEP was to provide a level of security comparable to wired LANs, he states that the IEEE 802.11 Working Group is currently working on improvements to WEP to incorporate better security into the next version of the standard.

IEEE 802.1X

Security in 802.11 networks can be broken down into three components: authentication framework, authentication algorithm/protocol, and encryption. IEEE 802.1X is trying to address the authentication framework part of the puzzle. Although still in development, 802.1X provides a scalable, centralized framework for authentication. 802.1X may deploy a variety of authentication protocols (currently Cisco's *Lightweight Extensible Authentication Protocol* [LEAP] and Microsoft's *Extensible Authentication Protocol – Transport Layer Security* [EAP-TLS] are available), and it works with both wired and wireless LANs. The widely used *Remote Access Dial-In User Service* (RADIUS) protocol is also used in the 802.1X framework. 802.1X/LEAP is available with the Cisco Aironet 350 Series of wireless LAN devices; EAP-TLS is supported in Windows XP. Although it is still a draft, 802.1X may one day become the solution to the authentication issues of 802.11.

IEEE 802.11i

Task Group I of the IEEE Working Group 802.11 is currently defining MAC enhancements to provide enhanced security for 802.11. This is a work in progress, and no IEEE 802.11i draft exists at the time of writing.

Cisco's Solution

Cisco Systems has responded to both papers on the security of the WEP^[10]. Cisco agrees that the WEP has serious shortcomings, and states that its Aironet series of wireless networking products offers many solutions to these problems: dynamic WEP keys, secure key derivation, and mutual authentication using LEAP^[13]. However, Cisco agrees that improvements are needed in the standard itself.

RC4 Fast Packet Keying for WEP

In a Document Nr 550r2, "Temporal Key Hash," submitted by Russ Housley of RSA Security and Doug Whiting of Hifn to the IEEE 802.11 Working Group, they describe a solution to the WEP problem that uses a hashing technique that rapidly generates a unique RC4 key for each packet of data sent over the wireless network. This technique addresses the performance aspect of the security solution as well—the hash algorithm used in *Fast Packet Keying* (FPK) is much faster than traditional hash algorithms such as MD5 and SHA1 because of the special caching approach. The IEEE 802.11 Working Group has decided to include this technique in the IEEE 802.11i as an informative document. In most cases, FPK may be implemented as a firmware upgrade for the existing hardware. It is possible that when released, IEEE 802.11i may use FPK as the solution—but this decision is yet to be made. No definite plans are announced at the time of writing. For more information, see:

<http://www.rsasecurity.com/rsalabs/technotes/wep-fix.html>.

Health and IEEE 802.11

Concerns about safety and health effects of various wireless solutions such as mobile phones and wireless network devices periodically surface in the media. In particular, the question of whether mobile phones are linked to brain cancer and other diseases is still open. However, in response to these concerns regarding wireless networking equipment health effects, Cisco Systems has published a white paper entitled "Cisco Systems Spread Spectrum Radios and RF Safety," which explains why these devices do not present a threat to human health when correctly used. The bottom line is that devices certified as compliant with U.S. Federal Communications Commission or Industry Canada's regulations are safe to use because of their low emitted power.

Practical Uses

Many companies, such as MobileStar, Wayport, Surf&Sip, and Airwave, have begun providing IEEE 802.11b Internet access at numerous locations throughout the United States. Several international airports also provide 802.11b service free of charge to travelers. No doubt more such services will continue to appear all over the world, maybe making a dream—Internet anywhere—a reality.

Summary

IEEE Standard 802.11 brought the long-awaited standardization to wireless LAN networking. Unfortunately, it also brought various security problems. Despite that, IEEE 802.11 is widely used, and with the coming of IEEE 802.11a, it can only gain in popularity. What now remains to be done is more effective and truly secure privacy and authentication for 802.11 wireless networks.

The IEEE 802.11 Working Group is actively working to improve what has been done to date. The most improvements are obviously needed in the area of security, where Working Groups 802.1X and 802.11i are working to define better security mechanisms. In particular, 802.11 WG is working on a new release of 802.11, which will include improvements over 802.11-1999. In the meantime, consider your wireless LAN as an external, insecure network—just like the Internet—and employ additional security measures, such as Virtual Private Networks, Transport Layer Security, SSH, and IP Security Architecture—in addition to WEP.

References

- [1] IEEE Standard 802-1990: “IEEE Standards for Local and Metropolitan Area Networks: Overview and Architecture,” ISBN 1-55937-052-1.
- [2] IEEE Standard 802.11-1999: “Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications.”
- [3] IEEE Standard 802.11a-1999: “Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications (5 GHz).”
- [4] IEEE Standard 802.11b-1999: “Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications (2.4 GHz).”
- [5] “IEEE 802.11b Wireless Equivalent Privacy (WEP) Security,” February 19, 2001, Wireless Ethernet Compatibility Alliance (WECA).
- [6] Nikita Borisov, Ian Goldberg, and David Wagner, “Intercepting Mobile Communications: The Insecurity of 802.11.”
<http://www.isaac.cs.berkeley.edu/isaac/wep-draft.pdf>
- [7] William A. Arbaugh, Narendar Shankar, and Y.C. Justin Wan, “Your 802.11 Wireless Network Has No Clothes,”
<http://www.cs.umd.edu/~waa/wireless.pdf>

- [8] William A. Arbaugh, "An Inductive Chosen Plaintext Attack Against WEP/WEP2," IEEE Document 802.11-01/230.
- [9] J. R. Walker, "Unsafe at Any Key Size; An Analysis of the WEP Encapsulation," IEEE Document 802.11-00/362.
- [10] "Cisco Comments on Recent WLAN Security Paper from University of Maryland," Cisco Systems, Product Bulletin 1327.
- [11] Fluhrer S., Mantin L., and Shamir A., "Weaknesses in the Key Scheduling Algorithm of RC4," Eighth Annual Workshop on Selected Areas in Cryptography, August 2001.
- [12] Stuart J. Kerry et al, "Response from the IEEE 802.11 Chair on WEP Security," IEEE 802.11 Working Group.
<http://www.ieee802.org/11/>
- [13] "Cisco Aironet Security Solution Provides Dynamic WEP to Address Researchers' Concerns," Cisco Systems, Product Bulletin 1281.
- [14] Adam Stubblefield, John Ioannidis, and Aviel Rubin, "Using the Fluhrer, Mantin, and Shamir Attack to Break WEP, Revision 2," AT&T Laboratories Technical Report TD-4ZCPZZ, August 21, 2001.

EDGAR DANIELYAN is a Cisco Certified Network, Design, and Security Professional, as well as member of IEEE, ACM, USENIX, SAGE, and the IEEE Computer Society. Currently self-employed, he consults and writes on internetworking, UNIX, and security. His book, *Solaris 8 Security*, was published by New Riders Publishing in October 2001. The author is not affiliated with any of the organizations (except the IEEE) mentioned in this article. E-mail: edd@danielyan.com

Code Signing

by Eric Fleischman, The Boeing Company

Code signing is a common mechanism that authors of executable code use to assert their authorship of that code and to provide integrity assurance to the users of the code that an unauthorized third party has not subsequently modified the code in any way. Code signing is widely used to protect software that is distributed over the Internet. It is also widely used for mobile code security, being a core element of the mobile code security systems of both Microsoft's ActiveX and JavaSoft's Java applet systems. Despite this widespread use, common misunderstandings have arisen concerning the actual security benefits provided by code signing. This article addresses this issue. It explains how code signing works, including its dependence upon underlying *Public Key Infrastructure* (PKI) technologies.

Motivation for Code Signing

Code signing, which is also known as *object signing* in certain programming environments, is a subset of electronic document signing. In many ways code signing is a simplification of the more generic technology in that generally only a single signature is permitted and that signature pertains to the entire file. That is, code signing usually does not support multiple signatures, encryption of (data) content, dynamic data placement, or sectional signing, which are commonly available in many document-signing systems. As a result, code signing provides only authenticity and integrity for *electronic executable files*—it does not provide privacy, authentication, or authorization, which are supported by several electronic document-signing approaches.

A signature provides authenticity by assuring users as to where the code came from—who really signed it. If the certificate originated from a trusted third-party *Certificate Authority* (CA), then the certificate embedded in the digital signature as part of the code-signing process provides the assurance that the CA has certified that the code signer is who he or she claims to be. Integrity occurs by using a signed hash function as evidence that the resulting code has not been tampered with since it was signed.

In the pre-Internet era, software was distributed in a packaged manner via branding or trusted sales outlets. It frequently came in a shrink-wrapped form directly from the vendor or a trusted distributor. In the Internet era, software is often distributed via the Web, by e-mail, or by file transfer. Code signing provides users with a similar level of assurance as to software authenticity in this comparatively anonymous—and comparatively insecure—new distribution paradigm as was previously offered by packaged software in the pre-Internet era.

In all cases, what is assured is the authorship of the software, including the verification that third parties have not subsequently modified the code. In no case does the user receive any assurance that the code itself is safe to run or actually does what it claims. Thus, the actual value of code signing remains a function of the reliability and integrity of its author. Code signing, therefore, is solely a mechanism for software creators to assert their authorship of the product and validate that it has not been modified. In no case does it provide the end user with any claim as to the quality, intent, or safety of the code.

How Code Signing Works

Code signing appends a digital signature to the executable code itself. This digital signature provides enough information to authenticate the signer as well as to ensure that the code has not been subsequently modified.

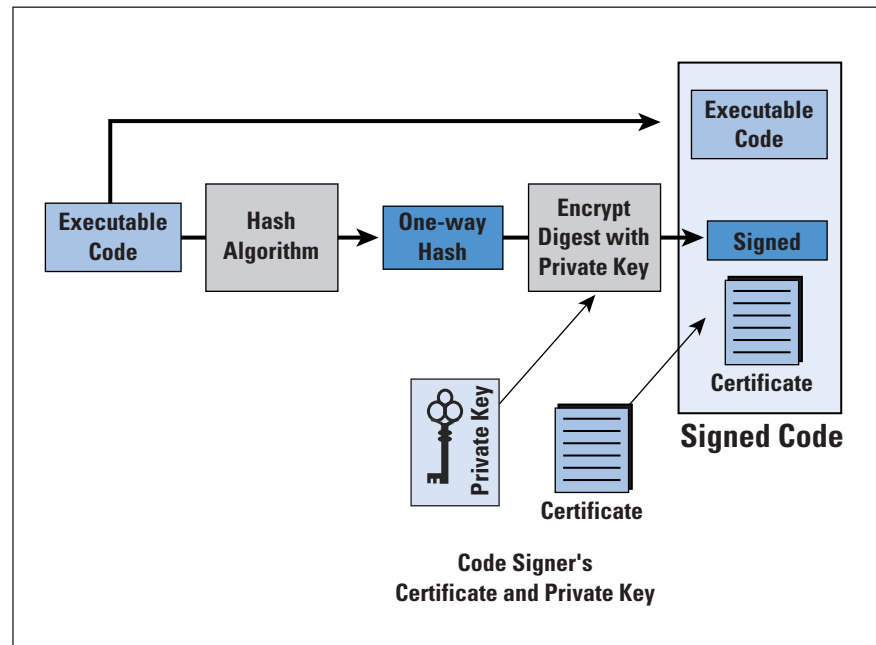
Code signing is an application within a PKI system. A PKI is a distributed infrastructure that supports the distribution and management of public keys and digital certificates. A digital certificate is a signed assertion (via a digital signature) by a trusted third party, known as the Certificate Authority (CA), which correlates a public key to some other piece of information, such as the name of the legitimate holder of the private key associated with that public key. The binding of this information then is used to establish the identity of that individual. All system participants can verify the name-key binding coupling of any presented certificate by merely applying the public key of the CA to verify the CA digital signature. This verification process occurs without involving the CA.

A *public key* refers to the fact that the cryptographic underpinnings of PKI systems rely upon asymmetric ciphers that use two related but different keys, a public key, which is generally known, and a *private key*, which should be known only by the legitimate holder of the public key. This approach is known as *public-key cryptography* and directly contrasts to symmetric ciphers, which contrastingly require the two entities to share an identical secret key in order to encrypt or decrypt information.

The certificates used to sign code can be obtained in two ways: They are either created by the code signers themselves by using one of the code-signing toolkits or obtained from a CA. The signed code itself reveals the certificate origin, clearly indicating which alternative was used. The preference of code-signing systems (and of the users of signed code) is that the certificates come from a CA, and CAs, to earn the fee they charge for issuing certificates, are expected to perform “due diligence” to establish and verify the identity of the individual or institution identified by the certificate. As such, the CA stands behind (validates) the digital certificate, certifying that it was indeed issued only to the individual (or group) identified by the certificate and that the identity of

that individual (or group) has been verified as stated. The CA then digitally signs the certificate in order to formally bind this verified identity with a given private and public key pair, which is logically contained within the certificate itself. This key pair will subsequently be used in the code-signing process. Self-created certificates, by contrast, are unconstrained as to the identities they may impersonate.

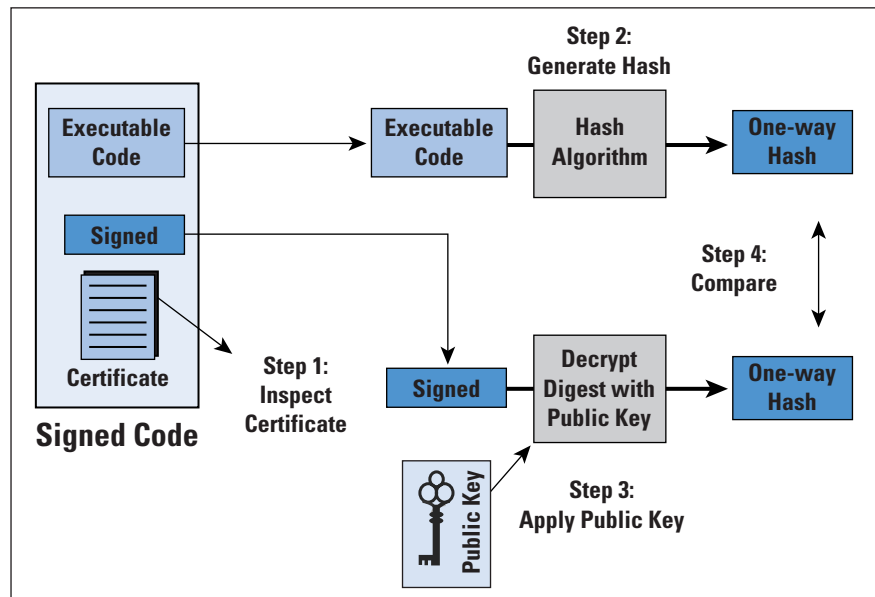
Figure 1: Code-Signing Process



Code signing itself is accomplished as follows: Developers use a *hash function* on their code to compute a *digest*, which is also known as a *one-way hash*. The hash function securely compresses code of arbitrary length into a fixed-length digest result. The most common hash function algorithms used in code signing are the *Secure Hash Algorithm* (SHA), *Message Digest Algorithm 4* (MD4), or MD5. The resulting length of the digest is a function of the hash function algorithm, but a common digest length is 128 bits. The digest is then encrypted using the developer's private key, which is part of the developer's certificate. A package containing the encrypted digest and the developer's Digital Certificate is encapsulated into a special structure called the *signature block*. The signature block is then appended to the executable code to form the signed code.

In a Java context, the signed Java byte code is called a JAR file. First introduced in the *Java Developer's Kit* (JDK) version 1.1, this capability was greatly expanded with Java 2.

Figure 2: Code Verification Process



At some subsequent time, this signed code will be presented to a recipient, usually through the agency of a code-signing verification tool on the recipient's computer. This tool will inspect the signature block to verify the authenticity and integrity of the received code. This inspection is done in the following manner, as shown in Figure 2:

1. The certificate is inspected from the signature block to verify that it is recognizable to the code-signing verification system as a correctly formatted certificate.
2. If it is, the certificate identifies the hash function algorithm that was used to create the signed digest within the received signature block. With this information, the same hash algorithm code that was used to create the original digest is then applied to the received executable code, creating a digest value, which then is temporarily stored. If it is not a correctly formatted certificate, then the code-signing verification process fails.
3. The signed digest value is then taken from the signature block and decrypted with the code signer's public key, revealing the digest value, which was originally computed by the code signer. Failure to successfully decrypt this signed digest value indicates that the code signer's private key was not used to create the received signature. If this is the case, then that signature is a fraud and the code-signing verification process fails.
4. The recomputed digest of Step 2 is then compared to the received digest that was decrypted in Step 3. If these two values are not identical, then the code has subsequently been modified in some way and the code-signing verification process fails. If any such anomaly occurs, then the verification system alerts the recipient concerning the nature of the failure, indicating that the resulting code is suspect and should not be trusted. However, if the digests are identical, then the identity of the code signer is established.

5. If establishment occurs, then the code signer's certificate is copied from the signature block and presented to the recipient. The recipient then has the option to indicate whether or not he or she trusts the code signer. If so, then the code is executed. If not, then it is not executed.

Types of Code Signing

Code signing is a mechanism to sign executable content. The term executable content refers to presenting executable programs in a manner so that they could be run locally—regardless of whether the executable file originated locally or remotely. Code signing is commonly used to identify authorship within several distinct usage scenarios:

- Applications can be code signed to identify their ownership within comparatively anonymous software distribution mechanisms using the Web, the *File Transfer Protocol* (FTP), or e-mail. This type of code signing establishes the origin for downloadable JAR, tar, zip, or CAB file software distributions, for example.
- Code signing can provide Web users more control over mobile code that is available to their Web browsers. Mobile code is code that travels a network in its lifetime in order to execute on a destination machine. The term is usually associated today with active Web content that executes on the client's machine via technologies such as Java, JavaScript, VBScript, ActiveX, and MS Word macros.
- Device drivers can be code signed to inform an operating system of the authorship of that driver. For example, the device drivers for Windows 98, Windows ME, and Windows 2000 operating systems should preferentially be certified by Microsoft's device driver certification laboratory^[25]. The entity signs the device driver executable in order to certify that the device driver in question has indeed been successfully demonstrated by a Microsoft certification laboratory to correctly run on that operating system.
- A recent news report^[20] has stated that Microsoft will be using code signing as a security mechanism within its forthcoming Windows XP operating system. The article stated: "Microsoft is to incorporate a 'signed application' system in Whistler [that is, Windows XP], the intention being to furnish users with a super-secure mode of operation that just plain stops [unsigned] code executing on the machine."

Code Signing Does Not Provide Total Security

A fundamental problem with code signing is that it cannot provide any guarantee about the good intentions of the signer or the quality, intent, operations, or safety of the code. The VeriSign and Thawte CAs, for example, combat this limitation somewhat for executables signed by certificates they issue by requiring the entities receiving their certificates to sign a "software publisher's pledge" not to sign a piece of malicious software. If they subsequently learn of violations of this agreement, they ask the owner to correct the problem.

If the owner refuses, then they cancel the owner's digital certificate and potentially bring a lawsuit against the offender. The code-signing literature has documented that the latter has occurred at least once^[21].

Another problem is that the digital signing by even a reputable entity can be forged if the private key of the signer becomes known. This forging can occur when the criminally minded exploit any of numerous potential vulnerabilities, including hacking into the key store on the signer's machine, carelessness on the part of the signer exposing this information, or an error in a CA PKI key distribution system.

Perhaps the best summary of these issues is provided by Schneier, who wrote:

"Code signing, as it is currently done, sucks. There are all sorts of problems. First, users have no idea how to decide if a particular signer is trusted or not. Second, just because a component is signed doesn't mean that it is safe. Third, just because two components are individually signed does not mean that using them together is safe; lots of accidental harmful interactions can be exploited. Fourth, "safe" is not an all-or-nothing thing; there are degrees of safety. And fifth, the fact that the evidence of attack (the signature on the code) is stored on the computer under attack is mostly useless: The attacker could delete or modify the signature during the attack, or simply reformat the drive where the signature is stored." (Quoted from page 163 of [17]).

Mobile Code Security

Mobile code security is a two-edged sword: it seeks to protect computer systems receiving potentially hostile mobile code and it also seeks to protect mobile code from potentially hostile users of those computer systems.

Code signing has emerged as a major adjunct to mobile code security. Because mobile code probably represents the dominant use of code signing that occurs today, this section examines how code signing assists mobile code security.

There is substantial and growing literature on mobile code security (for example, see [3] through [16]). The literature identifies four distinct approaches to mobile code security, together with a few hybrids that merge two or more methods. Each of the four approaches has an inherent trust model that identifies the assumptions upon which the approach is based. Rubin and Geer^[4] list these four approaches as being:

- The *sandbox approach*, which restricts mobile code to a small set of safe operations. This is the historic approach used by Java applets. In the approach, each Java interpreter implementation attempts to adhere to a security policy, which explicitly describes the restrictions that should be placed on remote applets. "Assuming that the policy

itself is not flawed or inconsistent, then any application that truly implements the policy is said to be secure. ... The biggest problem with the Java sandbox is that any error in any security component can lead to a violation of the security policy. ... Two types of applets cause most of the problems. Attack applets try to exploit software bugs in the client's virtual machine; they have been shown to successfully break the type safety of JDK 1.0 and to cause buffer overflows in HotJava. These are the most dangerous. Malicious applets are designed to monopolize resources, and cause inconvenience rather than actual loss.”^[4] The trust model assumed by the sandbox approach is that the sandbox is trustworthy in its design and implementation but that mobile code is universally untrustworthy.

- In code signing, the client manages a list of entities that it trusts. When a mobile code executable is received, the client verifies that it was signed by an entity on this list. If so, then it is run; otherwise it does not run. This approach is most commonly associated with Microsoft's ActiveX technology. “Unfortunately, there is a class of attacks that render ActiveX useless. If an intruder can change the policy on a user's machine, usually stored in a user file, the intruder can then enable the acceptance of all ActiveX content. In fact, a legitimate ActiveX program can easily open the door for future illegitimate traffic, because once such a program is run, it has complete access to all of the user's files. Such attacks have been demonstrated in practice.”^[4] The trust model for this approach assumes that it is possible to distinguish untrustworthy authors from trustworthy ones and that the code from trustworthy authors is dependable.
- The *firewalling approach* involves selectively choosing whether or not to run a program at the very point where it enters the client domain. “Research shows that it may not always be easy to block unwanted applets while allowing other applets ... to run. The firewalling approach assumes that applets can somehow be identified. ... This approach is fundamentally limited, however, by the halting problem, which states that there is no general-purpose algorithm that can determine the behavior of an arbitrary program.”^[4]

A related and more viable alternative is the playground architecture that has been used to separate Java classes that prescribe graphics actions from all other actions. The former are loaded on the client, whereas the latter are loaded on a “sacrificial” playground machine for execution and then reporting of the results to the browser. Because this approach requires byte-code modification, it cannot be used in conjunction with the usual approach to code signing.

- The *Proof-Carrying Code* (PCC) technique is a theoretical approach that statistically checks code to ensure that it does not violate safety policies. “PCC is an active area of research so its trust model may change. At present, the design and implementation of the verifier are considered trustworthy but mobile code is universally untrustworthy.”^[4]

The most common hybrid approach occurs for Java's JDK 1.1 and Java 2. Each combines the sandbox approach, which was the security mechanism for JDK 1.0, with code signing. This hybrid originated from the realization that the inherent restrictions of the sandbox model kept applications from doing "interesting and useful things." Therefore, a mechanism for running applications outside of the sandbox, code sharing, was devised to supplement the sandbox-based original. Specifically, in JDK 1.1 a signed applet enjoys unlimited access to system resources, just like local applications do, provided that the corresponding public key is trusted in the executing environment. This system evolved within Java 2 to optionally provide a consistent and flexible policy for applets and applications, determined by the policies established within a protection domain.

The literature is unanimous that the net result of this hybrid version "introduces the same security problems [as those] inherent in the ActiveX code-signing approach."^[4] For this reason, Bernard Cole^[11] has stated "neither [the sandbox nor the code signing] model is appropriate to the new environment of small information appliances, connected embedded devices, numerous web-enabled wireless phones and set-top boxes."^[11] Indeed, several articles (for example, perhaps the best collection is contained in^[13]) contained worrying descriptions of how to compromise specific sandbox and code-signing products.

The literature (see [3] through [16]) is also clear that despite the demonstrable weaknesses of both the sandbox and code-signing approaches as mechanisms for securing mobile code, they are the best practical alternatives available today. In the meantime, researchers are currently exploring enhanced mobile code security by making hybrids containing three—or all four—of the above mechanisms.

Researchers have also begun to investigate alternative techniques. For example, Zhao^[16] reports that "Additional innovative authentication functions are needed for mobile code. One approach is to apply digital fingerprinting to authenticate mobile code. Analogous to 'biometric authentication' for access control, a digital fingerprint of mobile code is a unique authentication code that is an integral and intrinsic part of the thing being authenticated. It is placed into the mobile code during its development by using digital watermarking techniques."

Major Code-Signing Systems

Code-signing systems are often functions of specific applications. For example, Thawte^[22] is a CA that provides the following certificate types:

- The *Apple Developer Certificate* is used by Apple MacOS-based application developers to sign software for electronic distribution.
- The *JavaSoft Developer Certificate* can be used with JavaSoft's JDK 1.3 and later to sign Web applets.
- A *Marimba Channel Signing Certificate* is used to sign Castanet channels on the Marimba platform.

- A *Microsoft Authenticode Certificate* is used with the Microsoft InetSDK developer tools to sign Web applets (for instance, ActiveX controls) as well as .CAB, .OCX, .CLASS, .EXE, .STL, and .DLL files, and other potentially harmful active content on Microsoft OS platforms. These Authenticode certificates work only with Microsoft IE 4.0 and later browsers.
- *VBA Developer Certificates* are identical to the Microsoft Authenticode certificates. They are used by developers to sign macros in Office 2000 and other VBA 6.0 environments.
- *Netscape Code-Signing Certificates* are used to sign Java applets, browser plug-ins, and other active content on the Netscape Communicator platform.

Despite this diversity, the clearly dominant code-signing systems today come from Microsoft, Netscape, and JavaSoft. Although these three systems generally adhere to the same set of standards, their approaches are highly diverse from each other. Each has its own certificate type. Each system approaches code signing with different orientations, goals, and expectations.

Interoperability Problems

Although all code signing uses similar technology, interoperability problems currently impact code signing. These problems may originate from interoperability problems within the underlying PKI infrastructure, from certificate differences, or from different (vendor) approaches to code signing itself.

PKI Infrastructure Interoperability

The PKI Forum has identified ten impediments to the widespread adoption of PKI^[23], the most significant being the “lack of interoperability” between PKI products. Because of this, the technical working group of the PKI Forum is currently concentrating on addressing PKI interoperability problems: “The Technical Working Group continues its focus on multi-vendor interoperability projects. Over the last six months, it has sponsored monthly interoperability “bake-offs” based on the *Certificate Management Protocol* (CMP) standard, with participation from a growing number of vendors. In addition, two workshops have been held to date on application-level interoperability through the use of digital certificates, with remote testing ongoing. Looking forward, the Technical Working group plans to initiate two new interoperability projects in the areas of Smart Card/Token Portability and CA interoperability, and it will be defining a large-scale, multi-vendor interoperability project for public demonstration in the first quarter of 2001.”^[24]

Certificate Interoperability

Numerous potential interoperability issues stem from the certificates themselves because certain certificates are themselves tied to specific types of applications.

However, not every certificate is a code-signing certificate. Rather, code-signing certificates are special certificates whose associated private keys are used to create digital signatures. In addition, the `id-kp-codesigning` value within the extended key usage field of the certificate itself (see Section 4.2.1.13 of RFC 2459) needs to be set to indicate that the certificate can be used for code signing.

In any case, code-signing certificates must be packaged in the appropriate format [*Public Key Cryptographic Standards* (PKCS)], and the various code-signing approaches (for example, Microsoft, Netscape, JavaSoft) expect both the signing certificates and the code that is to be signed to conform to different file format requirements.

These differences between code-signing systems introduce opportunities for incompatibility, even if each approach otherwise rigorously adheres to the same basic certificate standards.

Not all certificates can be used to support all potential certificate uses, even if they originate from the same CA. For example, the Java Developer Certificates are not interoperable (exchangeable) with any other certificates at this time. Fortunately, it is possible to buy certificates that can be used for many (but not all) potential uses. For example, a single certificate can support Microsoft Authenticode, Microsoft Office 2000/VBA Macro Signing, Netscape Object Signing, Apple Code Signing, and Marimba Channel Signing.

Code Signing System Interoperability

Probably the least understood of the potential interoperability problems are due to different vendor approaches to code signing itself. Perhaps McGraw and Felten have provided the best insight to code-signing system interoperability within Appendix A of their book *Securing Java*^[15]. Unfortunately, those insights were in regard to an earlier version of Java, which has evolved considerably since then.

Certificate Issues

Each of the three major code-signing systems (Microsoft, Netscape, JavaSoft) has its own certificates. Each provides its own certificate stores to house certificates within its system.

Each of the three systems supports mechanisms by which certificates may be exported from a given user's certificate store and imported into a different user's certificate store on the same or on a different machine. The Microsoft and Netscape systems also have provisions for importing certificates between code-signing systems.

Certificates are usually exported between PKI systems or certificate stores in the PKCS-12 format (`.p12` files if Netscape or `.pfx` files if Microsoft Authenticode), which contains both certificate and key pair information within the same file. Certificates can also be exported in the PKCS-7 format (for example, `.cer` or `.spc` files).

The latter approach lacks information to permit the certificate to be used for code signing by the importing system unless the missing elements can be retrieved via other mechanisms.

Code-Signing Certificates

The Netscape certificate utility (that is, *signtool -L*) indicates which of the certificates located within a certificate store can be used for code signing. By contrast, all certificates (except for those explicitly prohibited from doing code signing according to the provisions of RFC 2459 Section 4.2.1.13) within a Microsoft certificate store can be used for code signing within the Microsoft system. This means that a certificate that is unable to be used for code signing in a Netscape system can be imported into the Microsoft system and be successfully used for code signing there.

This difference stems from RFC 2459 Section 4.2.1.13, which deals with the extended key usage field. The relevant text of the standard is as follows:

“If the extension is flagged critical, then the certificate MUST be used only for one of the purposes indicated. If the extension is flagged non-critical, then it indicates the intended purpose or purposes of the key, and may be used in the correct key/certificate of an entity that has multiple keys/certificates. It is an advisory field and does not imply that usage of the key is restricted by the certification authority to the purpose indicated. Certificate using applications may nevertheless require that a particular purpose be indicated in order for the certificate to be acceptable to that application.”

What has occurred is that Netscape has implemented its system such that certificates can be used only for the purposes specified in the extended usage field. Netscape does this for both critical and noncritical markings. Microsoft, by contrast, provides that restriction solely to certificates that have been marked “critical,” permitting certificates without a critical marking to be used for any activity possible. Both approaches are legal, and both fully conform to the standard.

Code Signing from an End User’s Perspective

The results obtained when you try to execute signed code is a function of your underlying operating system, the browser you are using, and whether or not the executable is a Java applet. This should not be surprising, because similar differences also occur with unsigned code. For example, a Microsoft executable file will execute on a Microsoft Windows operating system but is unlikely to execute on operating systems that do not recognize that format. Similarly, a Java applet cannot be directly invoked on a Windows operating system, because that operating system does not recognize the **.jar** file extension. However, it will cleanly execute when accessed off of a Web page, regardless of the underlying operating system.

References

- [1] “A Closer Look at the E-signatures Law,” by Linda Rosencrance, *Computer World*, October 5, 2000.
- [2] “Standards Issue Mars E-signature,” by Jaikumar Vijayan and Kathleen Ohlson, *Computer World*, July 10, 2000.
- [3] “Mobile Code and Security,” by Gary McGraw and Edward Felten, *IEEE Internet Computing*, Volume 2, Number 6, November/December 1998.
- [4] “Mobile Code Security,” by Aviel Rubin and Daniel Geer, *IEEE Internet Computing*, Volume 2, Number 6, November/December 1998.
- [5] “Securing Systems Against External Programs,” by Brant Hashii, Manoj Lal, Raju Pandey, and Steven Samorodin, *IEEE Internet Computing*, Volume 2, Number 6, November/December 1998.
- [6] “Secure Web Scripting,” by Vinod Anupam and Alain Mayer, *IEEE Internet Computing*, Volume 2, Number 6, November/December 1998.
- [7] “Secure Java Class Loading” by Li Gong, *IEEE Internet Computing*, Volume 2, Number 6, November/December 1998.
- [8] “Mobile Code Security: Taking the Trojans out of the Trojan Horse,” by Alan Muller, University of Cape Town. April 5, 2000.
<http://www.cs.uct.ac.za/courses/CS400W/NIS/papers00/amuller/essay1.htm>
- [9] “Understanding the keys to Java Security—The Sandbox and Authentication” by Gary McGraw and Edward Felten, *JavaWorld Magazine*, May 1997.
- [10] “Repair Program or Trojan Construction Kit?” by Greg Guerin, September 7, 1999.
<http://www.amug.org/~glguerin/opinion/crypto-repair-kit.html>
- [11] “Security, Reliability Twin Concerns in Net Era,” by Bernard Cole, *Electrical Engineering Times*, July 24, 2000.
- [12] “Java Security: From HotJava to Netscape and Beyond,” by Drew Dean, Edward Felten, and Dan Wallach, Proceedings of 1996 IEEE Symposium on Security and Privacy, May 1996.
- [13] “Formal Aspects of Mobile Code Security,” by Richard Drews Dean, PhD thesis, Princeton University, January 1999.
<http://www.cs.princeton.edu/sip/pub/ddean-dissertation.php3>
- [14] “A Flexible Security Model for Using Internet Content,” by Nayeem Islam, Rangachari Anad, Trent Jaeger, and Josyula Rao, IBM Thomas J Watson Research Center, June 28, 1997.
<http://www.ibm.com/java/education/flexsecurity/>
- [15] *Securing Java—Getting Down to Business with Mobile Code*, by Gary McGraw and Edward Felten, ISBN 0-471-31952-X, John Wiley & Sons, 1999.

- [16] “Mobile Code: Emerging Cyberthreats and Protection Techniques,” by Dr. Jian Zhao, Proceedings of the Workshop on Emerging Threats Assessment—Biological Terrorism, July 7–9, 2000, Dartmouth College, Hanover, NH.
- [17] *Secrets and Lies—Digital Security in a Networked World*, by Bruce Schneier, ISBN 0-471-25311-1, John Wiley and Sons, 2000.
- [18] Telephone conversation between Bob Moskowitz and Eric Fleischman on September 26, 2000.
- [19] E-mail correspondence between Joseph M. Reagle, Jr., of the W3C and Eric Fleischman on December 6, 2000.
- [20] <http://www.theregister.co.uk/content/4/14592.html>
- [21] <http://www.halcyon.com/mclain/ActiveX/Exploder/FAQ.htm>
- [22] <https://www.thawte.com/cgi/server/step1.exe?zone=devel>
- [23] <http://pkiforum.org/About/Overview/sld037.htm>
- [24] <http://www.pkiforum.org/News/2000/PKI-Forum-third-meeting-20000919.htm>
- [25] <http://www.microsoft.com/hwtest/Signatures/>

[A longer version of this article can be obtained from the author.]

ERIC FLEISCHMAN has university degrees from Wheaton College (Illinois), the University of Texas at Arlington, and the University of California at Santa Cruz. He currently works in data communications security. He is employed as an Associate Technical Fellow by The Boeing Company. Eric was formerly employed by the Microsoft Corporation, AT&T Bell Laboratories, Digital Research, and Victor Technologies. He can be contacted at Eric.Fleischman@boeing.com

Book Review

Internet Performance Survival Guide

Internet Performance Survival Guide: QoS Strategies for Multiservice Networks, by Geoff Huston, ISBN 0-471-37808-9, John Wiley & Sons, 2000.

Many readers of IPJ are familiar with the name Geoff Huston. He contributes articles frequently. I find his style to be very lucid and his writings to be very well structured and organized.

I have need at my job to begin implementation of *Quality of Service* (QoS) strategies to deal with an ever-increasing demand for *Virtual Private Network* (VPN) tunnels over shared media. So, when I came across the title of this book and saw who wrote it, I jumped at the opportunity to review it for IPJ.

Organization

This book is organized more like a textbook than a reference manual. If you are looking for a quick and dirty guide that simply lists all the tricks of the trade and gives examples of how to implement them on specific equipment, then this book is not for you. If, however, you are looking for a well-written text that will help you to understand the issues, the practices that address them, and the theory that underlies these practices, then this is an excellent book.

The book begins with a chapter that explains in detail the problems that administrators and engineers on heterogeneous, multiprotocol networks face today. There is a quick historical survey of the evolution of networking and how that has shaped the nature of the problem. In a very topical fashion, this introduction covers the basic techniques that can be used to implement QoS, but also explains the complexity involved with these techniques, their limitations, and why they are not widely deployed yet. The book continues from there, starting with a low-level view of the building blocks of the network and gradually building to higher- and higher-level topics.

The second chapter begins with some details about the performance features built into the Internet Protocol, and in particular IPv6. This chapter continues into TCP and covers all the well-known performance features that are built into it, and then moves on to routing, switching, and *Multiprotocol Label Switching*, or MPLS. MPLS is a unified approach to switching across large networks, and it has particular applications to QoS. This topic is one of the main reasons I sought for this book, and I am glad it was covered in such detail. The second chapter ends with a survey of the various transmission systems that are available today, and discusses in detail the performance characteristics and problems that are peculiar to each.

The third chapter is a well-organized exposition of the various types of performance-tuning techniques that are available. The author keeps the discussion at a reasonably abstract level, yet is not afraid to discuss the details of the application of these techniques to the specifics of the network when such details are important. In particular, the use of QoS techniques in conjunction with the *Open Shortest Path First* (OSPF) routing protocol is discussed.

The fourth chapter combines the building blocks of Chapter 2 and the techniques of Chapter 3 into an architectural view that spans the network. The author discusses the metrics that can be used to analyze network performance, the protocols that can be used to implement service strategies, the tradeoffs that are inherent in the problem, and the policy choices that need to be made in order to come up with a clear design. In particular, the Integrated Service and Differentiated Service models are discussed separately, and then the author shows how these can be combined into an end-to-end network design. As with Chapter 3, the author explains important specific cases such as the use of the *Resource Reservation Protocol* (RSVP) with ATM.

The fifth chapter moves on to explain how the architectures that have been described can be used to attack the various kinds of problems that exist on real networks. The emphasis is clearly on the end user of the system and how to measure the levels of service being provided and to bring into play the techniques already discussed to assure a consistent level of service. The organization of this chapter seemed less clear than that of the previous chapters, but that is perhaps due more to the nature of the complexity of the problems being discussed than to the author's limitations or inattention.

The sixth chapter provides little new material, per se, and is more of a perspective on the material already provided. However, it contributes highly to the content of the book in two important ways. First, it provides more of a top-down view of QoS to complement the material in the preceding four chapters, which present a mostly bottom-up view. Secondly, it acts as a natural bookend for the first chapter. The first chapter raises the issues and poses the questions. The middle of the book examines the protocols, techniques, and architectures in detail. The last chapter then attempts to answer the questions that were initially raised.

The author does an excellent job of presenting material that is complex, vast, and is still in the process of evolving in the field. He is very diligent about managing the level of detail, and is careful to first cover the material topically before diving into the details. The examples are appropriate and have been carefully chosen.

One of the features of the material that is most appreciated is the practical perspective that the author brings to his work. The theory never gets out of hand, and is always balanced by a real-life approach to problems that, unfortunately, can never be completely solved. And, the author's observations always seem in tune with the experiences of the reader.

The material is well organized, and readers will appreciate the effort expended on the textual conventions that help to organize and structure the material. The diagrams that accompany the text are clear and well-placed, and they contribute to the reader's comprehension.

A glossary in the back helps a reader who has not thoroughly read the preceding sections of the book. The index is also well done, and the reference material is copious and pertinent.

Recommended

Overall, I would recommend this book to any professional who manages large, integrated networks, particularly those professionals who work for Internet Service Providers in an engineering capacity. I think this reflects the particular interests of the author, but that is as it should be.

—David P. Feldman, Tudor Investment Corporation

David.Feldman@Tudor.com

Would You Like to Review a Book for IPJ?

We receive numerous books on computer networking from all the major publishers. If you've got a specific book you are interested in reviewing, please contact us and we will make sure a copy is mailed to you. The book is yours to keep if you send us a review. We accept reviews of new titles, as well as some of the "networking classics." Contact us at **ipj@cisco.com** for more information.

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Fragments

ICANN Considers Structural Reform

Stuart Lynn, President and CEO of *The Internet Corporation for Assigned Names and Numbers* (ICANN) recently proposed a sweeping series of structural reforms designed to lead ICANN towards attainment of its core mission. “The current structure of ICANN was widely recognized as an experiment when created three years ago,” noted Board Chairman Vint Cerf. “The rapid expansion of and increasing global dependence on the Internet have made it clear that a new structure is essential if ICANN is to fulfill its mission.”

ICANN was formed three years ago as an entirely private global organization designed to assume responsibility for the DNS root from the U.S. government and to coordinate technical policy for the Internet’s naming and address allocation systems. In the new proposals, the basic mission remains intact, but the means of achieving that mission changes. “What has become clear to me and others is that a purely private organization will not work,” said Lynn. “The Internet has become too important to national economic and social progress. Governments, as the representatives of their populations, must participate more directly in ICANN’s debates and policymaking functions. We must find the right form of global public-private partnership—one that combines the agility and strength of a private organization with the authority of governments to represent the public interest.”

Noting that current organizational inertia and obsession with process over substance has impeded agility, Lynn laid out a roadmap designed to instill confidence in key stakeholders and to ensure that ICANN can be more effective. This roadmap entails restructuring the Board of Directors into a Board of Trustees composed in part of trustees nominated by those governments who participate in the ICANN process; in part by the chairs of proposed new “policy councils” that would replace the existing supporting organizations and that would provide expert advice; and in part by trustees proposed by a broadly-based nominating committee and appointed by the Board itself. The roadmap is designed to bring all critical stakeholders to the table, something that has been difficult to achieve with the present structure and has slowed ICANN’s progress and its ability to fulfill its responsibilities. It is also designed to establish a broad-based funding mechanism sufficient to support the critical mission of ICANN.

A paper written by Lynn that explains the reasons for change and the roadmap for reform is posted on the ICANN web site:

<http://www.icann.org/general/lynn-reform-proposal-24feb02.htm>

“We need to build a stronger organization, supported by our key stakeholders, led by the best team that can be assembled, and properly funded,” Lynn said. “We must be structured to function effectively in this fast-paced global Internet environment.” “A key requirement is to keep the best of the present ICANN,” added Cerf, “in ensuring transparency, openness, and participation, while creating an ICANN that can act responsibly and quickly. That will mean rejecting practices that have emphasized process over achievement. Above all, ICANN must be—and be seen to be—effective and supportive of technical innovation and of a reliable Internet.”

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Technology
WorldCom, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 2002 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-10/5
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSR STD
U.S. Postage
PAID
Cisco Systems, Inc.

The Internet Protocol Journal

June 2002

Volume 5, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
BEEP	2
ENUM	13
DHCP	24
Book Review	32
Call for Papers	35
Fragments	36

FROM THE EDITOR

The networking industry is full of acronyms, as the table of contents for this issue clearly illustrates. According to the dictionary, an acronym is "...a word formed from the initial letter or letters of each of the successive parts or major parts of a compound term." While neither BEEP nor ENUM are strictly speaking acronyms, these "short names" are becoming ever more prevalent and difficult to keep track of. We promise to continue to provide acronym expansion whenever possible.

BEEP is an example of a technology that came to life in a very short time. While IETF standards often take years from initial idea to protocol specification, BEEP seems to have happened in just over a year. There is already a textbook on BEEP from which our first article is adapted. Marshall Rose gives an overview of the BEEP framework and explains how you can get involved in its further development.

ENUM refers to the use of the *Domain Name System* (DNS) to look up telephone numbers and subsequently route telephone calls to the right destination using the Internet as the underlying routing fabric. This integration of the traditional telephone network with the Internet is becoming a reality and several standardization bodies are working on technologies to make this as seamless as possible. Geoff Huston explains the mechanisms and politics behind ENUM.

Our series "One Byte at a Time" examines the *Dynamic Host Configuration Protocol* (DHCP). This protocol is widely used to provide IP address and other basic routing information to clients. This is particularly useful for mobile devices, but it can be used in any network environment. Since the IP addresses are assigned as leases with a configurable time limit, DHCP also provides for effective address management. Douglas Comer explains the details of DHCP and its predecessor BOOTP.

As always, we appreciate your feedback. Send your comments and questions to ipj@cisco.com

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

An Overview of BEEP

by Marshall Rose, Dover Beach Consulting

The *Blocks Extensible Exchange Protocol* (BEEP) is something like “the missing link between the application layer and the *Transmission Control Protocol* (TCP).”

This statement is a horrific analogy because TCP is a transport *protocol* that provides reliable connections, and it makes no sense to compare a protocol to a layer. TCP is a highly-evolved protocol; many talented engineers have, over the last 20 years, built an impressive theory and practice around TCP. In fact, TCP is so good at what it does that when it came to survival of the fittest, it obliterated the competition. Even today, any serious talk about the transport protocol revolves around minor tweaks to TCP. (Or, if you prefer, the intersection between people talking about doing an “entirely new” transport protocol and people who are clueful is the empty set.)

Unfortunately, most application protocol design has not enjoyed as excellent a history as TCP. Engineers design protocols the way monkeys try to get to the moon—that is, by climbing a tree, looking around, and finding another tree to climb. Perhaps this is because there are more distractions at the application layer. For example, as far as TCP is concerned, its sole reason for being is to provide a full-duplex octet-aligned pipe in a robust and network-friendly fashion. The natural result is that while TCP’s philosophy is built around “reliability through retransmission,” there isn’t a common mantra at the application layer.

Historically, when different engineers work on application protocols, they come up with different solutions to common problems. Sometimes the solutions reflect differing perspectives on inevitable tradeoffs; sometimes the solutions reflect different skill and experience levels. Regardless, the result is that the wheel is continuously reinvented, but rarely improved.

So, what is BEEP and how does it relate to all this? BEEP integrates the best practices for common, basic mechanisms that are needed when designing an application protocol over TCP. For example, it handles things like peer-to-peer, client/server, and server/client interactions. Depending on how you count, there are about a dozen or so issues that arise time and time again, and BEEP just deals with them. This means that you get to focus on the “interesting stuff.”

BEEP has three things going for it:

- It’s been standardized by the *Internet Engineering Task Force* (IETF), the so-called “governing body” for Internet protocols.
- There are open source implementations available in different languages.
- There’s a community of developers who are clueful.

The standardization part is important, because BEEP has undergone a lot of technical review. The implementation part is important, because BEEP is probably available on a platform you're familiar with. The community part is important, because BEEP has a lot of resources available for you.

Application Protocols

An application protocol is a set of rules that says how your application talks to the network. Over the last few years, the *Hypertext Transfer Protocol* (HTTP) has been pressed into service as a general-purpose application protocol for many different kinds of applications, ranging from the *Internet Printing Protocol* (IPP)^[1] to the *Simple Object Access Protocol* (SOAP)^[2]. This is great for application designers: it saves them the trouble of having to design a new protocol and allows them to reuse a lot of ideas and code.

HTTP has become the reuse platform of choice, largely because:

- It is familiar.
- It is ubiquitous.
- It has a simple request/response model.
- It usually works through firewalls.

These are all good reasons, and—if HTTP meets your communications requirements—you should use it. The problem is that the widespread availability of HTTP has become an excuse for not bothering to understand what the requirements really are. It's easier to use HTTP, even if it's not a good fit, than to understand your requirements and design a protocol that does what you really need.

That's where BEEP comes in. It's a toolkit that you can use for building application protocols. It works well in a wide range of application domains, many of which weren't of interest when HTTP was being designed.

BEEP's goal is simple: you, the protocol designer, focus on the protocol details for your problem domain, and BEEP takes care of the other details. It turns out that the vast majority of application protocols have more similarities than differences. The similarities primarily deal with “administrative overhead”—things you need for a working system, but aren't specific to the problem at hand. BEEP mechanizes the similar parts, and lets you focus on the interesting stuff.

Application Protocol Design

Let's assume, for the moment, that you don't see a good fit between the protocol functions you need and either the e-mail or the Web infrastructures. (We'll talk more about this later on in the section “The Problem Space”.) It's time to make something new.

First, you decide that your protocol needs ordered, reliable delivery. This is a common requirement for most application protocols, including HTTP and the *Simple Mail Transfer Protocol* (SMTP).^[3] The easiest way to get this is to layer the protocol over TCP.

So, you decide to use TCP as the underlying transport for your protocol. Of course, TCP sends data as an octet stream—there aren't any delimiters that TCP uses to indicate where one of your application's messages ends and another one begins. This means you have to design a framing mechanism that your application uses with TCP. That's pretty simple to do—HTTP uses an octet count and SMTP uses a delimiter with quoting.

Since TCP is just sending bytes for you, you need to not only frame messages, but have a way of marking what's in each message. (For example, a data structure, an image, some text, and so on.) This means you have to design an encoding mechanism that your application uses with the framing mechanism. That's also pretty simple to do—HTTP and SMTP both use *Multipurpose Internet Mail Extensions* (MIME).^[4]

Back in the early 1980s, when I was a young (but exceptionally cynical) computer scientist, my advisor told me that protocols have two parts: *data* and *control*. It looks like the data part is taken care of with MIME, so it's onto the control part. If you are fortunate enough to know ahead of time every operation and option that your protocol will ever support, there's no need for any kind of capabilities negotiation. In other words, your protocol doesn't need anything that lets the participants tell each other which operations and options are supported. (Of course, if this is the case, you have total recall of future events, and really ought to be making the big money in another, more speculative, field.)

The purpose of negotiation is to find common ground between two different implementations of a protocol (or two different versions of the same implementation). There are lots of different ways of doing this and, unfortunately, most of them don't work very well. SMTP is a really long-lived, well-deployed protocol, and it seems to do a pretty good job of negotiations. The basic idea is for the server to tell the client what capabilities it supports when a connection is established, and then for the client to use a subset of that.

Well, that's just the first control issue. The next deals with when it's time for the connection to be released. Sometimes this is initiated by the protocol, and sometimes it's required by TCP because the network is unresponsive. To further complicate things, if the release is initiated by the protocol, maybe one of the computers hasn't finished working on something, so it doesn't want to release the connection just yet.

Some application protocols don't do any negotiation on connection release, and just rely on TCP to indicate that it's time to go away—even though this is inherently ambiguous. Is ambiguity a good thing in a protocol? Computers lack subtlety and nuance, so in protocols between computers, ambiguity is a bad thing. For example, in HTTP 1.0 (and earlier), you often didn't know whether a response was truncated or not. For a more concrete example, interested readers will be amused by page 2 of RFC 962.^[5]

The final control issue deals with what happens between connection establishment and release. Most application protocols tend to be client/server in nature: one computer establishes a connection, sends some requests, gets back responses, and then releases the connection. But, are the requests and responses handled one at a time (in lock-step), or can multiple requests be outstanding, either in transit or being processed, at the same time (asynchronously)?

In the original SMTP, the lock-step model was implicitly assumed by most implementors; later on, SMTP introduced a capability to allow limited pipelining. Regardless, as soon as we move away from lock-stepping, it looks as though we'll need some way of correlating requests and responses.

Although this is a step in the right direction, some application protocols need even more support for asynchrony. The reasoning is a little convoluted, but it all comes down to performance. There's a lot of overhead involved in terms of establishing a connection and getting the right user state, so it makes sense to maximize the number of transactions that get done in a single connection. While this helps in terms of overall efficiency, if the transactions are handled serially, then transactional latency—the time it takes to transit the network, process the transaction, and then transit back—isn't reduced (and may even be increased); a transaction might be blocked while waiting for another to complete. The solution is to be able to handle transactions in parallel.

Earlier I mentioned how, back in the 1980s, protocols had two parts, *data* and *control*. Today, things have changed. First of all, I'm still cynical, but more comfortable with it, and—perhaps as important—many might argue that protocols now have a third part, namely *security*.

The really unfortunate part is that security is a moving target on two fronts:

- When you deploy your protocol in different environments, you may have different security requirements.
- Even in the same environment, security requirements change over time.

This introduces something of a paradox: modern thinking is that security must be tightly integrated with your protocol, but at the same time, you have to take a modular approach to the actual technology to allow for easy upgrades. Worse, it's very easy to get security very wrong. (Just ask any major computer vendor!) Few applications folks are also expert in protocol security, and obtaining that expertise is a time-consuming, thankless task, so there's a lot of benefit in having a security mechanism menu, developed by security experts, that applications folk can pick from.

Now the good news: there's already something around designed to meet just those requirements. It's called the *Simple Authentication and Security Layer* (SASL), and a lot of existing application protocols have been retrofitted over the last four years to make use of it.

Well, let's see what all this means. Without ever having talked about what your application protocol is going to do to earn a living, we have to develop solutions for:

- Framing messages
- Encoding data
- Negotiating capabilities (versions and options)
- Negotiating connection release
- Correlating requests and responses
- Handling multiple outstanding requests (pipelining)
- Handling multiple asynchronous requests (multiplexing)
- Providing integrated and modular security
- Integrating all these things together into a single, coherent framework

So, going back to the question "Why use BEEP?", the answer is pretty simple: if you use BEEP, you simply don't have to think about any of these things. They automatically get taken care of.

Now maybe you're the kind of hardcore engineer that really wants to solve these problems yourself. Okay, go right ahead! But first, I'll let you in on a little secret: engineers have been solving these problems since 1972. In fact, they keep solving them over and over again. For each problem, there are usually two or three good solutions, and while individual tastes may vary, the sad fact is that you can make any of them work great if you're willing to put in the hours. But why put in the hours if they have nothing to do with the primary reason for writing the application protocol to begin with? Isn't there something more productive that you'd care to do with your life than design yet another framing protocol?

So, what's really *new* about BEEP? The short answer is: not much. The innovative part is that some folks sat down, did an analysis of the problems and solutions, and came up with an integrated framework that put it all together. That's not really innovation, but it's really good news if you're already familiar with the building blocks that BEEP uses.

Doesn't all this stuff add a lot of overhead? The short answer is: nope. The reason is a little more complex. BEEP is fairly minimalistic—it provides a simple mechanism for negotiating things on an à la carte basis. If you don't want privacy, no problem; don't turn it on. If you don't want parallelism, that's easy; just say “no” if the other computer asks for it. The trick here is two-fold:

- BEEP's inner mechanisms (for example, framing) are pretty lightweight, so you don't incur a lot of overhead using them (even if you don't use all the functionality they provide).
- BEEP's outer mechanisms (for example, encryption) are all controlled via bilateral negotiation, so you can decide exactly what you want to get and pay for.

There's no free lunch, but if you want to start with something “lean and mean,” BEEP doesn't slow you down, and when you want to bulk up (say, by adding privacy), BEEP lets you negotiate it. You incur only the overhead you need. (This overhead *will* show up, regardless of whether you use BEEP or grow your own mechanisms.)

It turns out that this philosophy can yield some interesting results. For example, take a look at this high-level scripting fragment:

```
::init -server example.com -port 10288 -privacy strong
```

This fragment is invoking a procedure to establish a BEEP session. With the exception of the last two terms, it looks pretty conventional.

The last two terms tell the procedure to “tune” the session by looking at the security protocols supported in common, selecting one that supports “strong privacy,” and then negotiating its use. What's interesting here is that neither the person who designed the application protocol nor the person who wrote the application making the procedure call has to be a security expert. The choice to use strong privacy, and how it gets transparently used, is all an issue of provisioning. Of course, the application protocol designer may still provide security guidelines to the implementor; naturally, the implementor may bundle a wide range of security protocols with the code. However—and this is key—everyone got to focus on what they do best (even the security guys), and it still comes together into a working system.

The cool part here is how easily this all integrates into an evolving protocol. Back in the good ol' days (say the mid-1980s) when the *Post Office Protocol* (POP)^[6] was defined, this kind of flexibility wasn't available. Whenever someone wanted to add a new security mecha-

nism for authentication or privacy, you had to muck with the entire protocol. With BEEP's framework, you just add a module that works seamlessly with the rest of the protocol. This means less work for everyone, and presumably fewer mistakes getting the work done.

Now we've come full circle: the reason for using BEEP is because it makes it a lot easier to specify, develop, maintain, and evolve new application protocols.

The Problem Space

BEEP works for a large class of application protocols. However, you should always use the right tool for the right job. Before you start using BEEP for a project, you should ask yourself whether your application protocol is a good fit for either the e-mail or Web models.

Dave Crocker, one of the Internet's progenitors, suggests that network applications can be broadly distinguished by five operational characteristics:

- Server push or client pull
- Synchronous (interactive) or asynchronous (batch)
- Time-assured or time-insensitive
- Best-effort or reliable
- Stateful or stateless

For example:

- The World Wide Web is a pull, synchronous, time-insensitive, reliable, stateless service.
- Internet mail is a push, asynchronous, time-insensitive, best-effort, stateless service.

This is a pretty useful taxonomy.

So, your first step is to see whether either of these existing infrastructures meet your requirements. It's easiest to start by asking if your application can reside on top of e-mail. Typically, the unpredictable latency of the Internet mail infrastructure raises the largest issues; however, in some cases it's a non-issue. For example, in the early 1990s, some of the earliest business-to-business exchanges were operated over e-mail (for example, USC/ISI's FAST project). If you can find a good fit between your application and Internet e-mail, use it!

More likely, though, you'll be tempted to use the Web infrastructure, and there are a lot of awfully good reasons to do so. After all, when you use HTTP:

- There's lots of tools (libraries, servers, etc.) to choose from.
- It's easy to prototype stuff.
- There's already a security model.
- You can traverse firewalls pretty easily.

All of this boils down to one simple fact: it is pretty easy to deploy things in the Web infrastructure. The real issue is whether you can make good use of this infrastructure.

HTTP was originally developed for retrieving documents in a LAN environment, so HTTP's interaction model is optimized for that application. Accordingly, in HTTP:

- Each session consists of a single request/response exchange.
- The computer that initiates the session is also the one that initiates the request.

What needs to be emphasized here is that this is a perfectly fine interaction model for HTTP's target application, as well as many other application domains.

The problem arises when the behavior of your application protocol doesn't match this interaction model. In this case, there are two choices: make use of HTTP's extensibility features, or simply make do. Obviously, each choice has some drawbacks. The problem with using HTTP's extensibility features is that it pretty much negates the ability to use the existing HTTP infrastructure; the problem with "just making do" is that you end up crippling your protocol. For example, if your application protocol needs asynchronous notifications, you're out of luck.

A second problem arises due to "the law of codepaths." The HTTP 1.1 specification, RFC 2616^[10] is fairly rigorous. Even so, few implementors take the time to think out many of the nuances of the protocol. For example, the typical HTTP transaction consists of a small request, which results in a (much) larger response. Talk to any engineer who's worked on a browser and they'll tell you this is "obvious." So, what happens when the "obvious" doesn't happen?

Some time ago, folks wanted a standardized protocol for talking to networked printers. The result was something called the *Internet Printing Protocol (IPP)*^[11]. IPP sits on top of HTTP. At this point, the old "obvious" thing (small request, big response) gets replaced with the new "obvious" thing—the request contains an arbitrarily large file to be printed, and the response contains this tiny little status indication. A surprising amount of HTTP software doesn't handle this situation particularly gracefully (that is, long requests get silently truncated). The moral is that even though HTTP's interaction model doesn't play favorites with respect to lengthy requests or responses, many HTTP implementors inadvertently make unfortunate assumptions.

A third problem deals with the unitary relationship between sessions and exchanges. If a single transaction needs to consist of more than one exchange, it has to be spread out over multiple sessions. This introduces two issues:

- In terms of stateful behavior, the server computer has to be able to keep track of session state across multiple connections, imposing a significant burden both on the correctness and implementation of the protocol (for example, to properly handle time-outs).
- In terms of performance, TCP isn't designed for dealing with back-to-back connections—there's a fair amount of overhead and latency involved in establishing a connection. This is also true for the security protocols that layer on top of TCP.

HTTP 1.1 begins to address these issues by introducing persistent connections that allow multiple exchanges to occur serially over a single connection, but still the protocol lacks a session concept. In practice, implementors try to bridge this gap by using “cookies” to manage session state, which introduces ad-hoc (in)security models that often result in security breakdowns (as a certain Web-based e-mail service provider found out).

This brings us to a more general fourth problem: although HTTP has a security model, it predates SASL. From a practical perspective, what this means is that it's very difficult to add new security protocols to HTTP. Of course, that may not be an issue for you.

If you can find a good fit between your application and the Web infrastructure, use it! (For those interested in a more architectural perspective on the reuse of the Web infrastructure for new application protocols, consider RFC 3205^[7].)

Okay, so we've talked about both the e-mail and Web infrastructures, and we've talked about what properties your application protocol needs to have in order to work well with them. So, if there isn't a good fit between either of them and your application protocol, what about BEEP?

BEEP's interaction model is pretty simple, with the following three properties:

- Each session consists of one or more request/response exchanges.
- Either computer can initiate requests or notifications.
- It's connection-oriented.

By using BEEP, you get an amortization effect with respect to the cost of connection establishment and state management. This is largely derived from the first property. Similarly, the second property gives BEEP its ability to support either peer-to-peer or client-server interactions. What we really need to explain is the connection-oriented part.

To begin, all three of the interaction models we've looked at (BEEP, e-mail, and the Web) are connection-oriented. (Although e-mail may get delivered out of order, the commands sent over each e-mail “hop” are processed in an ordered, reliable fashion.) The connection-oriented model is the most commonly used for application protocols, but it does introduce some restrictions.

A connection-oriented interaction model means that data is delivered reliably and in the same order as it was sent. If you don't require ordered, reliable delivery, you don't need a connection-oriented interaction model. For example, Internet telephony applications don't fit this model, nor do traditional multicast applications.

So, BEEP is suitable for unicast application protocols (two computers are talking to each other). However, not all unicast applications need a connection-oriented model—for example, the *Domain Name System* (DNS) manages name-to-address resolutions just fine without it. In fact, if your protocol is able to limit each session to exactly one request/response exchange with minimalist reliability requirements, and also limit the size of each message to around 65K octets, then it's probably a good candidate for using the *User Datagram Protocol* (UDP) instead.

The IETF and BEEP

BEEP is an emerging standard from the *Internet Engineering Task Force* (IETF). The IETF is a voluntary professional organization that develops many of the protocols running in the Internet. (Of course, anyone is free to develop their own protocols to run in their own little part of the Internet, but if you want multi-vendor support, you need an organization like the IETF.) So why does the IETF care about BEEP?

The answer is that the largest area in the IETF deals with application protocols. There are usually over two dozen working groups developing different application protocols. And, the IETF has been doing this for a long, long time. It turns out that even though there are well-engineered solutions to the different overhead issues, BEEP is the first time that the IETF decided to develop a standard approach that integrates the best practices for each issue. Before BEEP, each working group would spend endless hours arguing about different solutions, and then, if any time was remaining, they might sit down and look at the actual problem domain. (Okay, this is an exaggeration... but not by much!)

So, here's the process by which BEEP got designed:

- Identify the common domain-independent problems.
- Determine the best solution for each problem.
- Integrate the solutions into a consistent framework.
- Declare victory.

Now, the obvious question is: how do you determine what's "best?"

The truth is that in some cases, the answer is obvious, and in other cases, the answer is arbitrary. (Protocol experts hate to admit this, but in some cases, there is no clear winner, and it's simply better to pick *one* and order another drink.) Since most of what BEEP does is hidden from the application designer and implementor, there's really not a lot of mileage in going through it here.

beepcore.org

Where can you find out more about BEEP? To start, you can always consult the two RFCs: the BEEP core framework^[8] and the BEEP's mapping onto TCP^[9]. However, it's probably better to start with the BEEP community Web site <http://beepcore.org> where you'll find:

- News about BEEP meetings and events
- Information about BEEP projects, programmers, and consultants
- Information about beepcore (open source) and commercial software
- BEEP-related RFCs, Internet-Drafts, and whitepapers

[This article is adapted from *Beep—The Definitive Guide*, by Marshall T. Rose, ISBN 0-596-00244-0, O'Reilly & Associates, 2002. Used with permission. <http://www.oreilly.com/catalog/beep/>]

References

- [1] Herriot, R., Ed., Butler, S., Moore, P., Turner, R., "Internet Printing Protocol/1.0: Encoding and Transport," RFC 2565, April 1999.
- [2] <http://www.w3.org/TR/SOAP/>
- [3] Postel, J., "Simple Mail Transfer Protocol," RFC 821, August 1982.
- [4] Freed, N., Borenstein, N., "Multipurpose Internet Mail Extensions (MIME) Part One: Format of Internet Message Bodies," RFC 2045, November 1996.
- [5] Padlipsky, M. A., "TCP-4 prime," RFC 962, November 1985.
- [6] Rose, M. T., "Post Office Protocol: Version 3," RFC 1081, November 1988.
- [7] Moore, K., "On the use of HTTP as a Substrate," RFC 3205, February 2002.
- [8] Rose, M., "The Blocks Extensible Exchange Protocol Core," RFC 3080, March 2001.
- [9] Rose, M., "Mapping the BEEP Core onto TCP," RFC 3081, March 2001.
- [10] Fielding, R., Gettys, J., Mogul, J., Frystyk, H., Masinter, L., Leach, P., Berners-Lee, T., "Hypertext Transfer Protocol — HTTP/1.1," RFC 2616, June 1999.

MARSHALL T. ROSE is the prime mover of the BEEP Protocol. In his former position as the Internet Engineering Task Force (IETF) area director for network management, he was one of a dozen individuals who oversaw the Internet's standardization process. Rose was responsible for the design, specification, and implementation of several Internet-standard technologies, and wrote more than 60 of the Internet's Requests For Comments (RFCs). With a Ph.D. in information and computer science from the University of California, Irvine, Rose is the author of several professional texts.

E-mail: mrose@dbc.mtview.ca.us

ENUM—Mapping the E.164 Number Space into the DNS

by Geoff Huston, Telstra

Many communications networks are constructed for a single form of communication, and are ill suited to being used for any other form. Although the Internet is also a specialized network in terms of supporting digital communications, its relatively unique flexibility lies in its ability to digitally encode a very diverse set of communications formats, and then support their interaction over the Internet. In this way many communications networks can be mapped into an Internet application and in so doing become just another distributed application overlaid on the Internet. From this admittedly Internet-centric perspective, voice is just another Internet application. And for the growing population of *Voice over IP* (VoIP) users, this is indeed the case. Being able to transmit voice over the Internet is not enough. Allowing one Internet handset to connect to any other Internet handset is still not enough. In the same way that walkie-talkies became ubiquitous mobile phones only when there was a seamless integration with the telephone network, a truly useful VoIP approach will be one that supports seamless integration with the telephone network.

The basics of the telephony world are very simple indeed. Telephone handsets are little more than a speaker and a microphone. When a call is made, the network connects the microphone of one party to the speaker of the other, and vice versa. Of course you don't need a specialized telephone network to support the carriage of voice. As any user of a desktop computer would confirm, there are now a plethora of applications that can deliver a voice signal across the network. For an application to support a voice conversation, a conventional approach is to use a network base of the *User Datagram Protocol* (UDP) transport protocol, with a *Real-Time Protocol* (RTP) overlay, and the RTP payload is an encoded version of the original analogue voice signal. Carrying voice signals in real time across an Internet is a well-understood network service, with an accompanying set of existing protocols and associated applications.

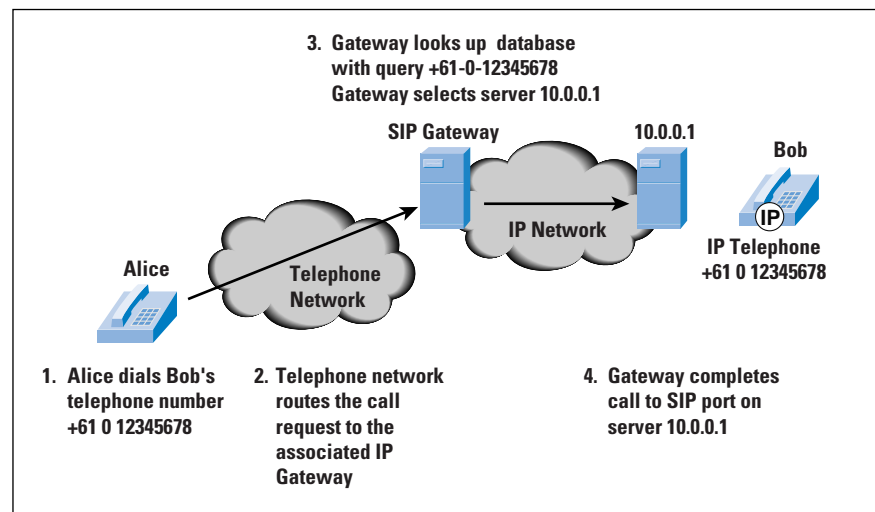
E.164 Addresses and IP Services

However, being able to transmit voice signals across a network is not enough. It was Strowger's step-by-step switching system of the late 19th century that transformed the telephone into a truly useful communications network, allowing any telephone subscriber to initiate a conversation with any other subscriber. This has evolved today into a global numbering plan where every device connected to the telephone network is assigned a unique numerical address. This numbering plan is administered by the *International Telecommunication Union* (ITU), and the plan, Recommendation E.164, involves the assignment of number prefixes to each country code administrator^[1].

If the Internet voice domain interoperates seamlessly with the telephone network, supporting this E.164 numbering plain into the realm of the Internet is a critical step. To make Internet telephony truly useful, the Internet telephony world has to be able to interface to the telephone network by allowing Internet-connected telephone devices to make and receive calls to any other telephone device, whether the other device is connected to the Internet, connected to the telephone network, or connected to any other network that seamlessly interoperates with the telephone network. For this to work, one of the preconditions is that every Internet device that supports telephone operation needs to also have an alias in the form of a unique telephone address. But there's a bit more to it than simple numbering.

Each Internet telephone is also an IP device, and, for the Internet component of the end-to-end path, the voice traffic will be carried by IP packets. These packets obviously require the IP address of the Internet telephone device. So each Internet telephone requires both an Internet address and a telephone address. It is the mapping from a telephone number to an IP address that is the crucial part of this function.

Figure 1: Calling an IP Telephone



Consider an example. When Alice, on a normal telephone, wants to call Bob, on an Internet phone, all Alice needs to do is simply dial Bob's telephone number, or his E.164 address (Figure 1). Of course, because Bob's phone is connected to the Internet and can't directly receive Alice's call request, a gateway is necessary. The telephone system should be able to map Alice's call request to the Internet telephony gateway that is configured to act as Bob's gateway agent. The gateway then needs to translate Bob's E.164 phone number into an IP address. Then the gateway has to map the telephone network signals associated with Alice's call request to corresponding signals within an Internet session initiation protocol, and then send these IP packets to Bob's Internet phone. If Bob answers the call, the phone uses the same protocol to inform the gateway, which then sends a corresponding telephone call code across the telephone network to Alice.

When Bob accepts the call, the gateway can then pass all data originating from Alice to Bob's IP address, and all data received from Bob's IP address across to the telephone connection to Alice for the duration of the call. Alice never needs to know that Bob is using an Internet device. Alice dialed a phone number, heard it ring, and then heard Bob answer the call. For Alice, nothing has changed. Bob heard the phone ring, picked it up, and talked to Alice. For Bob, nothing has changed.

The simplest way to configure each gateway is to load each gateway with a configured list of E.164 phone numbers and corresponding IP addresses. This approach is currently very common, but, like all statically configured approaches, has its weaknesses. But what happens when the IP device is numbered dynamically using the *Dynamic Host Configuration Protocol* (DHCP), or if it's mobile, and moves from one service provider's IP network to another, or when the end subscriber changes providers and that subscriber's network is renumbered, or when the primary gateway fails and the providers want to switch to a secondary device? In other words, how can this mapping be dynamic rather than static?

The way a dynamic domain name-to-IP address mapping can be maintained on the Internet is through the Internet *Domain Name System* (DNS). The telephony gateway can use the E.164 address as the DNS query, and request the DNS to return the corresponding IP address. In our example, when Alice rings Bob, the gateway can use the DNS to obtain Bob's current IP address. The gateway can then use the *Session Initiation Protocol* (SIP) to send to Bob's Internet phone a call request, which then starts Bob's phone ringing. If Bob changes IP address, then the corresponding change is a change in the DNS, not in the gateway itself. If the primary gateway fails and a secondary gateway is used, the secondary system can already access all necessary mappings through the DNS.

So the general approach of using the DNS to contain this mapping is one with some merit, but, as always, the devil is in the details. There are two parts to mapping a E.164 number into the DNS. The first is the nature of the transforms to be applied to the E.164 address to obtain a DNS query string, and the second is the form of the DNS response to this query.

Mapping E.164 Addresses into DNS Query Strings

One possible approach to mapping an E.164 number into the DNS is to simply place numbers as text blocks into the DNS. In this way, the number +61-0-12345678 could be mapped to the DNS string **61012345678.example.com**. If this method were to be used for a sizable number of E.164 numbers, there are obvious DNS performance implications associated with the size of this DNS zone file, together with the issue of frequency of update of the zone and its cache characteristics.

There are also a large number of E.164 country code delegated authorities and, consequently, a large number of entities who would like to be the authority for parts of such a monolithic unstructured DNS zone file.

In order to avoid these issues, some structure in the E.164 address space has to be used to map into the hierarchical name structure used in the DNS. One helpful observation is that E.164 numbers and Internet domain names use opposite ordering. Whereas a fully qualified domain name, such as **test.example.com**, has the more specific parts to the left and the most general part, the root, on the right of the name, a telephone number code has the most general part, the reference to the country code prefix “+” to the left and the more specific parts to the right. If one were to reverse the order of E.164 symbols, then the two address domains would have a similar structure.

One of the first efforts to provide a mapping between E.164 number and the DNS was part of the TPC fax gateway service, started in 1993^[2]. This approach uses a reversed E.164 number, and treats every digit as a node on the DNS name hierarchy. In our example, the E.164 address +61 0 12345678 would map to the DNS query string **8.7.6.5.4.3.2.1.0.1.6.tpc.int.** (in the TPC service, the parent DNS zone of this mapping is **tpc.int.**)

This mapping has some very convenient properties. Each country code corresponds to a delegatable DNS domain, so that the international country code for Australia, +61, can have a corresponding DNS delegation for the zone **1.6.tpc.int.** Within the country code the DNS can be further delegated to operators in a manner that parallels the further delegation of E.164 common prefix number blocks.

This same mapping is used by ENUM, using a DNS name parent of **e164.arpa**. The mapping entails taking a complete E.164 address (including the country code), and then removing all nondigit symbols from the address. The digit string is reversed and a “.” is placed between each pair of digits. The string **.e164.arpa.** is then appended to make a complete DNS query string. Using this process, our example number +61-0-12345678 is transformed into the DNS query:

8.7.6.5.4.3.2.1.0.1.6.e164.arpa.

Although this form of mapping is technically well suited to the DNS, it does mean that the DNS equivalent of the E.164 address is not very easily adapted to our conventional use of telephone numbers. The implication is that it is likely that Internet-based telephony applications will continue to present E.164 numbers in their user interfaces as conventional telephone numbers, and manipulate the DNS equivalent strings as internal objects.

The DNS Response

The telephone network supports more than simple voice conversations, and any serious attempt to bridge the telephone network and the Internet also should be able to handle various forms of text messaging and paging services as well as document transmission undertaken as faxes. The desired outcome is that the interface between the telephone network and the Internet should be able to seamlessly redirect the telephone service to the appropriate Internet service. In other words, we are seeing a requirement that a set of services associated with the same E.164 address should be able to be mapped to a set of IP servers, rather than a single server with a single IP address.

The implication is that the DNS response to an ENUM query should have a richer functionality than simply returning a single IP address. In DNS terms, associating a conventional “A” DNS resource record with each ENUM domain name is not sufficiently flexible for our purposes.

The approach adopted by the TPC fax gateway service was to map a fax in the telephone environment to an e-mailed multimedia message in the Internet environment. To support this mapping, telephone numbers were mapped to DNS *Mail Exchange* (MX) resource records, and these records were mapped to a mail server’s IP address in a second DNS lookup.

ENUM attempts to solve a more general model of providing mappings for any relevant service. One possible approach is to use a collection of DNS name roots, one for each mappable service. Thus, for example, **fax.e164.arpa.** could hold mappings for the fax service, while **voice.e164.arpa.** could hold mappings for voice services, and so on. However, this approach is not consistent with the generic architecture of the DNS, and the distribution of service information has the potential to lead to synchronization errors. Usefully, the DNS allows a collection of resource records to be associated with a DNS name, and this set of records is returned as the answer to a query. It is then left to the application to determine which particular record to use, with perhaps some preference hints provided in the DNS response. The approach used by ENUM takes advantage of this DNS capability, and ENUM uses the DNS to map an **e164.arpa** number onto a collection of service-specific *Uniform Resource Identifiers* (URIs)^[3].

A gateway that uses ENUM to query the DNS will receive the complete collection of service-specific URIs in response to a request to translate an E.164 address to a URI. Depending on the type of service being requested, the gateway can then select the most appropriate URI and use the DNS a second time to translate the domain name part of the URI to an IP address using the URI-specific DNS resource record as a query term. The gateway can then use the full URI specification to open an IP session with the selected service port and complete the service transaction.

The URI resource records used by ENUM are *Naming Authority Pointers* (NAPTR) records^[4]. This form of use of the DNS allows for entries where the entry itself can be decomposed into further delegations, using name formats that use URI syntax^[5].

NAPTR fields contain numerous components:

- An *Order* field to specify the order in which multiple NAPTR records must be processed
- A *Preference* field to determine the processing order when multiple NAPTR records have the same order value
- A *Service* field to specify the resolution protocol and service
- *Flags* to modify the actions of further DNS lookups
- A *regular expression* to allow the query client to rephrase the original request in a DNS format
- A *Replacement* field to define the next DNS query object

The intended operation of ENUM is to first take the E.164 number and convert it to a query in the **e164.arpa** domain. The resultant set of services is specified by the returned collection of NAPTR records. The agent selects a service that matches the service characteristics of the original request, and takes the corresponding URI for further resolution by the DNS. The elements of this URI are further decomposed as per any rewrite rules in the NAPTR record. DNS queries are generated as per the sequence of preferred NAPTR rewrite operations. The ultimate result of this sequence of DNS queries is the specification of a protocol, an associated port address, and the IP address for a preferred server for the service.

An Example of the Use of ENUM

Let's say Bob's Internet telephone services are mapped to the E.164 address +61-0-12345678. When Alice tries to call Bob, the telephone network routes the call request toward the Internet gateway that is the nominated service agent for this E.164 number. The Internet gateway takes the call setup request with Bob's number and first reverses the digits, then inserts a "." between each digit, and finally appends **e164.arpa**. The resultant DNS string is the fully qualified domain name **8.7.6.5.4.3.2.1.0.1.6.e164.arpa**. This name is then passed as a query to the DNS, to retrieve all associated NAPTR DNS resource records.

Bob has specified that he prefers to receive calls using SIP addressed to user **bob** at the server **telebob.au** by placing the following in the DNS:

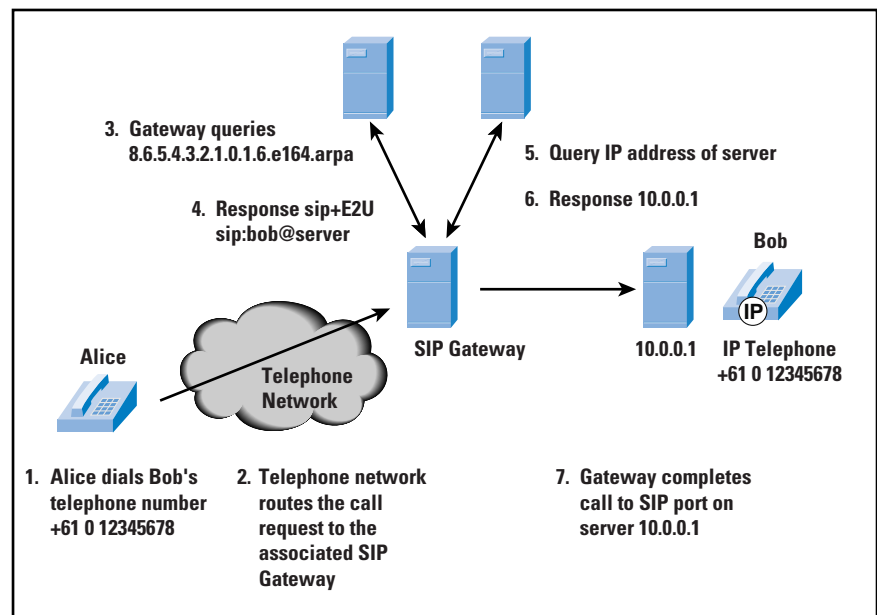
```
$ORIGIN 8.7.6.5.4.3.2.1.0.1.6.e164.arpa.
IN NAPTR 100 10 "u" "sip+E2U" "!.^.*$!sip:bob@sip.telebob.au!" .
```

In this case the DNS entry uses an order value of 100 and a preference of 10. The “u” flag indicates that the rule is terminal and that the specified URI is to be used. The service field specifies that the SIP protocol is to be used, in conjunction with the E.164 to URI (E2U) resolution service^[6]. The operation of the regular expression produces the URI of the form **sip:bob@telebob.au**.

For this call request, the gateway picks the **sip+E2U** service and performs the associated regular expression transform using the original E.164 number and the regular expression. This produces the **sip:** URI. The gateway then uses the DNS a second time to translate the domain part of the URI, **sip.telebob.au**, into an IP address using a DNS A record.

The gateway then opens up a session with UDP port 5060 on this SIP server to complete the call setup, requesting a voice session with the user Bob on this server. (Figure 2).

Figure 2: Using ENUM to Call an IP Phone



If, on the other hand, Alice is sending Bob a short text message, then Bob may want this to be delivered to him as mail. Bob would add the following entry into the DNS:

```
$ORIGIN 8.7.6.5.4.3.2.1.0.1.6.e164.arpa.
IN NAPTR 100 10 "u" "sip+E2U" "!^.*$!sip:bob@sip.telebob.au!" .
IN NAPTR 102 10 "u" "mailto+E2U" "!^.*$!mailto:bob@mail.pobob.au!" .
```

In this case the gateway would use this **mailto:** URI and use the domain part of the URI as a MX DNS query. The DNS responses are a list of mail server names and associated preferences. The gateway then selects this more preferred server and resolves this name to an IP address by a further query to the DNS for an A address record.

The gateway can complete the original text message delivery request by opening a TCP session on port 25 of the mail server and sending the message as mail addressed to user **bob@mail.pobob.au**.

Services in ENUM

Other URIs can also be associated with an E.164 number, even services not normally associated with a mapping of a telephone function. These may include **http:** URIs, even other E.164 telephone numbers, specified by **tel:** URIs.

Let's complete the example of Bob, who wants his SIP phone, mail address, Web page, and mobile telephone to be referenced from a single telephone number.

```
$ORIGIN 8.7.6.5.4.3.2.1.0.1.6.e164.arpa.
IN NAPTR 100 10 "u" "sip+E2U"      "!^.*$!sip:bob@sip.telebob.au!" .
IN NAPTR 100 10 "u" "mailto+E2U"   "!^.*$!mailto:bob@mail.pobob.au!" .
IN NAPTR 100 10 "u" "http+E2U"     "!^.*$!http://www.webhostbob.au" .
IN NAPTR 103 10 "u" "tel+E2U"      "!^.*$!tel:+61-4-12341234" .
```

Alice can enter the phone number *61012345678* into her browser and retrieve Bob's Web page in response. She can address e-mail to this number and thereby send mail to Bob. Or she can make a telephone call to Bob's SIP phone, and if it does not answer she can try Bob on his mobile phone. And she can do all this from a single number.

Numerous interesting technical issues still need to be resolved, such as the necessity and level of cacheing within the global ENUM system and the creation of a standard registry scheme for ENUM service definition.

The Politics of ENUM

There is quite some depth in the capabilities of the regular expression rewrite rules in ENUM, but the basic functionality is one of mapping a telephone number to a collection of service points that are associated with the telephone customer who was assigned that telephone number.

Despite this apparent functional simplicity, ENUM appears to have a powerful set of attractors for regulatory and social controversy.

A key benefit of moving into ENUM and the associated realm of IP-based voice communications is that service creation becomes a function of the edge and not the network. What were seen as telephone network functions such as no answer and busy redirect, call forwarding, number translation, and conference calls can all be implemented as edge applications driven by user scripts, rather than what we now see in the telephone network as value-added network-based services. One way of viewing this ENUM approach is that the DNS is functionally capable of assuming the role of service control point for telephone services, taking over the role undertaken by *Signaling System 7/Channel 7* (SS7/C7).

Service creation and signaling are slipping away from the hands of network operators into the hands of enterprises and eventually consumers, in much the same way that the Internet has redefined other services in terms of edge-based function instead of network mediation.

There is also the issue of ownership of these ENUM DNS zones, or to put it another way: who gets to populate the **e164.arpa** domain with all these URIs? It could be that this is a responsibility of existing telephone service providers, because after all these entities operate the E.164 address space in each country. It could also be that this is a responsibility of Internet Service Providers (ISPs), because the data in the resource records is describing Internet-based services. Or maybe the end subscribers get to populate the DNS with their own entries, based on a collection of services that may be sourced from a set of providers.

It is quite conceivable that we could see ISPs that have no direct role in carrying voice traffic wanting access to a country's E.164 number plan in order to provide various forms of ENUM services. Given that each element of an ENUM service collection can use URIs that refer to different ISP services, it is possible that the one ENUM record can be populated by URIs referring to numerous different service providers. This model of multi-agent access to such infrastructure resource records is a novel concept to many regulatory and operating regimes, where a single operator manages the entire associated infrastructure elements that are needed to deliver a service.

Some of the discussion about ENUM has been on more subtle aspects of this mapping. There's the choice of **e164.arpa** as the common DNS root for ENUM DNS entries. At an international level there's a lingering perception that "**arpa**" is too American and that a name root of "**int**" appears to be more neutral.

But there's something else lurking here, which has surfaced within the regulatory debate in the United States. North America has the .164 country code of "1," implying that under ENUM there is a single DNS domain for ENUM, namely **1.e164.arpa**. Single domains imply single operators, and single operators have an implication of a noncompetitive monopoly service regime. There has been a call for multiple E.164 DNS root locations for North America, allowing for two or more competing service operators using different DNS hierarchies to locate their ENUM services.

On the one side there is the view that such attempts to create multiple partially populated ENUM name hierarchies to support competitive service provision in ENUM-based services are no more than an incitement to address and service chaos. This chaos would, in turn, seriously hamper the uptake of ENUM services.

On the other hand, the competitive provision proponents of multiple DNS root domains argue that a regulatory-sanctioned monopoly is still a monopoly, and this monopoly situation will likely lead to high service prices for ENUM services. This escalated pricing structure would, in turn, seriously hamper the uptake of ENUM services.

As we have seen with the use of multiple services for an **e164.arpa** entry, the proponents of ENUM envisage a single telephone number as being an alias not only for your Internet phone service, but also for instant messaging, e-mail, your Web page, and any other service that is associated with you. One identifier is all that would be required to reach you, using a service protocol and service provider of your choice. The implication of such a use of a telephone number is, on a personal level, no more business cards cluttered with phone numbers, fax numbers, mobile numbers, e-mail addresses, Web addresses, and instant-messaging handles. Phone numbers are still the most widely used naming scheme in communications, and the use of these numbers as a universal locator has the advantage of being linguistically neutral as well as enjoying almost ubiquitous use. There are no international character set issues within this particular number space. All we need is just one ENUM address, or just one number, for all these services.

“One number to rule them all, one number to find them, one number to bring them all and in the darkness bind them,” is the ENUM version of Tolkien’s saga^[7].

But one person’s ease of use is often another’s opportunity to exploit. To be *Lord of the Numbers* would indeed be a powerful role if such uses of ENUM were to become widespread. In addition to the commercial opportunity in operating ENUM registries, ENUM can be seen as yet another erosion of personal privacy on the Internet. It can be viewed as one more step toward the use of single individual digital identity that could be used to track individuals within the Internet. On a more immediate and mundane level of concern it opens up the opportunity for spammers to use a wealth of new ways to drive you to complete distraction.

It appears that the technical components of ENUM are generally the most straightforward part. The regulatory and social implications of ENUM are more of a concern, and it is here that with ENUM we are entering into “the Land of Mordor where the shadows lie.”

Further reading:

- [1] List of ITU-T Recommendation E.164 Assigned Country Codes, available online at:
http://www.itu.int/itudoc/itu-t/ob-lists/icc/e164_717.pdf
- [2] Malamud, C., and Rose, M., “Principles of Operation for the TPC.INT Subdomain: Remote Printing—Technical Procedures,” RFC 1530, October 1993.
- [3] Fälström, P., “E.164 Number and DNS,” RFC 2916, September 2000.
- [4] Mealling, M., and Daniel, R., “The Naming Authority Pointer (NAPTR) DNS Resource Record,” RFC 2915, September 2000.
- [5] Berners-Lee, T., Fielding, R., and Masinter, L., “Uniform Resource Identifiers (URI): Generic Syntax,” RFC 2396, August 1998.
- [6] Handley, M., Schulzrinne, H., Schooler, E., and Rosenberg, J., “SIP: Session Initiation Protocol,” RFC 2543, March 1999.
- [7] Tolkien, J. R. R., *The Lord of the Rings*, George Allen and Unwin, London 1955.
- [8] <http://www.enum.org> has a good overview of ENUM and its potential application as well as references to further ENUM resources.
- [9] “Interim Approval for ENUM Provisioning,” see the Fragments section in this issue of *The Internet Protocol Journal*, page 37.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. Huston is currently the Chief Scientist in the Internet area for Telstra. He is also a member of the Internet Architecture Board, and is the Secretary of the APNIC Executive Committee. He is author of *The ISP Survival Guide*, ISBN 0-471-31499-4, *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, and coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. All three books are published by John Wiley & Sons.
E-mail: gih@telstra.net

by Douglas Comer, Purdue University

The process of starting a computer system is known as *bootstrapping*. In most systems, the initial bootstrap sequence begins with code in ROM, which the CPU executes. The ROM code only contains a first step—it merely loads an image into the computer’s RAM and branches to the image. There are two approaches used to obtain an image:

- *Embedded system*: On a diskless computer, the ROM code contains sufficient support software to permit network communication. The ROM code uses the network support to locate and download an image.
- *Conventional computer*: On a computer that has secondary storage (for instance, a PC), the ROM code loads the image from a well-known place on disk. Typically, the loaded image consists of an operating system that then controls the computer.

In either case, the image loaded by ROM is not tailored to the specific physical hardware. Instead, an image is *generic*, which means that before it can be used, it must be configured for the local hardware. In particular, the image does not contain such networking details as the computer’s IP address, address mask, or domain name. Each of these items must be supplied before applications can use TCP/IP.

Early in the history of TCP/IP, designers chose to provide a separate mechanism for each item of configuration information. Thus, the *Reverse Address Resolution Protocol* (RARP) only allowed a computer to obtain its IP address. When subnet masks were introduced, ICMP Address Mask messages were added to allow a computer to obtain a subnet mask. The chief advantage of such an approach lies in flexibility—a computer can decide which items to obtain from a local file on disk and which to obtain over the network. The chief disadvantage becomes apparent when one considers the network traffic and delay. A given computer must issue a series of small request messages. More important, each response returns a small value (for instance, a 4-octet IP address). Because networks enforce a minimum packet size, most of the space in each packet is wasted.

BOOTP

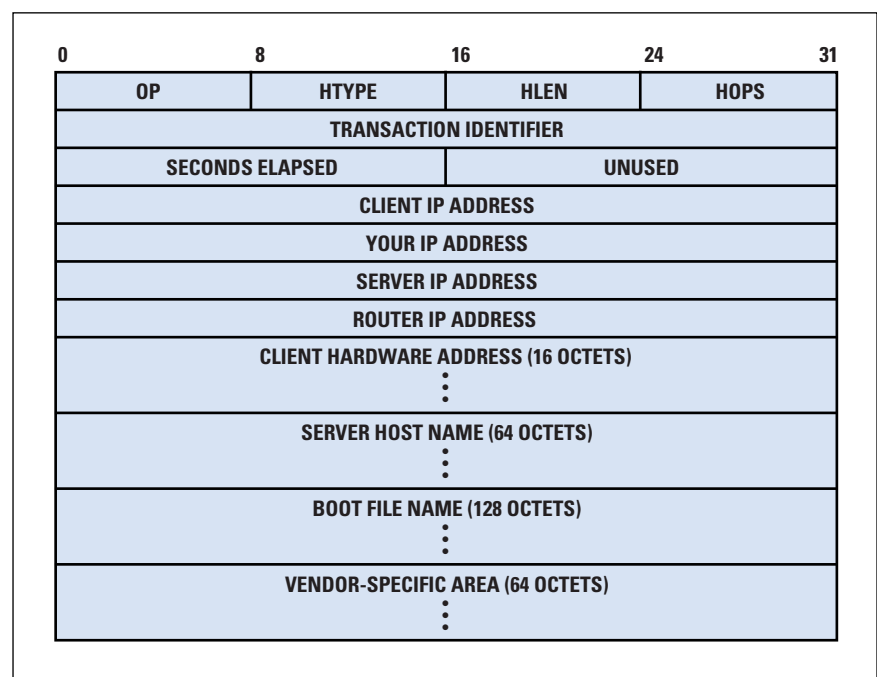
As the complexity of configuration grew, TCP/IP protocol designers observed that many of the configuration steps could be combined into a single step if a server was able to supply more than one item of configuration information. To provide such a service, the designers invented the *BOOTstrap Protocol* (BOOTP). To obtain configuration information, protocol software broadcasts a *BOOTP Request* message.

A BOOTP server that receives the request looks up several pieces of configuration information for the computer that issued the request, places the information in a single *BOOTP Response* message, and returns the reply to the requesting computer. Thus, in a single step, a computer can obtain information such as the computer's IP address, the server's name and IP address, and the IP address of a default router.

Like other protocols used to obtain configuration information, BOOTP broadcasts each request. Unlike other protocols used for configuration, BOOTP appears to use a protocol that has not been configured: BOOTP uses IP to send a request and receive a response. How can BOOTP send an IP datagram before a computer's IP address has been configured? The answer lies in a careful design that allows IP to broadcast a request and receive a response before all values have been configured. To send a BOOTP datagram, IP uses the all-1's limited broadcast address as a *DESTINATION ADDRESS*, and uses the all-0's address as a *SOURCE ADDRESS*. If a computer uses the all-0's address to send a request, a BOOTP server either uses broadcast to return the response or uses the hardware address on the incoming frame to send a response via unicast. (The server must be careful to avoid using ARP because a client that does not know its IP address cannot answer ARP requests.)

Thus, a computer that does not know its IP address can communicate with a BOOTP server. Figure 1 illustrates the BOOTP packet format. The message is sent using UDP, which is encapsulated in IP.

Figure 1: BOOTP Packet Format



Each field in a BOOTP message has a fixed size. The first seven fields contain information used to process the message. The *OP* field specifies whether the message is a *Request* or a *Response*, and the *HTYPE* and *HLEN* fields specify the network hardware type and the length of a hardware address. The *HOPS* field specifies how many servers forwarded the request, and the *TRANSACTION IDENTIFIER* field provides a value that a client can use to determine if an incoming response matches its request. The *SECONDS ELAPSED* field specifies how many seconds have elapsed since the computer began to boot. Finally, if a computer knows its IP address (for instance, the address was obtained using RARP), the computer fills in the *CLIENT IP ADDRESS* field in a request.

Later fields are used in a response message to carry information back to the computer that is booting. If a computer does not know its address, the server uses field *YOUR IP ADDRESS* to supply the value. In addition, the server uses fields *SERVER IP ADDRESS* and *SERVER HOST NAME* to give the computer information about the location of a computer that runs servers. Field *ROUTER IP ADDRESS* contains the IP address of a default router.

In addition to protocol configuration, BOOTP allows a computer to negotiate to find a boot image. To do so, the computer fills in field *BOOT FILE NAME* with a generic request (for instance, the computer can request the UNIX operating system). The BOOTP server does not send an image. Instead, the server determines which file contains the requested image, and uses field *BOOT FILE NAME* to send back the name of the file. Once a BOOTP response arrives, a computer must use a protocol like the *Trivial File Transfer Protocol* (TFTP) to obtain a copy of the image.

Automatic Address Assignment

Although it simplifies loading parameters into protocol software, BOOTP does not solve the configuration problem completely. When a BOOTP server receives a request, the server looks up the computer in its database of information. Thus, even a computer that uses BOOTP cannot boot on a new network until the administrator manually changes information in the database.

Can protocol software be devised that allows a computer to join a new network without manual intervention? Yes—several such protocols exist. For example, IPX and IPv6 can generate a protocol address from the computer's hardware address. To make automatic generation work correctly, the hardware address must be unique. Furthermore, if the hardware address and protocol address are not the same size, it must be possible to translate the hardware address into a protocol address that is also unique.

The AppleTalk protocols use a *bidding* scheme to allow a computer to join a new network. When a computer first boots, the computer chooses a random address. For example, suppose computer *C* chooses address 17. To ensure that no other computer on the network is using the address, *C* broadcasts a request message and starts a timer. If no other computer is using address 17, no reply will arrive before the timer expires; *C* can begin using address 17. If another computer is using 17, the computer replies, causing *C* to choose a different address and begin again.

Choosing an address at random works well for small networks and for computers that run client software. However, the scheme does not work well for servers. To understand why, recall that each server must be located at a well-known address. If a computer chooses an address at random when it boots, clients will not know which address to use when contacting a server on that computer. More important, because the address can change each time a computer boots, the address used to reach a server may not remain the same after a crash and reboot.

A bidding scheme also has the disadvantage that two computers can choose the same network address. In particular, assume that computer *B* sends a request for an address that another computer (for example, *A*) is already using. If *A* fails to respond to the request for any reason, both computers will attempt to use the same address, with disastrous results. In practice, such failures can occur for a variety of reasons. For example, a piece of network equipment such as a bridge can fail, a computer can be unplugged from the network when the request is sent, or a computer can be temporarily unavailable (for instance, in a hibernation mode designed to conserve power). Finally, a computer can fail to answer if the protocol software or operating system is not functioning correctly.

DHCP

To automate configuration, the *Internet Engineering Task Force* (IETF) devised the *Dynamic Host Configuration Protocol* (DHCP). Unlike BOOTP, DHCP does not require an administrator to add an entry for each computer to the database that a server uses. Instead, DHCP provides a mechanism that allows a computer to join a new network and obtain an IP address without manual intervention. The concept has been termed *plug-and-play networking*. More important, DHCP accommodates computers that run server software as well as computers that run client software:

- When a computer that runs client software is moved to a new network, the computer can use DHCP to obtain configuration information without manual intervention.
- DHCP allows nonmobile computers that run server software to be assigned a permanent address; the address will not change when the computer reboots.

To accommodate both types of computers, DHCP cannot use a bidding scheme. Instead, it uses a client-server approach. When a computer boots, the computer broadcasts a *DHCP Request* to which a server sends a *DHCP Reply*. (The reply is classified as a DHCP *offer* message that contains an address the server is offering to the client.)

An administrator can configure a DHCP server to have two types of addresses: permanent addresses that are assigned to server computers, and a pool of addresses to be allocated on demand. When a computer boots and sends a request to DHCP, the DHCP server consults its database to find configuration information.

If the database contains a specific entry for the computer, the server returns the information from the entry. If no entry exists for the computer, the server chooses the next IP address from the pool, and assigns the address to the computer.

In fact, addresses assigned on demand are not permanent. Instead, DHCP issues a *lease* on the address for a finite period of time. (When the administrator establishes a pool of addresses for DHCP to assign, the administrator must also specify the length of the lease for each address.)

When the lease expires, the computer must renegotiate with DHCP to extend the lease. Normally, DHCP will approve a lease extension. However, a site may choose an administrative policy that denies the extension. (For example, a university that has a network in a classroom might choose to deny extensions on leases at the end of a class period to allow the next class to reuse the same addresses.) If DHCP denies an extension request, the computer must stop using the address.

Optimizations in DHCP

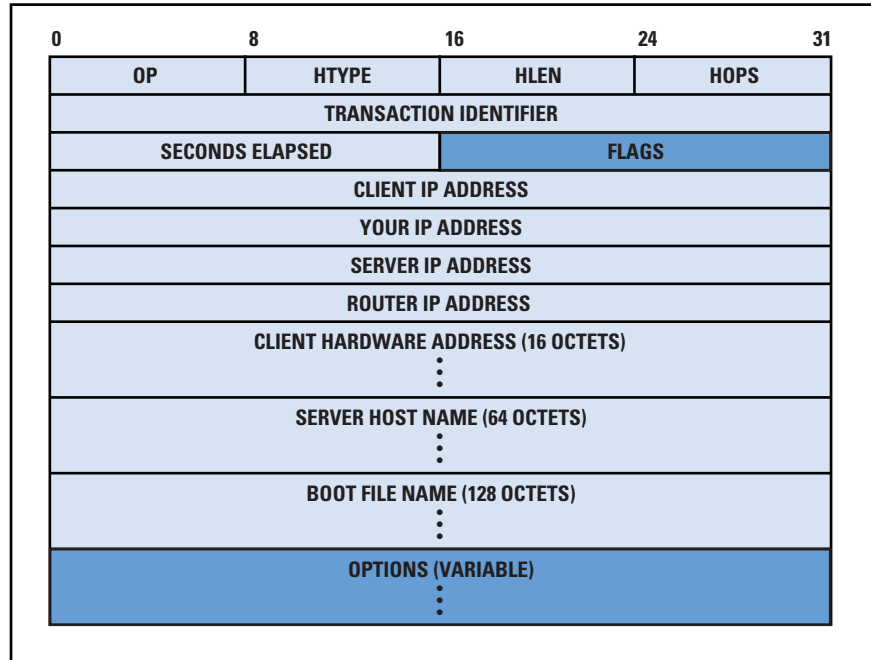
If the computers on a network use DHCP to obtain configuration information when they boot, an event that causes all computers to restart at the same time can cause the network or server to be flooded with requests. To avoid the problem, DHCP uses the same technique as BOOTP: each computer waits a random time before transmitting or retransmitting a request.

The DHCP protocol has two steps: one in which a computer broadcasts a *DHCP Discover* message to find a DHCP server, and another in which the computer selects one of the servers that responded to its message and sends a request to that server. To avoid having a computer repeat both steps each time it boots or each time it needs to extend the lease, DHCP uses *caching*. When a computer discovers a DHCP server, the computer saves the server's address in a cache on permanent storage (for example, a disk file). Similarly, once it obtains an IP address, the computer saves the IP address in a cache. When a computer reboots, it uses the cached information to revalidate its former address. Doing so saves time and reduces network traffic.

DHCP Message Format

Interestingly, DHCP is designed as an extension of BOOTP. As Figure 2 illustrates, DHCP uses a slightly modified version of the BOOTP message format.

Figure 2: DHCP Message Format



Most of the fields in a DHCP message have the same meaning as in BOOTP; DHCP replaces the 16-bit *UNUSED* field with a *FLAGS* field, and uses the *OPTIONS* field to encode additional information. For example, as in BOOTP, the *OP* field specifies either a *Request* or a *Response*. To distinguish among various messages that a client uses to discover servers or request an address, or that a server uses to acknowledge or deny a request, DHCP uses a *message type option*. That is, each message contains a code that identifies the message type.

DHCP and Domain Names

Although DHCP makes it possible for a computer to obtain an IP address without manual intervention, DHCP does not interact with the Domain Name System. As a result, a computer cannot keep its name when it changes addresses. Interestingly, the computer does not need to move to a new network to have its name change. For example, suppose a computer obtains IP address **192.5.48.195** from DHCP, and suppose the domain name system contains a record that binds the name **x.y.z.com** to the address. Now consider what happens if the owner turns off the computer and takes a two-month vacation during which the address lease expires. DHCP may assign the address to another computer. When the owner returns and turns on the computer, DHCP will deny the request to use the same address. Thus, the computer will obtain a new address. Unfortunately, the *Domain Name System* (DNS) continues to map the name **x.y.z.com** to the old address.

For several years, researchers have been considering how DHCP should interact with the DNS. Although a dynamic DNS update protocol has been defined, it has not been widely deployed. Thus, many sites that use DHCP do not have a mechanism to update a DNS database. From a user's perspective, the lack of communication between DHCP and DNS means that when a computer is assigned a new address, the computer's name changes.

Summary

The bootstrapping sequence loads a generic image into a computer, either from secondary storage or over the network. Before application software can use TCP/IP protocols, the image must be configured by supplying values for internal parameters such as the IP address and subnet mask, and for external parameters such as the address of a default router; the process is known as *configuration*. Initially, separate protocols were used to obtain each piece of configuration information. Later, the *BOOTstrap Protocol*, BOOTP, was invented to consolidate separate requests into a single protocol. A BOOTP response provides information such as the computer's IP address, the address of a default router, and the name of a file that contains a boot image.

The *Dynamic Host Configuration Protocol* (DHCP) extends BOOTP. In addition to permanent addresses assigned to computers that run a server, DHCP permits completely automated address assignment. That is, DHCP allows a computer to join a new network, obtain a valid IP address, and begin using the address without requiring an administrator to enter information about the computer in a server's database. When DHCP allocates an address automatically, the DHCP server does not assign the address forever. Instead, the server specifies a lease during which the address may be used. A computer must extend the lease, or stop using the address when the lease expires.

For Further Study

Details about BOOTP can be found in reference [1], which compares BOOTP to RARP and serves as the official protocol standard. Reference [2] tells how to interpret the vendor-specific area, and reference [3] recommends using the vendor-specific area to pass the subnet mask. Most uses of BOOTP have been replaced by DHCP. Reference [4] contains the specification for DHCP, including a detailed description of state transitions. A related document, [5], specifies the encoding of DHCP options and BOOTP vendor extensions. Finally, reference [6] discusses the interoperability of BOOTP and DHCP. The chair of the DHCP working group, Ralph Droms, and Ted Lemon have written a book about DHCP [7].

References

- [1] W. J. Croft, J. Gilmore, “Bootstrap Protocol,” RFC 951, September 1985.
- [2] J. K. Reynolds, “BOOTP Vendor Information Extensions,” RFC 1084, December 1988.
- [3] R. Braden (ed), “Requirements for Internet Hosts—Application and Support,” RFC 1123, October 1989.
- [4] R. Droms, “Dynamic Host Configuration Protocol,” RFC 2131, March 1997.
- [5] S. Alexander, R. Droms, “DHCP Options and BOOTP Vendor Extensions,” RFC 2132, March 1997.
- [6] R. Droms, “Interoperation between DHCP and BOOTP,” RFC 1534, October 1993.
- [7] R. Droms and T. Lemon, *The DHCP Handbook: Understanding, Deploying, and Managing Automated Configuration Services*, ISBN 1578701376, MacMillan, 1999.

[This article is adapted from *Computer Networks and Internets, with Internet Applications, 3rd edition*, by Douglas Comer, with CD by Ralph Droms, ISBN 0130914495, Prentice Hall, 2001.]

Dr. DOUGLAS COMER is a professor of Computer Science at Purdue University, consultant to industry, and an internationally recognized authority on TCP/IP. He has written numerous research papers and textbooks, including the classic three-volume reference series *Internetworking with TCP/IP*, and currently heads research projects. He designed and implemented X25NET and Cypress networks, and the Xinu operating system. He was a principal on the CSNET project, is director of the Internetworking Research Group at Purdue, editor of the journal *Software—Practice and Experience*, a former member of the IAB, and a Fellow of the ACM.
E-mail: comer@cs.purdue.edu

Book Review

The Elements of Networking Style

The Elements of Networking Style, by M. A. Padlipsky, originally published by Prentice-Hall, 1985, ISBN 0132681110; now available from iUniverse, 2000, ISBN 0595088791.

Sometime in the autumn of 1986, I read Padlipsky on a flight from Boston to San Francisco, and about 15 minutes into it I began to get enraged. A few minutes later, I was snickering. By the time the attendants came around with profferings of alleged comestibles, I was laughing aloud, and a gentleman sitting near the window was grateful that there was a vacant seat between us.

Padlipsky brought together several strands that managed to result in the perfect chord for me over 15 years ago. I reread this slim volume (made up of a Foreword, 11 chapters (each a separate arrow from Padlipsky's quiver) and three appendixes (made up of half a dozen darts of various lengths and a sheaf of cartoons and slogans) several months ago, and have concluded that it is as acerbic and as important now as it was 15 years ago.

The instruments Padlipsky employs are a sharp wit (and a deep admiration for François Marie Arouet), a sincere detestation for the ISO Reference Model, a deep knowledge of the *Advanced Research Projects Agency Network* (ARPANET)/Internet, and wide reading in classic science fiction.

Arouet is better known by his pen name, Voltaire. He was a social rebel, a political agitator, and an acerbic satirist comparable to Swift. Isaiah Berlin, in a lecture published in *Salmagundi* 27 [1974], remarks:

“Voltaire is the central figure of the Enlightenment, because he accepted its basic principles and used all his incomparable wit and energy and literary skill and brilliant malice to propagate the principles and spread havoc in the enemy's camp. Ridicule kills more surely than savage indignation...”

Padlipsky is pungent and sharp and witty ... and knowledgeable. His critiques of X.25, of the *International Organization for Standardization* (ISO) seven-layer cake, and of the standards process in general, are still relevant.

History

In the early 1970s, the CCITT (now the ITU), made up of PTTs and monolithic telcos, fixed upon a putative standard for a network interface protocol, X.25. First approved in 1976, and revised in 1977, 1980, 1984, 1988, and 1992, X.25 was unsatisfactory in its original form and remains less than effective.

One of the greatest drawbacks is that it is basically a store-and-forward mechanism, meaning that it has an intrinsic delay and (as noted by Sangoma Technologies) this delay is typically 0.6 seconds. It also requires a great deal of buffering space.

Padlipsky's "Critique of X.25" (Mitre Corporation Report, M82-50, September 1982; RFC 874 12 August 1983) is revised as Chapter 9 in *The Elements of Networking Style*. Padlipsky has restored, however, his original title: "Low Standards."

Flush with the failure of X.25, the *Consultative Committee for International Telegraph and Telephone* (CCITT) moved ahead.

In 1977, the British Standards Institute proposed to ISO that an architecture was needed to define the communications infrastructure. To me, this, as with *International Federation for Information Processing* (IFIP), CCITT, and similar efforts, shows how "the road to hell is paved with good intentions." Because X.25 was unsatisfactory, the IFIP Working Group was set up in the hope that that the technological community could forestall the highly political arena of ISO. (It didn't.)

ISO set up a technical committee [ISO/TC 97/SC 16]. The next year (1978), ISO published its "Provisional Model of Open Systems Architecture" [ISO/TC 97/SC 16 N 34]. This was labeled a "Reference Model," and referred to as the *Open Systems Interconnection Reference Model* (OSIRM or ISORM—pronounced "eye-sorm"—by Padlipsky).

In general, it was based on work done by Mike Canepa's group at Honeywell Information Systems, which came up with a seven-layered architecture, which itself owed a great deal to IBM's proprietary *Systems Network Architecture* (SNA). SNA had been announced in 1974, and its seven layers do not correspond exactly to OSI/ISORM's. TC 97/SC 16 turned over proposal development to the *American National Standards Institute* (ANSI), to which Canepa and his technical lead, Charlie Bachman, presented their layered model.

This, in turn, was the only proposal presented to the ISO subcommittee at a meeting in Washington in March 1978. It was accepted and published immediately. A "refined" version of the ANSI submission to ISO appeared in June 1979. This published version is nearly identical to Honeywell's of 1977.

Rage and Ridicule

While he eschews the history I've outlined here, Padlipsky is enraged by the standards process and its results. As Dave Walden and Alex McKenzie (both then at BBN, both now retired) pointed out in 1979, both virtual circuit and datagram services are valuable. "An international standard would do well to support both." [*IEEE Computer*, September 1979].

The 1977–1979 models were such that extant host-host protocols did not fit ISORM. ISO was trying to construct a set of geometric figures that would be a “tidy model.” The ARPANET workers, of whom Padlipsky was one, were interested in getting things to actually work. They were into pushing bits around the system.

The irascible Padlipsky has described the OSI system as two high rises with parking garages. The two high-rises are seven-story buildings; the parking garages are the three-story X.25 structures.

John Quarterman once pointed out:

“OSI specified before implementation. So specification took forever and implementation never happened, except for bits and pieces. In addition, heavy government backing (by the EC, now the EU, and various national governments) led some OSI participants to attempt to substitute official authority for technical capability. OSI and TCP/IP started at about the same time (1977). OSI wandered off into the weeds and TCP/IP won the race. Those governments that backed OSI bet on the wrong horse.”

TCP/IP had clearly “won the race” by the early 1980s; it took till 1994 for the U.S. government to recognize the de facto standard by rescinding its *Federal Information Processing Standards* (FIPS). At that time, too, the *Defense Data Network* (DDN) was made up of IP router nets, not X.25-based nets.

In a totally different vein, there’s Chapter 11: “An Architecture for Secure Packet-Switched Networks” (based on a presentation to the Third Berkeley Workshop on Distributed Data Management and Networking, August 1978). Here, Padlipsky suggests per-host processes. It was a really good notion.

Padlipsky’s rants—and many of the chapters are just that—precede Quarterman’s remarks by nearly a decade. But they are worth reading (and rereading).

I’m glad *The Elements of Networking Style* is available again.

—Peter H. Salus
peter@matrix.net

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

Stephen D. Crocker Receives 2002 IEEE Internet Award

The *Institute of Electrical and Electronics Engineers* (IEEE) has named Stephen D. Crocker, chief executive officer of Shinkuro, Inc. in Bethesda, Md., as recipient of the 2002 *IEEE Internet Award*. The award recognizes Crocker for his leadership in the creation of key Internet protocols. It will be presented on 19 June, at INET 2002, in Arlington, Va.

In the formative days of the Internet and its predecessor, the ARPANET, Crocker led the development of crucial technologies, processes and organizations that continue to support the Internet today. At the University of California at Los Angeles, Crocker and his team developed protocols for the ARPANET such as the *Network Control Protocol*. NCP laid the groundwork for today's *Transmission Control Protocol* (TCP). Crocker also founded and led the *Network Working Group* (NWG), which has evolved to become the *Internet Engineering Task Force* (IETF).

In organizing the notes from the first few meetings of NWG, Crocker was anxious to expand the community and invite further discussion and responses, and thus named the series *Requests for Comments*. RFCs remain a mainstay of Internet protocol publishing today, and have played a big part in creating the environment of open and evolving standards of the Internet.

"The Internet Society is honored that INET 2002 was chosen as the venue to present this year's prestigious IEEE Internet Award," said Lynn St. Amour, president and CEO of the Internet Society (ISOC). "Dr. Stephen Crocker is highly regarded throughout the international Internet community and we're pleased that his contributions will be recognized at INET 2002 in front of his peers."

Crocker's many contributions to the Internet also include extensive work organizing the standards process of the IETF, where he has served as area director of security and on the Internet Architecture Board. Crocker previously worked for the University of Southern California Information Sciences Institute in Marina del Rey, the Aerospace Corporation in El Segundo, Calif., and at Trusted Information Systems, Inc., in Glenwood, Md. In 1994, he co-founded CyberCash of Reston, Va., and served as its senior vice president for development and chief technology officer. He also has started other ventures including Steve Crocker Associates in Bethesda, Md.; Executive DSL in Bethesda, Md.; and Longitude Systems in Chantilly, Va.

He has served on the Council of Visitors at the Marine Biological Laboratory, as part of the National Research Council Study of Information Systems Trustworthiness and currently chairs the ICANN Security and Stability Advisory Committee and the ISOC 2002 Jonathan B. Postel Service Award Committee. The author of numerous papers, Crocker also holds patents in relation to his security and electronic commerce work.

He received his bachelor's degree in mathematics and doctoral degree in computer science, both from UCLA, he and studied artificial intelligence at the Massachusetts Institute of Technology.

The IEEE is the world's largest technical professional society with more than 377,000 members in approximately 150 countries. Through its members, the IEEE is a leading authority on areas ranging from aerospace, computers and telecommunications to biomedicine, electric power and consumer electronics. Additional information is available at <http://www.ieee.org>

The Internet Society <http://www.isoc.org/> is a non-profit, non-governmental, open membership organization whose worldwide individual and organization members make up a veritable "who's who" of the Internet industry. It provides leadership in technical and operational standards, policy issues, and education. ISOC is the organizational home of the International Engineering Task Force, the Internet Architecture Board, the Internet Engineering Steering Group, and the IETF—the standards setting and research arms of the Internet community. For information about INET 2002 please visit <http://www.inet2002.org>

Interim Approval for ENUM Provisioning

The *International Telecommunication Union* (ITU) and the *Internet Architecture Board* (IAB) recently announced interim approval for a single domain for ENUM, a technology that builds a bridge between the public switched telephone network and the Internet.

Voice on IP networks today operate by translating telephone numbers to IP addresses and placing an H.323 or SIP call to the device. The interchange format and translation record has not heretofore been standardized, limiting the possibility of deployment of multi-corporate and international Voice on IP services. Under the ENUM proposal, E.164 numbers can be represented as Internet Domain Names, providing a scalable and standard way to translate the numbers, and opening the way to such services. ITU has begun approving delegations for the purposes of trials. "The lack of an interoperable standard way to turn a telephone number into an IP Address has been one factor limiting the deployment of Voice on IP services internationally," said Leslie Daigle, Chair of the IAB.

If desk-mounted computers or servers are given telephone numbers as well as mnemonic names, this system further enables common telephone handsets to place Voice or Video on IP calls to such computers. This is a significant step towards integrating Internet-based services with the global telephone network, and the current agreements between IAB and ITU will allow trials to take place.

Patrik Fältström, member of the *Internet Engineering Steering Group* (IESG), said that “the integration of the desktop telephone and computer allows corporations to simplify their internal networks.”

Roy Blane, Chair of ITU-T’s Study Group 2, concurred, saying that “In the long term this protocol may facilitate many new internet services. In the short term, countries wishing to trial the system can begin work on developing it.”

This interim approval is made possible due to cooperation between ITU, IAB and the IETF. As outlined in the ENUM specification document, RFC 2916, sub-domains from a single domain will be delegated after acceptance by the registries according to the existing assignment of country codes in the telephone address space. Information on how the ENUM registration requests will be processed can be found at:

<http://www.ripe.net/enum/>

The IETF is an international community of network designers, operators, vendors, and researchers concerned with the evolution of the Internet architecture and the smooth operation of the Internet. The definition of the ENUM protocol, as proposed by the IETF can be found at <http://www.ietf.org/rfc/rfc2916.txt> The IETF is an organized activity of the Internet Society.

The ITU is a global organization where the public and private sectors cooperate for the development of telecommunications and the harmonization of national telecommunications policies. Study Group 2 of the *ITU Telecommunication Standardization Sector* (ITU-T), where work on ENUM is being carried out, is the Lead Study Group on Service definition, Numbering, Routing and Global Mobility and is responsible for the operational aspects of service provision, networks and performance. More information on the ENUM protocol, and the issues related to it, can be found at <http://www.itu.int/ITU-T/worksem/enum/index.html>

Committee on ICANN Evolution and Reform posts Recommendations

Following the publication in February of “President’s Report: ICANN—The Case for Reform,” by Stuart Lynn, President and CEO of *The Internet Corporation for Assigned Names and Numbers* (ICANN), a committee of the board has been examining the details of the restructuring proposal, receiving input from the community at large, and publishing several documents with recommendations. You can find pointers to all of these documents in the “Announcements” section at <http://www.icann.org>

Upcoming Events

INET 2002, the annual conference of the Internet Society, will be held June 18–21, 2002 at the Crystal Gateway Marriott, in Arlington, Virginia (5 minutes from downtown Washington, DC).

<http://www.inet2002.org/>

The *IETF* will be meeting in Yokohama, Japan, July 15–19, 2002 and in Atlanta, Georgia, USA, November 17–22, 2002.

<http://www.ietf.org/meetings/meetings.html>

ACM SIGCOMM 2002 is the annual conference of the *Special Interest Group on Data Communication* (SIGCOMM), a vital special interest group of the *Association for Computing Machinery* (ACM). This year, SIGCOMM will be held in Pittsburgh, Pennsylvania, August 19–23.

<http://www.acm.org/sigcomm/sigcomm2002/>

ICANN will meet in Bucharest, Rumania, June 24–28, 2002 and in Shanghai, China, October 27–31, 2002.

<http://www.icann.org/meetings/>

The *Asia Pacific Network Information Centre* (APNIC) will hold its next Open Policy Meeting, September 3–6, 2002 in Kitakyushu, Japan. <http://www.apnic.net/meetings/index.html>

The next *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) will take place February 19–28 in Taipei, Taiwan. <http://www.apricot2003.net/>

Errata List

This is the 17th issue of *The Internet Protocol Journal*. Inevitably, some minor, and a few major errors have made their way into print since our June 1998 issue. We are planning to publish a list of corrections on our Web site in the near future. Since the online material is a reflection of the printed version, we feel it would be inappropriate to simply “silently” correct the online editions, thereby rewriting history. Instead, a list of the errors along with the corrections will be presented.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Technology
WorldCom, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 2002 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSR STD U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol Journal

September 2002

Volume 5, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Visitor Networks	2
Wireless Security	17
The Uncommon Carrier	23
Letters to the Editor	28
Book Review	31
Fragments	33

FROM THE EDITOR

The *Internet Protocol Journal* (IPJ) does not have a marketing department. New subscribers learn about IPJ through our Web page, or perhaps by picking up a copy at an Internet conference or meeting such as the IETF. Word of mouth is perhaps the most effective “marketing tool.” I was reminded of this in July when an article in IPJ was mentioned on the *SlashDot* Web site. Within a few days we received more than 900 new subscriptions, on the order of ten times the normal sign-up rate. I think this illustrates the power of the Web as a tool for information dissemination.

I am a big fan of visitor networks. Such networks, typically found in larger hotels, allow high-speed access to the Internet for a daily or weekly fee. Although most of the conferences and meetings I attend have purpose-built “terminal rooms,” it is still nice to be able to work in your hotel room at speeds orders of magnitude better than what can be obtained with a dialup modem. Dory Leifer explains how visitor networks are designed and operated in our first article.

In a previous article we explored the basics of IEEE 802.11 wireless networking. Such networks are growing at an amazing rate. Reports about wireless network “wiretapping” are frequently found in the trade press. Gregory R. Scholz describes an architecture for securing wireless networks, using a variety of technologies and protocols.

Geoff Huston is back with another opinion piece, this time discussing the role of the *Internet Service Provider* (ISP) as a “common carrier.” Many ISPs are finding themselves in the middle of disputes between customers, copyright owners, regulators and others. What role should an ISP play in this regard? Geoff provides some answers.

Please continue to provide your feedback to anything you read in this journal. Our “Letters to the Editor” section provides a sample of some of the correspondence we receive. As always, use ipj@cisco.com to contact us.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Visitor Networks

by Dory Leifer, DEL Communications Consulting

Visitor networks are LANs that are most often deployed in hotels, airports, cafés, college campuses, apartments, and other locations. They enable the public network access on an ad-hoc basis. Recently, 802.11 “hot spots” have gained increased attention; they represent one example of a visitor network.

Visitors attach devices such as a laptop or *personal digital assistant* (PDA) that they use only while traveling or, more often, they attach machines normally used in the office or home. These machines can be thought of as “visiting hosts.”

This article explores some of the technical issues with IP visitor networks and considers practical options for service provider deployment on wired Ethernet and wireless networks. In exploring deployment options, the article focuses mainly on solutions that do not require client software on the visiting host. These clientless techniques are based on heuristics and, although they do not work effectively under all circumstances, they have proven to be quite useful in practice.

For this discussion, it is assumed that the service provided by the visitor network is for access in one location at a time. Therefore, the article does not address network hand-off for mobile clients that are moving from one network attachment point to another while attempting to maintain connectivity.

Traditional LANs vs. Visitor Networks

Traditional LANs have been well optimized for enterprise networks. They provide high bandwidth and an economical and universal method of delivering network connectivity. In comparison, visitor networks are a rather curious hybrid of a LAN and a public network, such as one used for dial-in network access. Their objective is to physically use LANs to deliver what has normally been considered a public network service: *universal access*.

In enterprise networks, traditional LANs are usually carefully administrated. Normally the connected hosts are owned and administrated by the same enterprise that operates the network. Hosts that are connected to the network are configured according to the designated protocol and address schemes. They are often configured for at least *Simple Mail Transfer Protocol* (SMTP), *Post Office Protocol* (POP), file, and print sharing. On visitor networks, the hosts are typically owned and configured by the visitors, while the service provider administrates the network.

This difference in administration creates a serious challenge for the visitor network. The network must support a wide range of configurations because they will differ from one visiting host to another. For example, if a host had previously been configured for a static IP address, that address is likely to be from a different subnet, perhaps from a private network that the visitor normally uses at the office. Even if a host gets some of its configuration from *Dynamic Host Configuration Protocol* (DHCP), *Domain Name System* (DNS) and SMTP servers may refer to addresses or names on a private network that are not reachable on the visitor network.

Traditional wired LANs normally span physically secure areas, so any person who has access to the Ethernet wall jack for the building can connect anything to the network. With a visitor network it may be undesirable to allow everyone access. For example, a visitor network deployed in a university library may be available only to students. Similar to public dial-in access, visitor networks often rely on authentication and authorization before granting service.

Whereas LANs are excellent at facilitating peer-to-peer services such as file and print sharing between connected hosts, visitor networks often attempt to minimize these direct interactions between visitors, instead establishing a set of services that the service provider itself offers or simply routing the IP packets off the LAN to an Internet Service Provider. Minimizing interactions between visitors is desirable because service providers will want to reduce the risk of a visitor's machine being attacked by another visitor. On some occasions, however, visitors who do trust each other may want to use the visitor network for file sharing, printing, or even network gaming.

Going Clientless

One of the most difficult choices for service providers deploying visitor networks is to decide whether or not to rely on the installation of specialized client software on the visiting host.

Client software allows specific network protocols to be passed between the client and the visitor network. Protocols such as *Point-to-Point Protocol over Ethernet* (PPPoE)^[1], *Layer 2 Tunneling Protocol* (L2TP)^[2], and *Mobile-IP*^[3] support both authentication as well as IP tunneling to assist in routing and address assignment. On some wireless LANs and networks with high-end Ethernet switches, 802.1x (which will be discussed in more detail later) supports flexible authentication schemes and aids in data encryption^[4]. Although these protocols implemented on the client can present a significant technical advantage for implementing visitor networks, they require at least some modification to the configuration on the visiting host.

The lowest common denominator for traveling laptops is a simple TCP/IP stack and a browser. If the service can accommodate the visitor with only these items, the visitor network becomes much more suitable to the broadest audience. Of course without authentication, tunneling, and client configuration available from client software, the visitor network must rely on a set of heuristics or, said by some, hacks, to perform its tricks. Subsequent sections of this article illustrate technically how a visitor network can operate without relying on the installation of client software.

The service provider may choose to distribute client software in a situation where the visitor may use the service repeatedly. In many other situations, however, it is not feasible. For example, the last thing that travelers want to find in a hotel room upon arriving at midnight and needing a network connection is a CD-ROM full of new software drivers to drop on their laptop before using the hotel's in-room Ethernet. Even if the provided software does nothing but change the configurations, such as select a Web proxy server, it may have negative consequences when the laptop is returned to the office. Such added steps could also discourage visitors from using the visitor network again.

Visitor Network Basics

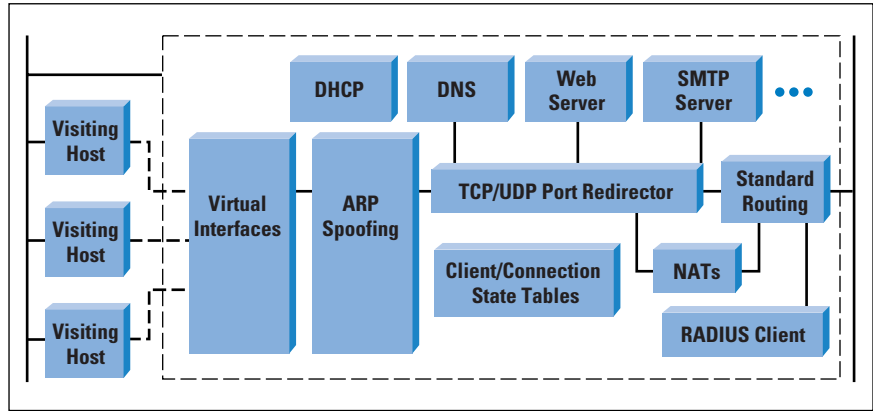
There are no hard guidelines or standards on what constitutes a visitor network. However, numerous vendors are selling devices that operate with wired and wireless networks, and act as gateways between the visitor network and the traditionally routed infrastructure. The typical visitor experience proceeds as follows (this is essentially a clientless example): the visiting host would not require the installation of special software, and in many cases would not require configuration changes:

- The visiting host is physically attached to the network by connecting to a twisted-pair Ethernet port.
- Visitors open their browser and attempt to load any page with the *Hypertext Transfer Protocol* (HTTP).
- Regardless of the specified URL, the browser loads a default page that requests authentication or billing information.
- When authenticated, the visitors now have general Internet access.
- An accounting record describing a visitor's session is generated and processed by the service provider's billing system, resulting in a charge on either the visitor's account or a corporate account.

Visitor Gateways

Visitor networks can be implemented with a special-purpose device called a "visitor gateway." Figure 1 illustrates the basic functional schematic of an example device. (Unfortunately, just about every vendor selling these devices uses a different name. This article uses the term in a generic sense and not to refer to any company's particular product.)

Figure 1: Visitor Gateway



The visitor gateway sits between the LANs used to provide service to the visitors and a standard routed interface. Physically, a visitor gateway is a device that appears much like a router or firewall, with minimally two Ethernet interfaces.

Hybrid of NAS and LAN

The following sections focus on the visitor gateway, specifically its operational model, its handling of various Internet packet types, *virtual LANs* (VLANs), authentication, and accounting.

Visitor gateways behave as a hybrid of a standard LAN and a *Network Access Server* (NAS). For illustration, one can compare the operation of the visitor gateway with the operation of a NAS. Like a NAS with individual modem ports, the visitor network gateway typically builds virtual port structures as new hosts are discovered on the connected LAN. These virtual interfaces are configured by the gateway to accommodate the IP addresses used and referred to by the visiting host. The visitor gateway may create a virtual port structure for every host based on its *Media Access Control* (MAC) address or VLAN identifier and treat every virtual interface as an independent subnet upon which the visiting host and the virtual interface of the visitor network are the only attachments. Think of the relationship as a logical point-to-point link.

Conversely, the NAS, using the *Point-to-Point Protocol* (PPP)^[5] on a dial-in connection, has a significant advantage over the visitor gateway in this scenario. PPP allows the NAS to negotiate an acceptable IP address for the dial-in client, set the client's default gateway, and even in some cases configure the client's DNS. The NAS normally has at least *Password Authentication Protocol* (PAP) and *Challenge Handshake Authentication Protocol* (CHAP) for authentication. If the visiting host requests configuration through DHCP^[6], the visitor network has an opportunity to assign private or public addresses that are mutually convenient for both parties. On the other hand, if the visiting host already has a static address configured for its native network, for example, then the visitor gateway must spoof or imitate the behavior of the configured subnet.

The appeal of PPP in the dial-in world led to the recent development of PPPoE for LANs. Although PPPoE has been used with service selection gateways to offer public *Digital Subscriber Line* (DSL), there has been little use of it on visitor gateways. This is likely to be true because of the lack of a ubiquitous client and the complexities of solving multilevel authentication and encryption involving the local link, local network, and private network. PPPoE certainly is worth future study for visitor networks.

ARP

Hosts learn Layer 2 MAC addresses using the *Address Resolution Protocol* (ARP). Although hosts and routers respond only when asked about the IP address of their interfaces or those on a proxy-ARP table, visitor gateways usually respond with their own MAC address to any ARP requests from the attached visiting hosts, effectively proxying for the host's default gateway (if one is configured). The visitor gateway can also configure the interface address of its virtual port based on the host's IP address. In this manner, the gateway auto-configures itself to accommodate the visitor, who can continue to use his/her configured address.

Used on a standard shared LAN, this technique only goes so far. If, for example, one host on the visitor network shared its default router configuration with the IP addresses of another host (not that uncommon for private network numbers), then when the first host attempted to get the MAC address of its default router, it would end up with two responses, one from the visitor gateway and one from the other host on the LAN.

TCP/UDP Port Redirector

The visitor gateway for each *Transmission Control Protocol* (TCP) and *User Datagram Protocol* (UDP) packet received from the visiting host decides whether to pass the packet through or direct it to a local service such as DNS, SMTP, or Web server. It makes this decision based on some configured policy from the service provider (such as to redirect all SMTP) and from authorization states of the visitors. For example, if the service provider wishes to charge visitors \$10 for daily access at a hotel, the port redirector could reflect HTTP requests to the local Web server that would, in turn, present the option to the visitor. Subsequent HTTP requests presumably would always be passed transparently through the gateway to the intended address.

The operation of the redirector is fairly simple. It works as a backwards network address-port translator. Instead of modifying the source, it modifies the destination and then applies standard IP forwarding on the resulting packets.

DNS

Visitor gateways typically implement proxies for domain name service requests and channel all DNS requests from the visiting host through the proxy. This serves at a minimum to reflect DNS requests to a closer DNS server, a useful performance advantage if the visitor's configured DNS server is a considerable distance away. Of greater significance is that it allows general Internet access by the visitor even if the configured DNS server is on a private network, which is now unreachable because the visitor's laptop has been moved from the office.

Redirecting to a DNS server not of the visitor's choosing may work smoothly until the visitor attempts to resolve domain names known only to the real DNS server on the private network. There is, of course, a limit to how well you can hide reality.

One common problem encountered by visitor networks is with a Web proxy on a private network. If the visitor refers to a Web proxy by name, the visitor gateway may choose to respond, inventing an IP address for the proxy and then assuming, by itself, operation of the proxy function. This technique has to be used with some care because hosts often cache DNS responses; these are effectively convenient lies that could end up being carried as "dirty entries" on the visitor's machine for longer than intended.

Rewriting DNS queries and responses does open the opportunity for the service provider to "assume" (some may say "hijack") sites. This opens the door to the possibility that, for example, **yahoo.com** is resolved to an address that is not Yahoo but rather a Web site with an affiliation to the service provider. Although this is a policy and business issue for the service provider, it is likely to irritate quite a number of visitors and reduce the perceived value of the service.

NATs

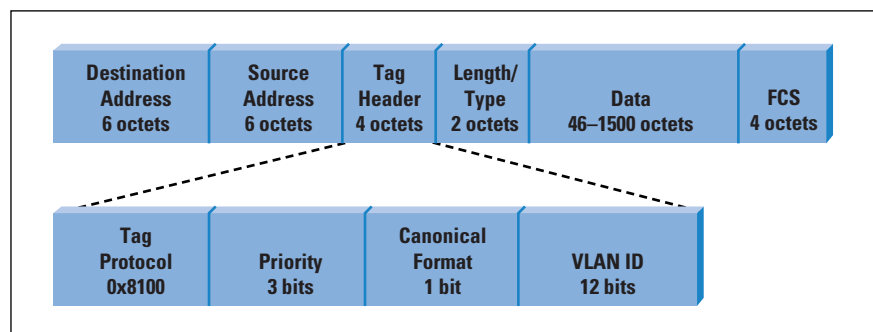
Visitor network gateways frequently use *Network Address Translation* (NAT), and often with port translation, in order to conserve IP addresses by sharing a small address pool with a large number of visitor hosts. In addition, NAT is required by the gateway if the source address used by the visiting host is not routable by the rest of the network back to the visitor gateway. This is almost always the case when the visiting host is using a static preconfigured IP address from another network. The gateway may choose its application of NAT based on policy. For example, two visitors may be configured for DHCP but one is assigned a private "Net 10" (RFC 1918) address that is passed through a NAT while another is assigned a routable address. In practice this flexibility is useful for service in apartments where the visitors are expected to "visit" for months. The service provider may choose to offer tiered services, one with a routable address suitable for the customers to run servers, and another with a private address suitable only for outgoing connections (e-mail, HTTP, and so on).

VLANs

The visitor gateway—modeling its relationship with visiting hosts as a virtual point-to-point link—may attempt to ignore the fact that hosts are on a shared network. However, certain interactions between hosts are inevitable on a shared LAN. For example, if a visitor’s Windows laptop is configured for file sharing with no security enabled, other visitors may see, or worse, have permission to write to, critical files.

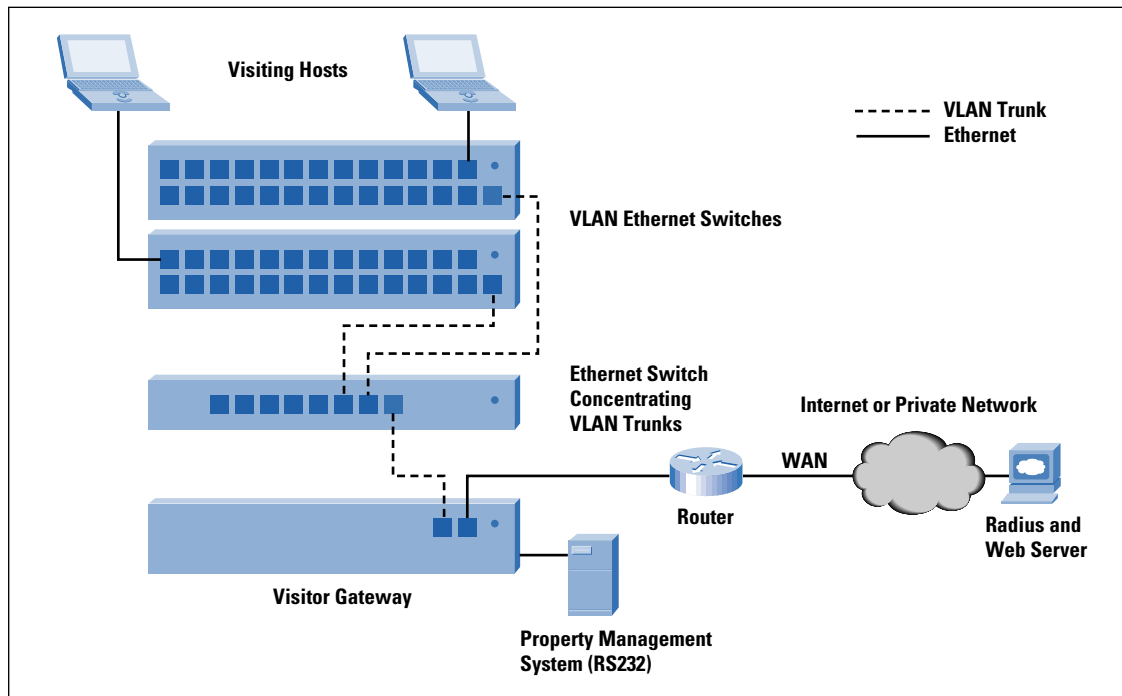
Virtual LANs provide a solution for isolating individual clients. On a wired Ethernet, many modern Ethernet switches can be configured to implicitly treat each port as a member of a different VLAN. For example, port 1 could be on VLAN 11; port 2 on VLAN 12; and so on. The visitor gateway is connected to one or more “trunk” ports that are configured as a member of all VLANs. This effectively allows another level of addressing so the visitor gateway can individually address a single Ethernet network connected to a port. The VLAN switches then act as simple concentrators. If a visiting host attempts to broadcast or multicast, these frames end up only traveling to the gateway and are not seen by other visiting hosts.

Figure 2: VLAN Frame Format



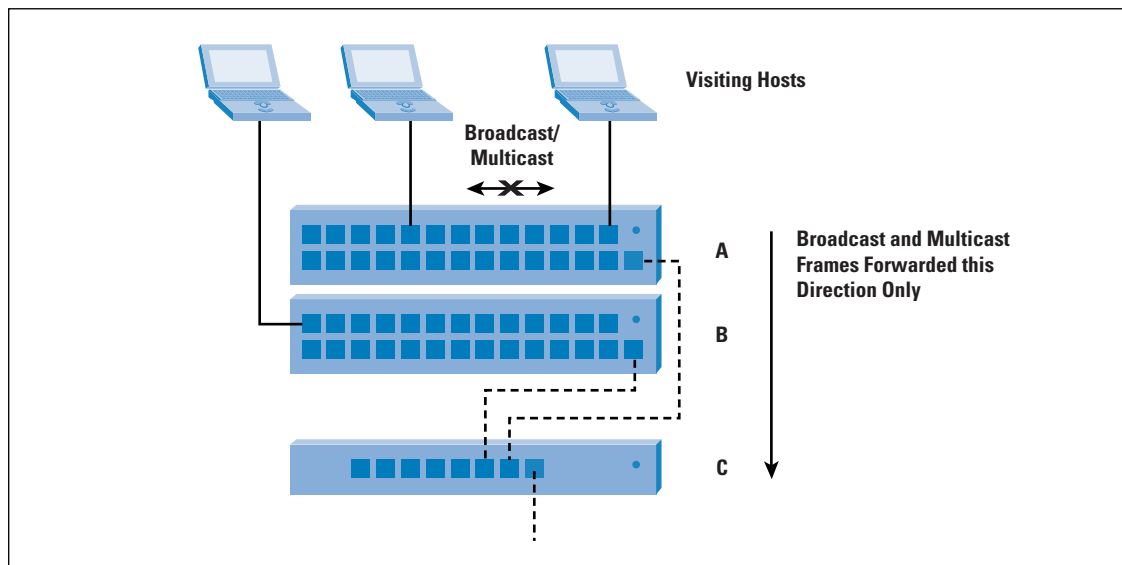
The VLAN frame format is shown in Figure 2. IEEE 802.1q defines the “tagging”^[8]. The VLAN-enabled Ethernet switch can add the appropriate headers to standard Ethernet frames, and it forwards these through the trunk port. Optionally, another Ethernet switch concentrates the trunk traffic and attaches to the visitor gateway, as seen in Figure 3. One potential catch is that some Ethernet switches will not pass the oversized (maximum 1504 octet) VLAN frames; others attempt to be “overly aware” of the VLAN membership rules and insist on configuration of each of the VLANs, a challenging prospect if you are concentrating thousands of ports, each with a unique VLAN identifier.

Figure 3: VLAN Configuration



Some Ethernet switch vendors have implemented a nonstandard technique whereby broadcasts and multicasts are forwarded exclusively to a designated port, the theory being that if a host's broadcast and multicast frames do not get forwarded to other hosts, the hosts effectively will not "see" each other because they do not see ARP requests or higher-layer service advertisements. In some ways, this is simpler than using VLANs and provides some isolation over standard Ethernet networks.

Figure 4: Switch Multicast Blocking



Combinations of these switches with normal ones can lead to some interesting frame forwarding scenarios. For example, as seen in Figure 4, Ethernet switches A and B are each connected to Ethernet switch C. Visiting hosts are attached to the ports on A and B. A and B are designed to have the forwarding restriction described, but C is a normal switch. This means that a broadcast from a visitor connected to A will not be seen by other visitors on A (by nature of the restriction) but it will be forwarded “upstream” to C, which will then forward it to B. Because B received it coming from the upstream, it will forward it to all the visiting hosts on B, causing the isolation technique to fail. A, B, and C all need to have the forwarding restriction.

Web-Based Authentication and Policy

Visitor networks often avail themselves of the one reliable way to converse with a human without additional client software: the Web browser. By selectively reflecting HTTP requests to the local gateway, the gateway can perform or facilitate several operations:

- Authenticate users with traditional username/password—The visitor gateway may, in turn, use a *Remote Access Dial-In User Service* (RADIUS)^[7] authentication request to validate the user.
- Provide links within a “walled garden”—sites that can be visited without authentication—These sites are implemented with either a Web proxy inside the gateway or access control lists effective on the individual visiting host’s virtual interface.
- Gather and validate credit card information through third-party credit card processing Web sites
- Offer visitors Web pages they can use to subscribe to services or to change service parameters

Using the browser can have a significant advantage, even over installed client software. The browser allows a conversation with a human user instead of a software client. This affords the network provider a wide variety of options, such as dealing politely with an authentication rejection, providing additional troubleshooting help, or confirming “conditions of use” before the user accepts charges. It is also a place for offering the user other products and services through Web links.

A central repository for visitor policy and configuration is especially important when a large number of gateways are deployed in disparate physical locations. An interesting option for visitor gateways is for them to learn policy by participating in the exchange of HTTP between the visiting host and an external Web server. The visitor gateway can piggyback the origin and state of a visiting host in a URL and refer the visitor’s browser to a Web site. This origin information when presented to a service selection application running in a provider’s data or operation center allows the application to determine which gateway the visitor is attached to as well as the visitor’s virtual port identification and MAC address.

With the origin information, the service selection gateway can present the visitor with any number of billing, quality of service, or IP addressing options that apply to his/her connection. When the service selection application needs to affect the policy information stored in the visitor gateway, it can use a similar piggyback technique in the return direction.

Accounting

Finding an easy-to-deploy accounting method is crucial for service providers to generate accurate billing. The visitor gateway may send RADIUS accounting records in response to connections and disconnections made by visiting hosts. Disconnections can be determined by *Simple Network Management Protocol* (SNMP) traps from the physical layer devices or by repeated interval polling of the visiting host using ARPs or pings. Because RADIUS has been widely deployed by service providers for dial-in or other networks, it is very possible that the existing accounting system would be able to support the visitor gateway if it, too, offers RADIUS.

In hotels, accounting information can be sent directly to the hotel's *Property Management System* (PMS), causing users to see an access charge on their folio. This is normally accomplished by connecting a standard low-speed serial interface between the visitor gateway and the PMS. The visitor gateway posts the charges by exchanging records with the PMS. A simple record format is used to identify a room and associated charge. Although the format and exchange protocols are usually simple, they are rarely standard. Interfacing to a PMS may require the vendor of the visitor gateway to pay a license fee to the company selling the PMS before it can implement a PMS protocol. Additionally, after implementation, the visitor gateway vendor may need to go through certification for each PMS to which the gateway will be connected. Even if the equipment vendor pays the license, service providers are rarely free to go to a hotel and attach to the hotel PMS—often the service provider is shocked that the hotel insists that they be reimbursed for “interface license fees” charged by the PMS vendor to “enable the protocol.”

The 802.1x Standard and Wireless LANs

Techniques of implementing visitor networks using wireless LANs (WLANs) have been both widely publicized and debated. Wireless 802.11 “hot spots” and the like have been the subject of great publicity because these WLANs are so convenient and cost-effective to deploy that they allow service providers to economically deploy them in areas that would be impractical to serve with wired networks. However, WLANs continue to be the topic of great debate because they have been plagued by the lack of compatibility and weaknesses in security architectures.

The 802.1x standard, recently ratified by the IEEE, holds the best promise in offering a standard authentication scheme for LANs. The 802.1x standard operates with client software. In one sample scenario, the visiting host, also known as the “supplicant,” receives an *Extensible Authentication Protocol* (EAP) request/identity message from the visitor network via an Ethernet switch, a WLAN access point, or a visitor gateway, any of which function as the “authenticator.” The authenticator then relays the client’s identification to an authentication server. The server then decides if the supplicant is to be allowed access and responds appropriately to the authenticator.

With WLANs, the *Wired Equivalent Privacy* (WEP) keys can be loaded as part of the exchange so the client and access points can operate without manual key selection. WEP has been used for several years as a method of encrypting user data over the air interface. Without WEP (or even with it, as we have seen), anyone with a laptop and a receiver can spy on the exchanged traffic^[10, 11].

Microsoft ships an 802.1x client in the standard distribution of Windows XP, an important move forward in making the protocol universal. Other software vendors are shipping or have announced product for older versions of Windows, Macintoshes, Linux, and some PDAs. The 802.1x standard client implementations, however, may need firmware support on the host adapters, and support may never be available on a large number of 802.11 cards already deployed. Furthermore, all 802.1x standards are not alike because they may implement different authentication schemes. Microsoft’s current implementation uses the *Extensible Authentication-Transport Level Security* (EA-TLS) protocol, which requires a *Public Key Infrastructure* (PKI)^[9]. Some critics contend that this creates additional deployment burdens on organizations with small networks. If a common provider, such as Boingo or T-Mobile, provides the visitor network in “hot spots” (that is, cafés and airports), the PKI requirement should not be an issue.

On Ethernet switches, 802.1x implemented directly on the switches may be adequate if the policy for visitor access is relatively simple. For example, if users on a particular network are all trusted employees working for the same business, the work of the authentication/authorization scheme is then to determine whether or not to allow someone to access the network, simply “port on” or “port off.” A more sophisticated approach would allow users to be classified as belonging to a set of classes. On some switches, 802.1x would allow each port to assume a set of VLAN memberships. For example, VLAN 120 would allow unrestricted Internet access, VLAN 119 would restrict access to a set of Web servers, and VLAN 118 would restrict access further to only an authentication server. The authentication system using 802.1x would direct the switch port configuration.

In practice, the control required by visitor networks needs to be far more flexible, and perhaps should be left to the visitor gateway. The gateway, as diagrammed in Figure 1, can control the routing system as well as higher-level protocol proxies based on policy. Besides, leaving the authentication behind the switches allows network implementors the flexibility of using virtually any Ethernet switch, or even other media such as Ethernet framing over xDSL.

The switch-based 802.1x approach, however, may have a significant advantage over the visitor gateway in that after the authentication is out of the way, the Ethernet switch can switch traffic simply at full speed without additional per-packet overhead.

Security Concerns—Better Just to Bootstrap?

Visitor networks are particularly vulnerable to hacking and snooping by virtue of their physical locations, especially if serviced by WLANs. Unfortunately, security is one of the few things that a service provider cannot deliver to visitors without their explicit cooperation and participation. The service providers face a difficult choice to either stay out of the solution or attempt to deliver adequate security through client configuration or special software distribution. The answer is difficult to determine; however, at least two factors to consider are whether the network is wired or wireless, and what the expectations of the visitors will be.

Weaknesses in WEP commonly offered on wireless LAN products have been very well publicized^[10,11]. These weaknesses involve the encryption protocols and the fact that most implementations use manually configured keys. The latter is of little use on a visitor network because the network provider would need to disclose the same keys to everyone. Better proprietary systems have been deployed using PKI, and 802.1x is also a possibility. WEP may be replaced by much stronger *Advanced Encryption Standard* (AES) in *Offset Codebook* (OCB) mode as part of the IEEE 802.1i working group^[12]. No solution has been both standardized and universally deployed. The lack of a standard and universal solution to replace WEP requires that the service provider who chooses another form of security customize a wireless solution. They may need to distribute specialized client software and/or restrict their service to supporting a set of wireless cards and drivers.

Simple Ethernet switches can provide some isolation between ports, but the learning bridge algorithms they use are designed to efficiently deliver Ethernet frames, not provide a secure service. With many switches, it takes one frame with a sham source MAC address to convince the switch to spill someone else's traffic onto the wrong port. "Man in the middle" attacks are often trivial after a visiting host is tricked into sending its traffic somewhere else; the opportunities of doing this to another machine on the same LAN are abundant.

As an end user of a visitor network, trusting an unfamiliar service provider in an unknown environment is a fundamentally insecure process. So, why not let the visitor network provide the basic IP connectivity in order to bootstrap the connection, and then let the visitors themselves implement the security on top? One reason is that unsuspecting users getting hacked at their favorite hotel chain does not bode well for the hotel if the incidents end up in the press. Guests probably feel pretty secure using the hotel phone for a dial-in network connection without any encryption; many also feel secure locking the door with the sliding chain.

One reasonable compromise is matching the security of a dial-in connection. A wired Ethernet, assuming that it cannot be easily coaxed to spill traffic between ports, could present an acceptable risk level. On the other hand, a poorly protected wireless network is like a hotel door without a lock.

If the visitor network offers no protection, then the burden is placed completely on the visitors to implement their own end-to-end security. Using *Virtual Private Network* (VPN) software that implements *IP Security* (IPSec) is one possibility. Unfortunately, even that is not always straightforward, given the complexities with using protocols such as IPSec over NATs^[13]. Other protocols such as *Transport Layer Security* (TLS) and *Secure Shell* (SSH), which operate above the network layer, may be a better option. In addition, several proprietary VPN protocols are designed to tunnel through NATs. Those without any security solution could compromise not only their personal data but also the security of their employer's networks.

Any long-term security solution is going to demand proper client configuration and compatible software. Ultimately, development of standards and client sophistication will make this possible, but in the meantime, we will need to choose between ease of connecting to an insecure network and dealing with the potential multiple layers of authentication and encryption before gaining access. Sadly, faced with this choice and looking forward to a 7 a.m. meeting, the trusty hotel phone and modem jack on the laptop might look pretty inviting.

Summary

Visitor networks allow service providers to provide access in public places. These networks can be implemented in a way that either may or may not require specialized client software on the visiting host. Client software allows service providers to more carefully control the behavior of the visiting host but, at the same time, may limit the user base to those who have the software installed.

Visitor networks often rely on a visitor gateway to perform functions generally not required on a traditional LAN. The gateway, which shares certain characteristics with a NAS, is responsible for routing, address assignment, translation, TCP/UDP redirection, authentication, accounting, and affecting policy.

The visitor gateway exchanges packets with the visiting hosts via LANs. On Ethernet, VLANs are often best suited to visitor networks because they allow the gateway to address each client separately providing the greatest level of isolation compared to other Ethernet options.

WLANs represent an important advance toward the universal deployment of visitor networks in “hot spots.” However, the lack of a common and effective solution may force service providers to choose between ease of access and security. Visitors may choose to implement a VPN or security scheme on top of the raw IP access offered by the visitor network.

References

- [1] L. Mamakos, K. Lidl, J. Evarts, D. Carrel, D. Simone, R. Wheeler, “A Method for Transmitting PPP over Ethernet (PPPoE),” RFC 2516, February 1999.
- [2] W. Townsley, A. Valencia, A. Rubens, G. Pall, G. Zorn, B. Palter, “Layer 2 Tunneling Protocol, L2TP,” RFC 2661, August 1999.
- [3] S. Glass, T. Hiller, S. Jacobs, C. Perkins, “Mobile IP Authentication, Authorization, and Accounting Requirements,” RFC 2977, October 2000.
- [4] IEEE Standards for Local and Metropolitan Area Networks: Port-Based Network Access Control, IEEE Std 802.1X-2001, June 2001.
- [5] W. Simpson, “The Point-to-Point Protocol (PPP),” STD 51, RFC 1661, July 1994.
- [6] R. Droms, “Dynamic Host Configuration Protocol,” RFC 2131, March 1997.
- [7] A. Rubens, W. Simpson, S. Willens, C. Rigney, “Remote Authentication Dial-In User Service (RADIUS),” RFC 2058, January 1997.
- [8] IEEE standard for local and metropolitan area networks: Virtual Bridged Local Area Networks, IEEE Std 802.1Q-1998.
- [9] B. Adoba, D. Simon, “PPP EAP TLS Authentication Protocol,” RFC 2716 (experimental), October 1999.

- [10] N. Borisov, I. Goldberg, D. Wagner, "Intercepting Mobile Communications: The Insecurity of 802.11,"
<http://www.isaac.cs.berkeley.edu/isaac/wep-faq.html>
- [11] E. Danielyan, "IEEE 802.11," *The Internet Protocol Journal*, Volume 5, Number 1, March 2002.
- [12] D. Whiting, R. Housley, "AES Encryption & Authentication Using CTR Mode with CBC-MAC," Status of Project IEEE 802.11i, July 2002,
<http://grouper.ieee.org/groups/802/11/Documents/DocumentHolder/2-001.zip>
- [13] B. Adoba, "IPSec-NAT Compatibility Requirements," Internet-Draft,
<http://www.ietf.org/internetdrafts/draft-ietf-ipsec-nat-reqts-01.txt>,
March 1, 2002.

DORY LEIFER is a principal with DEL Communications Consulting, Inc. He had co-founded PublicPort in 1998, an Ann Arbor, Michigan, startup that developed one of the first visitor network gateways. After Tut Systems acquired PublicPort, he held a director of marketing position with Tut until 2001. Leifer spent 11 years with the University of Michigan and Merit Network and during that time contributed to the *Internet Engineering Task Force* (IETF). He has taught tutorials in access technologies for various seminars and tutorials, including NetWorld+Interop. He holds a B.S. in Computer Science from Rensselaer Polytechnic Institute and an M.S.E. in Industrial and Operations Engineering from the University of Michigan. Leifer currently resides in the San Francisco Bay Area and can be reached by e-mail at leifer@del.com

An Architecture for Securing Wireless Networks

by Gregory R.Scholz, Northrop Grumman Information Technology

Wireless networks are described as both a boon to computer users as well as a security nightmare; both statements are correct. The primary purpose of this article is to describe a strong security architecture for wireless networks. Additionally, the reader should take from it a better understanding of the variety of options available for building and securing wireless networks, regardless of whether all options are implemented. The security inherent with IEEE 802.11 wireless networks is weak at best. The 802.11 standard provides only for *Wired Equivalent Privacy*, or WEP, which was never intended to provide a high level of security^[1]. For an overview of 802.11 and WEP, see reference^[2]. Wireless networks can, however, be highly secure using a combination of traditional security measures, open standard wireless security features, and proprietary features. In some regard, this is no different than traditional wired networks such as Ethernet, IP, and so on, which have no security built in but can be highly secure. The design described here uses predominantly Cisco devices and software. However, unless explicitly stated to be proprietary, it should be assumed that a described feature is either open standard or, at least, available from multiple vendors.

Customer needs

Customer needs range from highly secure applications containing financial or confidential medical information to convenience for the public “hot spot” needing access to the Internet. The former requires multiple layers of authentication and encryption that ensures a hacker will not be able to successfully intercept any usable information or use the wireless network undetected. The latter requires little or no security other than policy directing all traffic between the wireless network and the Internet. Security is grouped into two areas: maintaining confidentiality of traffic on the wireless network and restricting use of the wireless network. Some options discussed here provide both, whereas others provide for a specific area of security.

The level of security required on the wireless network is proportional to the skill set required to design it. However, the difficulty of routine maintenance of a secure wireless network is highly dependant on the quality of the design. In most cases, routine maintenance of a well-designed wireless network is accomplished in a similar manner to the existing administrative tasks of adding and removing users and devices on the network. It is also assumed that security-related services such as authentication servers and firewall devices are available on the wired network to control the wireless network traffic.

It is not necessarily the case that one can see the user or device attempting to use the wireless network. This is the most alarming part of wireless network security. In a wired network, an unauthorized connected host can often be detected by link status on an access device or by actually seeing an unknown user or device connected to the network. The term “inside threat” is often used to refer to authorized users attempting unauthorized access. This is the inside threat because they exist within the boundaries that traditional network security is designed to protect. Wireless hackers must be considered more dangerous than traditional hackers and the inside threat combined because if they gain access, they are already past any traditional security mechanisms. A wireless network hacker does not need to be present in the facility. This new inside threat may be outside in the parking lot. *War Driving*^[3] is the new equivalent to the traditional war dialing. All that is required to intercept wireless network communications is to be within range of a wireless access point inside or outside the facility.

Physical Wireless Network

In a highly secure environment, a best practice is to have the wireless access points connect to a wired network physically or logically separate from the existing user network. This is accomplished using a separate switched network as the wireless backbone or with a *Virtual LAN* (VLAN) that does not have a routing interface to pass its traffic to the existing wired network. This network terminates at a *Virtual Private Network* (VPN) device, which resides behind a firewall. In this manner, traffic to and from the wireless network is controlled by the firewall policy and, if available, filters on the VPN device. The VPN device will not allow any traffic that is not sent through an encrypted tunnel to pass through, with the exception of directed authentication traffic described later. With this model, the wireless clients can communicate among themselves on the wireless network, but there is no access to internal network resources unless fully encrypted from the wireless client to the VPN. This design may be further secured by configuring legitimate wireless-enabled devices to automatically initiate a VPN tunnel at bootup and by enabling a software firewall on the devices that does not allow communication directly with other clients on the local wireless subnet. In this manner, all legitimate communication is encrypted while traversing the wireless network and must be between authenticated wireless clients and internal network resources.

Authentication

Many security measures available relate to access controlled through individual user authentication. Authentication can be accomplished at many levels using a combination of methods. For example, Cisco provides *Lightweight Extensible Authentication Protocol* (LEAP)^[4] authentication based on the IEEE 802.1x^[5] security standard. LEAP uses *Remote Authentication Dial-In User Service* (RADIUS)^[6] to provide a means for controlling both devices and users allowed access to the wireless network.

Although LEAP is Cisco proprietary, similar functionality is available from other vendors. Enterasys Networks, for example, also uses RADIUS to provide a means for controlling *Media Access Control* (MAC) addresses allowed to use the wireless network. With these features, the access points behave as a kind of proxy, passing credentials to the RADIUS server on behalf of the client. When these features are properly deployed, access to the wireless network is denied if the MAC address of the devices or the username does not match an entry in the authentication server. The access points in this case will not pass traffic to the wired network behind them. For security, the authentication server should be placed outside the local subnet of the wireless network. The firewall and VPN devices must allow directed traffic between the access points and the authentication server further inside the network and only to ports required for authentication. This design protects the authentication server from being attacked directly.

In addition to authenticating users to the wireless network, the VPN authentication and standard network logon can be used to control access further into the wired network. In this solution, the VPN client has the ability to build its tunnel prior to the workstation attempting its network logon, but after the device has been allowed on the wireless network. After the tunnel is built, specific rules on the VPN and the firewall allow the traditional network logon to occur. A robust VPN solution also treats the users differently based on the group to which they are assigned. Different IP address ranges are assigned to each group, allowing highly detailed rules to be created at the firewall controlling access to internal network resources based on user or group needs. The policy on the firewall must be as specific as possible to restrict access to internal resources to only those clients for whom it is necessary. Building very specific policy for users' access will also allow an *Intrusion Detection System* (IDS) to better detect unauthorized access attempts.

Encryption

LEAP also provides for dynamic per-user, per-session WEP keys. Although the WEP key is still the 128-bit RC4 algorithm proven to be ineffective in itself⁷, LEAP adds features that maintain a secure environment. Using LEAP, a new WEP key is generated for each user, every time the user authenticates to use the wireless network. Additionally, using the RADIUS timeout attribute on the authentication server, a new key is sent to the wireless client at predetermined intervals. The primary weakness of WEP is due to an algorithm that was easy to break after a significant number of encrypted packets were intercepted. With LEAP, the number of packets encrypted with a given key can be tiny compared to the number needed to break the algorithm.

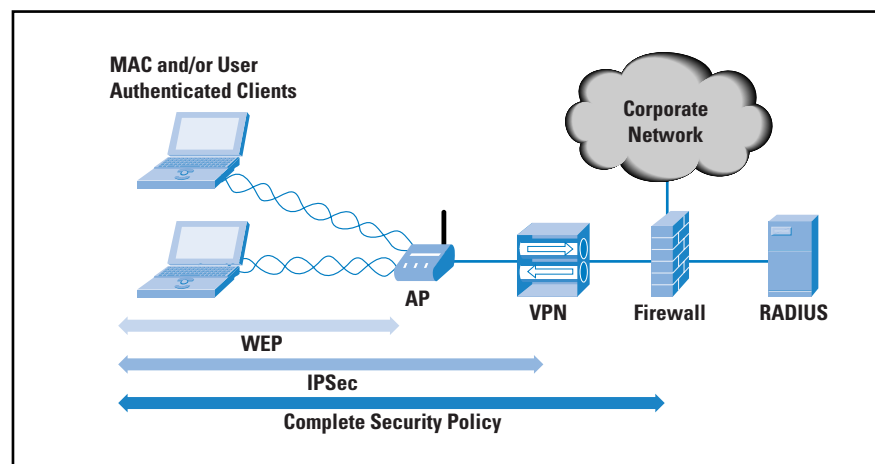
When using LEAP for user and device authentication, WEP encryption is automatically enabled and cannot be disabled. However, if added security is needed, a VPN, as described earlier, can provide any level of encryption desired. Using a VPN as the bridge between the wired and wireless network is recommended regardless of the underlying vendor or technology used on the wireless network. *IP Security* (IPSec) is a proven, highly secure encryption algorithm available in VPNs. By requiring all wireless network traffic to be IPSec encrypted to the VPN over the WEP-encrypted 802.11 Layer 2 protocol, any data passed to and from wireless clients can be considered secure. All traffic is still susceptible to eavesdropping, but will be completely undecipherable.

Aside from WEP and LEAP, some vendors provide other forms of built-in security. Symbol Technologies' Spectrum24 product provides Kerberos encryption when combined with a Key Distribution Center. Kerberos is more lightweight than IPSec and, therefore, may be better suited to certain applications such as IP phones or low-end *personal digital assistants* (PDAs). Other methods of automating the assignment and changing of WEP keys are also available, such as Enterasys' Rapid-Rekey^[8]. Wireless vendors have realized that security has become of critical importance and most, if not all, are working on methods for conveniently securing wireless networks. When available, most vendors seemingly prefer to use open-standard, interoperable security mechanisms with proprietary security being additionally available.

Bringing it all together

Numerous options are available to secure a wireless network. A highly secure design will include, at a minimum, an authentication server such as RADIUS, a high-level encryption algorithm such as IPSec over a VPN, and access points that are capable of restricting access to the wireless network based on some form of authentication. When all the security options are tied together, the wireless network requires explicit authentication to allow a device and the user on the wireless network, the traffic on the wireless network is highly encrypted, and traffic directed to internal network resources is controlled per user or group by an access policy at the firewall or in the VPN.

Figure 1: A Highly Secure Wireless Network



There is no substitute for experience and research when designing a network security solution. Using network security and design experience to exploit available technologies can further increase security of a wireless network. For example, grouping users into IP address ranges based on access requirements allows firewall access policy to help restrict unnecessary access. This can be accomplished using *Dynamic Host Configuration Protocol* (DHCP) reservations, assigning per-user or -group IP address ranges to the VPN tunnels or statically assigning addresses. Using a centralized accounts database for all authentication helps avoid inadvertently allowing an account that has been disabled in one part of the network to access resources through the wireless network. To use an existing user database for authentication while providing for dynamic WEP keys, use a LEAP-enabled RADIUS server that has the ability to query another server for account credentials. As with most network designs, a solid understanding of the available technologies is paramount to achieving a secure environment.

Utilizing all the security described in this article would yield the following design. When a device first boots up, it receives an IP address within a specified range on a segregated portion of the network. This IP range is based on the typical usage of the device and is most useful for machines dedicated to specific applications. As a user attempts to log onto a wireless device, a RADIUS server authenticates both the MAC address and the username of the device. If the user authentication is successful, access is granted within the wireless network. In order for traffic to leave the wireless network to access other network resources, a VPN tunnel must be established. Again, the IP address assigned to the tunnel can be controlled based on individual user authentication to help enforce access policy through the firewall. When the tunnel is established, firewall access policy will restrict access to resources on the network. Most, if not all, of the authentications required may be automated to use a user's existing network logon and transparently complete each authentication. This is not the most secure model, but it would be as secure as any single signon environment.

Summary

A secure wireless network is possible using available techniques and technologies^{[8] [9] [10]}. After researching needs and security requirements, any combination of the options discussed here, as well as others not discussed, may be implemented to secure a wireless network. With the right selection of security measures, one can ensure a high level of confidentiality of data flowing on the wireless network and protect the internal network from attacks initiated through access gained from an unsecured wireless network. At a minimum, consider the current level of network security and ensure that the convenience of the wireless network does not undermine any security precautions already in place in the existing infrastructure.

References

- [1] “Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications,” IEEE Standard 802.11, 1999 Edition.
- [2] “802.11,” Edgar Danielyan, *The Internet Protocol Journal*, Volume 5, Number 1, March 2002.
- [3] “War Driving,” Andrew Woods, <http://www.personaltelco.net/index.cgi/WarDriving>, last viewed August 11, 2002.
- [4] “Cisco Aironet® Product Overview,” Cisco Systems, http://www.cisco.com/univercd/cc/td/doc/product/wireless/airo_350/350cards/pc350hig/pc_ch1.htm, last viewed August 11, 2002.
- [5] “IEEE Standard for Local and Metropolitan Area Networks—Port-Based Network Access Control,” IEEE Standard 802.1X, 2001.
- [6] “Remote Authentication Dial-In User Service,” C. Rigney, S. Willens, A. Rubens, and W. Simpson, IETF RFC 2865, June 2000.
- [7] “Security of the WEP Algorithm,” Nikita Borisov, Ian Goldberg, and David Wagner, <http://www.isaac.cs.berkeley.edu/isaac/wep-faq.html>, last viewed August 11, 2002.
- [8] “802.11 Wireless Networking Guide,” Enterasys Networks, June 2002, http://www.enterasys.com/support/manuals/hardware/4042_08.pdf, last viewed August 11, 2002.
- [9] “Wireless LAN Security in Depth,” Sean Convery and Darrin Miller, Cisco Systems, http://www.cisco.com/warp/public/cc/so/cuso/epso/sqfr/safwl_wp.htm, last viewed August 11, 2002.
- [10] “Making IEEE 802.11 Networks Enterprise-Ready,” Arun Ayyagari and Tom Fout, Microsoft Corporation, May 2001, <http://www.microsoft.com/windows2000/docs/wirelessec.doc>, last viewed August 11, 2002.

GREGORY SCHOLZ holds a BS in Computer and Information Science from the University of Maryland. Additionally, he has earned a number of certifications from Cisco and Microsoft as well as vendor-neutral certifications, including a wireless networking certification. After serving in the Marine Corps for six years as an electronics technician, he continued his career working on government IT contracts. Currently he works for Northrop Grumman Information Technology as a Network Engineer supporting Brook Army Medical Center, where he performs network security and design functions and routine LAN maintenance. He can be reached at: gscholz@wireweb.net

Opinion: The ISP—The Uncommon Carrier

by Geoff Huston, Telstra

There is a long-standing role in the communications industry where a provider of public carriage services undertakes the role of a *common carrier*. What's so special about the role of a common carrier, and why is this role one that is quite uncommon in the *Internet Service Provider* (ISP) world?

Side comment: There once was a time when you could not trust the messenger. There once was a time when not only did you pay to have your message sent, but you paid to *receive* messages. And there was no guarantee that the message would not be read by the messenger. The contents of your note could have been used to determine how much the receiver should pay for the message. Your message could have been copied and sold to other parties. If you can't trust the messenger, then communications becomes a risky business.

The Messenger

Throughout history the position of a messenger has been a mixed blessing. To be the bearer of bad news was not an enviable role, and rather than being rewarded for the effort of delivering the message, the messenger might have been in dire straits, given the level of wrath of the recipient. The option of reading the message before delivering it could be seen as a personal survival strategy, as well as being a prudent business move—bad news could be discarded immediately, whereas good news could have the potential of extracting a higher delivery fee from the recipient. Although this scenario would have been good for the messenger, such a mode of operation was not beneficial to all. For the parties attempting to use the messenger service, message delivery could be a very haphazard affair. The message might or might not get delivered, the delivery time was variable, as was the cost of delivery, and if the message itself was intended to be a secret, then one could confidently anticipate that the messenger would compromise this secrecy.

The Common Carrier

For a communications network to be truly useful, numerous basic attributes must be maintained. These include predictability, so that a message passed to a communications carrier is delivered reliably to the intended recipient. Integrity is also necessary, because a message must not be altered by the carrier in any way. Privacy is also an essential attribute, because the message must not be divulged to any party other than the intended recipient, nor should even the existence of the message be made known to any other party. And above all there must be a solid foundation for trust between the carrier and the clients of the service. So in this form of social contract, what does the carrier get in return?

Apart from payment for the service, the carrier is absolved from liability regarding the content of the messages, and from the actions of the customers of the service. This form of social contract is the basis for the status of a common carrier.

It may have taken some time, but this role is well understood by the public postal network. And as many national postal operators encompassed the role of national telephone carrier, the common carrier role has been an integral part of the public telephone network.

The ISP's Role

But in the world of the ISP the position of common carrier is very uncommon indeed.

There once was a time when folk did not need to encrypt their letters nor speak in scrambled code to undertake a private conversation. The assumption, made law in many countries, was that the entity entrusted with public communications, the common carrier, was barred from deliberately inspecting the contents of the plain transmission, and various dire penalties were in place if a public carrier's employees or agents divulged anything they may have learned by virtue of being public carriers. Various measures were put in place to execute interception and monitoring, but these measures required due process and reference to some law enforcement agency and also the judiciary to ensure that the rights of the public user were adequately safeguarded.

The issues of the role of a common carrier and the current role of an ISP are clearly seen when looking at the reactions to unsolicited commercial e-mail, or spam. Every day ISPs receive strident demands of the form: "One of your users is sending unsolicited messages—disconnect him now!" Internet users are, in effect, holding the ISP responsible for the actions of its customers. A similar expectation of the ISP's responsibility for the actions of its customers is seen in response to various forms of hacking, such as port scanning. Similar messages are sent to ISPs, demanding the immediate disconnection of those customers who are believed to be originating such malicious attacks. From a small set of complaining messages some years back, the volume of such demands for ISP action is now a clamor that is impossible for any ISP to ignore.

What should the ISP do? Many responsible ISPs see it as appropriate to conduct an investigation in response to such complaints. ISPs often include provisions in their service contracts with their customers to allow them to terminate the service if they believe that their investigation substantiates the complaints on the basis of a breach of contract. When disconnected, such customers are often blacklisted by the ISP to ensure that they cannot return later and continue with their actions. Surely this is an appropriate response to such antisocial actions?

This may be the case, but it is not necessarily consistent with the role of the ISP as a common carrier. A common carrier is not a law enforcement agency, nor is it an agent of the judiciary. It may be entirely appropriate for a common carrier to investigate, under terms of strict privacy, a customer's activities and inspect the contents of traffic passed across the network if it has reasonable grounds to suspect that the integrity of the network itself is under threat. Equally, it is probably inappropriate for a common carrier to extend the scope of such investigations on the basis of external allegations of activities that are not related to the integrity of the service itself.

The assumption that an ISP is, in some way, responsible for the actions of its customers has been extended further in some countries, such that the ISP is, in part, responsible for the content carried over its network, including content that originates with a customer of its service. This expectation that ISPs should actively control and censor content passed across their network is not just an expectation of some Internet users. This expectation appears in numerous legislative measures enacted in many countries. The *Communications Decency Act* in the United States legislature is an example of such an expectation of the active role of the ISP in controlling content passed across its network.

Who Will You Call?

Perhaps the issue here is one of expediency. Where can a user direct a complaint after receiving yet another piece of unsolicited, and possibly highly offensive, e-mail, apart from the ISP of the sender of the message? Where else can users direct a complaint after being the subject of yet another port scan of their system, but to the ISP? And what else can an ISP do in response? The ISP often has little choice but to investigate such complaints in good faith, and take corrective action if the complaint is substantiated. In the absence of any effective regulatory framework that would allow such investigations to be undertaken by an appropriate external agency, the ISP is in a difficult position.

Whereas it may be the correct common carrier position to disclaim all responsibility for the actions of its customers together with the content passed across its network, to ignore such complaints marks the ISP as a haven for such antisocial activities. Adopting such a position often has a negative impact on the ISP's ability to interconnect with other ISPs, because ISPs also tend to hold each other responsible for the actions of their customers and the content passed across their network. ISPs tend to avoid extending interconnection services to those ISPs that disclaim any such responsibility. So the expedient response is for the ISP to assume some level of responsibility for its customers and the content of its network and act accordingly.

But short-term expedient measures should not be confused with long-term effective solutions. The problem with these short-term responses lies in the uniquely privileged position of the carrier. Even rudimentary forms of data mining of each customer's communications patterns and the content of their communications can yield vast quantities of valuable information. Such information can allow a carrier to discriminate between customers, compromise the integrity of the customer's use of the network, and actively censor the content passed across the network. Positions of privilege without accompanying checks and balances are readily abused. There is already the widespread expectation and acceptance that an ISP has the ability and duty to inspect network content and monitor customers' activities with respect to various forms of anti-social and often malicious activities. But how can checks and controls be enforced such that the information gained through such monitoring activities is not used for other purposes? Such monitoring is not without cost, and the option of recouping some revenue to balance this expenditure by regarding this information as a business asset is always present. The regulatory impost of a common carrier role is intended to be an economically efficient response to this issue. The common carrier role is intended to reduce the social power of public carriers and protect the public's open, uncensored, and equal access to the carrier's services.

It is often said that the road to hell is paved with the best of intentions—that the ultimate outcome of the solution is potentially far worse than the immediate problem being addressed. The ultimate outcome of erosion of the common carrier role is that public users of a public communications service can confidently expect their communications to be monitored, potentially stored and cross referenced, and possibly later acted on.

Public Policy

Today the short-term expedient measures abound. There is enormous pressure on ISPs from both the Internet's user base and numerous legislatures to take an active position of being responsible—and liable, for the content on the networks and the actions of their clients. If left unchecked, this will have severe longer-term consequences for free speech, basic personal privacy, and uncensored, nondiscriminatory, universal access to the Internet. And when the user base comes to recognize the debased value of such a compromised communications system, they will inevitably look to other means of communication that have retained their essential integrity as a common carriage service.

Perhaps it is time for the debate regarding the role and responsibilities of an ISP to be placed on the agenda of public policy makers. Perhaps it is time to recognize that ISPs are indeed common carriers, and that they have a clearly bounded set of responsibilities with respect to both content and the actions of clients of the service.

Perhaps it is time to consider how best to enforce social norms on the Internet without compromising the basic integrity of the carrier as a neutral party to the content being carried across the network. Perhaps it is time to recognize that in this domain the Internet is not entirely novel, and what we have learned from a rich history of carriage provision in society has direct relevance to the Internet today.

The Internet is simply too valuable a communications service to have its long-term potential as a universal communications service mindlessly destroyed on the altar of short-term expediency.

Disclaimer: I am by profession neither a lawyer nor a public policy maker. However, by virtue of working in the ISP industry, I have an increasing level of interest in the activities of these folk, for the reasons outlined above. I should also note that personal opinion comes in many forms. The above is one such form.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. Huston is currently the Chief Scientist in the Internet area for Telstra. He is also a member of the Internet Architecture Board, and is the Secretary of the APNIC Executive Committee. He is author of *The ISP Survival Guide*, ISBN 0-471-31499-4, *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, and coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. All three books are published by John Wiley & Sons.
E-mail: gih@telstra.net

Letters to the Editor

ENUM Ole,

I was looking at the June 2002 issue of *The Internet Protocol Journal*, and noticed what might be a misprint. In the story on ENUM, the next-to-last paragraph on page 21 has a sentence reading:

North America has the .164 country code of “1,” implying that under ENUM there is a single DNS domain for ENUM, namely **1.e164.arpa.**

I suspect it should read “... there is a single DNS domain for North America...” or something like that. (The “.164” should probably also be “E.164”—you don’t refer to it as just “.164” elsewhere in the article.)

A more substantive comment on Marshall Rose’s BEEP article in the same issue: It was a good overview, but I would have liked to see a mention of which application protocols are likely to use BEEP (assuming that none has already) in the near future. The middle of page 11 explains why the IETF thinks this is a good idea and why new application protocols need BEEP, but it was hard to tell whether it actually is being actively considered for use by any IETF working group.

Overall, I liked the issue, and particularly Peter Salus’s review of Padlipsky’s book—I came across it in the late 1980s, and actually met Michael sitting in a hallway at one of the Interop conferences before they got too big for Silicon Valley and I stopped attending. I still remember some of his cartoons and slogans (e.g., something to the effect “... the ITU is planning to have an 11-layer model because it’s a sacred number in Bali...”). I’ve also found the articles in some of the other recent issues of the IPJ—e.g., the articles on wireless LANs (particularly the discussion of security issues) and code signing/mobile code in the March 2002 issue—very helpful, and have pointed colleagues to them.

Best wishes.

—Eric M. Berg
Managing Director,
Technology Forecast Publications
PricewaterhouseCoopers Technology Centre
Eric.Berg@us.pwcglobal.com

Geoff Huston responds:

While we all try hard to eliminate various errors in manuscripts prior to publication, there are always a few author-mishaps that manage to sneak past the eagle eyes of the editor, and this is one of them.

The offending sentence should read:

North America has the .E164 country code of “1,” implying that under ENUM there is a single DNS domain for ENUM in North America, namely **1.e164.arpa.**

Thanks for pointing this out.

—Geoff

More about ENUM Ole,

In the June 2002 issue of IPJ (Volume 5, Number 2), Geoff Huston wrote an interesting article about ENUM. The technical side of ENUM (using DNS to map E164 numbers to services) seems rather straightforward. But its implications on both technical and social issues are much more complex and (in my opinion) interesting. I am not an expert on the subject, but I’d like to share a few thoughts about this. First, two technical issues come to mind.

The first one is about the use of the *Domain Name System* (DNS). The DNS has been very successful as a distributed replicated database of hostname-to-IP address (and reverse) mappings. Will it be able to handle gracefully all the stuff people intend to put in it? This is not certain, as shown by ICANN’s cautious attitude concerning the creation of new Top Level Domains. Content Distribution Networks, for example, often use lots of domain names with short TTLs, reducing the effectiveness of DNS caching (Geoff mentions this caching issue for ENUM). After all, DNS stands for “Domain Name System,” not “General Purpose Infinitely Scalable Distributed Dynamic Database.”

The second issue is about the status of addresses and names in the Internet. Simplifying things, we can say the following happens when somebody wants to access an Internet service with an E.164 number: The E.164 number is translated into a DNS name, and a DNS lookup gives back an URI. If the URI is a simple URL, the domain name in the URL is DNS-looked-up for an IP address, and then packets are sent to that IP address. If the URI is not a simple URL (such as a URN), some other resolving process implying the DNS occurs anyway.

That makes two levels of indirection, but, moreover, creates an “interesting” situation: IP addresses are “addresses,” i.e., network-friendly identifiers, whose structure is tied to the network topology.

Such identifiers are not user friendly, so user-friendly identifiers called “names” have been created, and a “domain name system” set up to translate names into addresses. E.164 numbers are really telephone addresses. They are tied to the telephone network topology and are surely not user friendly. There are no user-friendly names in the telephone system.

The strange thing is that with ENUM, E.164 numbers are not linked anymore to the network topology, but rather become names intended for user usage. In a sense, they even are “meta names,” since they translate to DNS names (that translate to addresses). But they obviously have not become user-friendly in the process.

I must admit I oversimplify a bit since I don’t distinguish between names and addresses identifying level 3 (network) resources (i.e., hosts) and those identifying level 7 (application) resources (e-mails, Web pages, etc.), but this doesn’t invalidate the idea.

Addresses are what the network needs, and names are what the users need. This brings me to the politics aspects of ENUM: who administers/controls/owns the namespace? A namespace is only partly technical; defining a namespace includes defining how and by whom the namespace is operated. The DNS is technically a big success, but the politics side is controversial, as shown by domain-name disputes or the setting up of alternative domain-name systems. It seems that social aspects are often more difficult to deal with than technical issues are to solve.

When I was studying networking we were taught how the technical differences between the Internet and the telephone network took their roots into a fundamental difference of culture. Now that the Internet culture seems to have won on the technical aspect (IP over broadband ISDN), wouldn’t it be a strange outcome for the Internet namespace to be owned by telephone companies?

To conclude, I think this ENUM stuff shows that the Internet community really needs to work on the namespace issue, to ensure a technically and socially sound namespace for the Internet.

—*Christophe Deleuze, Ph.D.*
R&D Senior Engineer
ActiVia Networks

Christophe.Deleuze@ActiVia.net

Book Review

Carrier-Scale IP Networks

Carrier-Scale IP Networks: Designing and Operating Internet Networks, edited by Peter Willis, ISBN 0-85296-982-1, The Institute of Electrical Engineers, London, United Kingdom, 2001

My heart jumped when I saw the nondescript brown box, about the thickness of a book, sitting by the receptionist. It was finally here! I had waited almost two months in great anticipation for this book to show up. Was it going to be the all-encompassing handbook for the network designers, operators, and managers in large-scale IP environments? The first few lines in the text indicated that it just might be: “The aim of this book is to give the reader an understanding of all the aspects of designing, building and operating a large global IP network.”

The definition of “large-scale” as given by the author and for the purposes of this review follows: Provides services for millions of end users, high-speed (greater than 100 Mbps) transit services, and is reliable, scalable, and manageable.

One thing to keep in mind is the way this book was constructed. The 16 chapters had 29 authors. Almost all authors came from some area of British Telecom (BT) and all were subject matter experts in the chapter they wrote. The 16 chapters are grouped roughly into four sections: Designing and building IP networks, transmission and access networks, operations, and development of future networks. Sadly, all of this is squeezed into 293 pages.

Designing and building IP networks

For the reader new to designing and building large-scale IP networks, the first few chapters are gold. For the reader already experienced in this area, it may bring back nostalgic feelings for the good old days of exponential growth. A lot of ground is covered, including the obligatory overview of IP, sufficient enough to give a nontechnical person the key concepts of IP routing, but can be skipped by those with even basic knowledge in this area. The examples given throughout this chapter (and the rest of the book) come directly from the design of BT’s and Concert’s backbone. A whole chapter, “The Art of Peering,” not to be mistaken for an excellent paper of the same name^[1], gives excellent key concepts in peering. Some coverage is even given to the logistics and difficulties in building points of presence globally, going so far as to mention earthquake bracing for equipment bays.

The next set of chapters give the reader detail about the transmission network (for some, be prepared to think *Synchronous Optical Network* [SONET] when you read *Synchronous Digital Hierarchy* [SDH]), and access networks, including various forms of broadband, wireless, dial, and satellite.

The technical information was squeezed into these chapters, not enough for a good technical treatise, but enough to give readers good grounding in a technology that is unfamiliar to them. The coverage was closer to being marketing material. These chapters alone are not enough to bring those new to the field up to speed if they are to design or operate such a network.

BT opened itself up and gave us a view into the operations of its network. Individuals who have worked in an environment like this will find something familiar. We get to see how BT structures the people, processes, and technologies. This is something that is not usually open to inspection by people outside of an organization. Planning and developing the operations side of the house is a difficult job. These chapters may give a kick-start to those coming into such a role.

I was disappointed with the two final chapters. Of course anything listed as being “the future” will one day become the present, but I digress. These two chapters seem like the odd couple that just did not fit with the rest of the chapters. The first chapter is on Traffic Engineering. It is really a primer on *Multiprotocol Label Switching Traffic Engineering* (MPLS TE). The second chapter covers *Virtual Private Networks* (VPNs), both the MPLS and *IP Security* (IPSec) types.

Recommendation

The authors set out with a lofty goal, and did not quite hit the mark. This book would be appropriate for someone trying to get a feel for what goes on inside of a carrier-scale network. People already in the business would be better served by just paying attention to what goes on around them.

Perhaps a small focused group could set out to create a book (or should I say tome) covering the elements of design, the foundation of support, and the basics of management. Something timeless is required here, independent of the protocol du jour, to develop the next generation of competent netheads.

—Kris Foster

kris.foster@telus.com

- [1] “The Art of Peering: The Peering Playbook,” William B. Norton, Equinix

Stephen Wolff receives Postel Service Award

In June 2002, Internet pioneer Stephen Wolff was honored by the *Internet Society* (ISOC) for his significant contributions on behalf of the Internet. A founding member of the ISOC, Wolff is considered one of the “fathers of the Internet” and was directly involved with its development and evolution.

Wolff received the *Postel Service Award*, named for Dr. Jonathan B. Postel, an Internet pioneer and head of the organization that administered and assigned Internet names, protocol parameters, and *Internet Protocol* (IP) addresses. He was the primary architect behind what has become the *Internet Corporation for Assigned Names and Numbers* (ICANN), the successor organization to his work. The recipient of the award receives a \$20,000 cash honoraria.

“We are pleased to recognize Steve with the Postel Award,” said ISOC President/CEO Lynn St.Amour, “especially as his contributions are well known to ISOC, having previously been commended by ISOC’s board for helping transform the Internet from an activity serving the particular goals of the research community to a worldwide enterprise which has energized scholarship and commerce in dozens of nations.”

The 1994 commendation from the ISOC board also states that “The personal leadership of Dr. Wolff, often under conditions of public controversy, has been an indispensable ingredient in surmounting a daunting array of technical, operational and economic challenges. His extraordinary commitment to the growth and success of the Internet reflect the highest standard of service to the networking community and command our respect and admiration.”

As Director of the Division of Networking and Communications Research and Infrastructure at the US National Science Foundation, he was responsible for NSNET, the *National Research and Education Network* (NREN), and for NSF’s support of basic research in networking and communications. While at the NSF he was among the founders of the interagency and international research networking management and advisory structure whose descendants today include the Large-scale Networking (LSN) working group and the PITAC.

Wolff left the federal government and joined Cisco Systems, Inc. in 1995, where he works in the University Research Program—Cisco’s program supporting academic investigators with unrestricted grants for research on computer networks.

Wolff was educated at Swarthmore College, Princeton University, and Imperial College. He taught electrical engineering at the Johns Hopkins University for ten years and subsequently spent fifteen years leading a computing- and network-related research group at the U.S. Army Research Laboratory. In 1983 he took a sabbatical half-year as a Program Director in the Mathematics Division of the U.S. Army Research Office.

ISOC is a not-for-profit membership organization founded in 1991 to be the international focal point for global cooperation and coordination in the development of the Internet. Through its current initiatives in support of education and training, Internet standards and protocol, and public policy, ISOC has played a critical role in ensuring that the Internet has developed in a stable and open manner. For 10 years ISOC has run international network training programs for developing countries which have played a vital role in setting up the Internet connections and networks in virtually every country that has connected to the Internet. For more information, please visit: <http://www.isoc.org/>

ISOC to Run .org?

Recently ICANN posted a preliminary Staff Report on the selection of a new registry operator to assume responsibility on January 1, 2003 for the .org registry. The report, which is subject to public comment and comment by all the bidders before being submitted for approval to the ICANN Board of Directors, recommends that the Board select the *Internet Society* (ISOC) as the successor registry operator for the .org registry, currently operated by VeriSign.

This preliminary report follows an extensive bid solicitation and evaluation process that was launched last April. Eleven bids were received in response to a Request for Proposals. These bids were analyzed and evaluated by three evaluation teams that operated independently of each other.

“We received eleven very strong and thoughtful proposals,” noted Stuart Lynn, President of ICANN. “We appreciate the response of the institutions behind these proposals. The ISOC proposal was the only one that received top ranking from all three evaluation teams. On balance, their proposal stood out from the rest.” Lynn also emphasized the openness and transparency of the solicitation and evaluation process.

Two evaluation teams focused on technical issues: one from Gartner, Inc., an international consulting and research organization that specializes in information technologies, and the other a team mainly composed of CIOs of major universities. Another team was provided by ICANN’s *Non Commercial Domain Name Holders* constituency; the NCDNHC team focused on the effectiveness of the proposals to address the particular needs of the .org registry. The staff report integrates these evaluations and other factors into the preliminary recommendation.

ISOC is an international not-for-profit organization of over 6,000 individual and 150 organizational members with chapters in over 100 countries. It provides leadership in addressing issues that confront the future of the Internet, as well as being a home for the *Internet Engineering Task Force* (IETF) and the *Internet Architecture Board* (IAB).

In operating the **.org** registry, ISOC will team with Afilias, an operating registry that recently launched the **.info** top level domain (TLD) that was authorized by ICANN as one of seven new TLDs over this past year.

“Afilias will provide ISOC with the necessary experience at operating a large registry,” said Lynn. “The **.info** registry already houses about 1 million domain names, which is on a scale that approaches the much older **.org** registry.”

ICANN is re-assigning the **.org** registry under a revised agreement among ICANN, VeriSign, and the U.S. Department of Commerce that was signed in May 2001. Under that agreement, VeriSign was permitted to keep its registrar business, NSI (that it was obligated to sell under the prior agreements) provided that it agreed to relinquish **.org** at the end of December 2002, and subject to other provisions of the revised agreements. As part of those revised agreements, VeriSign agreed to endow the new operator with US\$ 5 million to help fund operating costs, provided that the new operator was a not-for-profit organization.

Following an open and transparent process, ICANN has posted all eleven applications online together with all supplemental material and community comments received. The preliminary staff report and the evaluations are posted at:

<http://www.icann.org/tlds/org/preliminary-evaluation-report-19aug02.htm>.

Applicants and any member of the community are invited to send comments on the preliminary report and evaluations by e-mail to:

org-eval@icann.org

Upcoming Events

The *IETF* will meet in Atlanta, Georgia, USA, November 17–21, 2002.
<http://www.ietf.org/meetings/meetings.html>

ICANN will meet in Shanghai, China, October 27–31, 2002.
<http://www.icann.org/meetings/>

The next *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) will take place February 19–28, 2003 in Taipei, Taiwan. <http://apricot2003.net/>

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Technology
WorldCom, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 2002 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSR STD
U.S. Postage
PAID
Cisco Systems, Inc.

The Internet Protocol Journal

December 2002

Volume 5, Number 4

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Internet Multicast Tomorrow	2
Zero Configuration Networks	20
Book Reviews	27
Letters to the Editor	33
Fragments	35

FROM THE EDITOR

In December 1999 we published Part One of a two-part article on Internet Multicast. Some readers have asked “what happened to Part Two?” Finally, in this issue we are able to bring you the second article, “Internet Multicast Tomorrow.” Multicast remains a technology with limited Internet-wide deployment, but numerous research activities are underway that may change this situation. Ian Brown, Jon Crowcroft, Mark Handley and Brad Cain provide an overview of current developments in multicast.

If all computer networking was a simple matter of “plug-and-play,” I suppose this journal would not exist. Nevertheless, it is encouraging to see developments that aim to simplify configuration of network devices, particularly those that move around a lot. The Zeroconf working group of the *Internet Engineering Task Force* (IETF) has been developing standards for “configuration-free” networks. Edgar Danielyan explains the details in our second article.

We continue to receive numerous letters in response to our articles. Your feedback is very much appreciated, because it helps us develop material for future issues. Please keep your letters coming to ipj@cisco.com

The long-awaited online subscription system is now ready for deployment and you will be able to try it out in the very near future at www.cisco.com/ipj. With this system, you can update your mailing address as well as select delivery options, online notification of new issues and so on. As with any computer based system, I anticipate that we, with your help, will uncover a few bugs. Please report any problems you may encounter to ipj@cisco.com.

A new important resource is available from the *Internet Society* (ISOC). *The Internet Report* is a catalogue of IETF documents, including RFCs and Internet Drafts, that document the technology, protocols and operating procedures that form the Internet. The report includes RFCs, IETF Working Group drafts as well as individual drafts. The Internet Report is maintained by Geoff Huston. You can access the report online at <http://ietfreport.isoc.org/>

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Internet Multicast Tomorrow

by Ian Brown, *University College London*,
Jon Crowcroft, *University of Cambridge*,
Mark Handley, *ICIR*,
Brad Cain, *Storigen Systems*

This article is part of a pair, the first of which looked at the state of play in IP multicast routing^[0]. In this article, we look at the broader problems and future activities with multicast. We divide the areas into routing, addressing, transport, security, operations, and research.

There has been quite a bit of debate about the nature of compelling applications for multicast recently.^[44] It is certainly the case that we do not completely understand the “market” for multicast—this is at least in part because multicast does not yet provide a complete set of functions for all the applications and services we might imagine. This is a typical “chicken and egg” situation, though: To put an extreme version of the argument, the application writers do not see any multicast deployed; the *Internet Service Providers* (ISPs) do not see any multicast applications; and the router vendors do not see any multicast service demand from ISPs. (The same problem afflicts IPv6, Integrated and possibly Differentiated Services, and mobile IP, of course.)

As we discussed in the part I of this article^[0], this situation has been somewhat alleviated by streaming applications for audio and video from the classical content providers in the entertainment and news industries. And although we are still seeing some problems, we are also seeing broader interest and development.

The next section presents recent work on routing and addressing. After that we look at transport. Subsequently, we discuss security. Then we look at operations and management. Finally, we examine some of the research ideas that are available.

Routing and Addressing

The single biggest step recently in multicast routing and addressing has been the recognition that the demand for large-scale multicast is largely for one-to-many or single source. Combined with the ability to select sources at the receiver (as a means to prevent denial-of-service attacks) in the *Internet Group Management Protocol* (IGMP)v3, this has made a significant improvement to ISPs’ willingness to deploy the service^[42].

Source-Specific and Single-Source Multicast

The origins of the idea were thesis work at Stanford by Hugh Holbrook on Express multicast^[43]. This is a specialized multicast architecture for one-to-many multicast groups. In this way, Express is a subset of the current multicast model in that it allows only a single sender to a multicast group. The advantages of Express are that certain aspects of multicast routing and addressing are easier solved by ignoring the many-to-many case. Many feel that the most likely large-scale applications of multicast are one-to-many, a fact that explains why Express is becoming popular as a short-term solution.

Express addresses are *channels* that are 64-bit addresses (that is, source address plus group address). Express sources transmit to a channel and advertise that channel. Receivers learn about these channels through advertisements or through other means (that is, URL) and initiate an Express join. Routers propagate these joins directly toward the source, building a source rooted multicast forwarding tree.

The Express model offers two primary benefits. First, Express simplifies the complexity of multicast routing. Secondly, Express simplifies the assignment of multicast addresses for IPv4. Because Express channels are 64 bits, a source can select any lower 32 bits (any group address) for its channel and not collide with another.

In order to implement Express with IPv4 multicast protocols, a special range of multicast addresses was defined. The 232/8 address has been allocated by the *Internet Assigned Numbers Authority* (IANA) for single-source multicast experimentation. In this range, an address has meaning only when “coupled” with a source address. Another way to explain it is that this address range is reserved for the lower 32-bit Express addresses. With this scheme, Express requires no modification to multicast data packets.

Express can be implemented with two protocols that have already been developed: IGMPv3^[42] and *Protocol Independent Multicast Sparse Mode* (PIM-SM).

IGMPv3 extends IGMP to allow source-specific joins to a multicast address. This capability can be used to carry 64-bit (S,G) joins to a router. When a router receives the IGMPv3 join, it must be able to build the source-specific tree with a multicast routing protocol. PIM-SM, widely deployed in service provider networks, already possesses this capability. The combination of IGMPv3 and PIM-SM allows Express to be implemented without creating more protocols; this is one of the most powerful benefits of the Express model.

Interdomain Multicast

Currently there are four fairly widely deployed multicast routing protocols: *PIM Dense Mode* (PIM-DM), PIM-SM or *Source-Specific Multicast* (SSM), *Multicast OSPF* (MOSPF), and the *Distance Vector Multicast Routing Protocol* (DVMRP). Because of the different properties of these protocols, there are many difficulties in connecting heterogeneous routing domains together^[38]. In general, most problems arise when connecting explicit join type protocols with flood-and-prune protocols. With service providers rolling out multicast using PIM-SM, connecting DVMRP and PIM-DM flood-and-prune is becoming common.

In order to connect two multicast routing domains, a *Multicast Border Router* (MBR) needs to exist between the two domains. This router must implement a shared forwarding cache architecture^[39]. In this model, each multicast routing protocol running on a MBR submits its forwarding cache entries to a shared cache. This cache is the “bridge” between the trees in the different domains.

In order that the appropriate trees are created in each domain (on either side of a MBR), signaling must exist to bring sources from one domain to receivers in the other domain. This is part of the complication in connecting flood-and-prune protocol domains to explicit join protocol domains. In an explicit join protocol such as PIM-SM, joins are sent by edge routers to either a source or a *Rendezvous Point* when a host joins. A flood-and-prune protocol works quite differently, in a sense assuming that packets are desired; trees are pruned when edge routers receive new source packet but have no local listeners.

The signaling aspect of joining two domains can be accomplished with a variety of means. There are many options, but two stand out as providing the best methods of connecting domains. The first is to use *Domain Wide Reports* (DWRs)^[36] in flood-and-prune domains. DWRs are similar to IGMP reports except that they are sent on a domain-wide basis. When a border router receives a DWR report, it can join a group on behalf of an entire domain. The second solution is to use the *Multicast Source Discovery Protocol* (MSDP)^[37]. MSDP is currently used to send source lists between PIM-SM domains. It can also be used to connect domains by having the MBR also participate in MSDP. Sources can then be learned from an explicit join protocol domain; the MBR can then join the sources and flood them into attached flood-and-prune protocols domains.

Address Allocation

The schemes to provide dynamic distributed address allocation have not been successful to date. But with many multicast services being limited to either a single domain or a single source, the pressure is off. Instead, source-specific addresses are unique in any case. For many-to-many multicast (sometimes known as *Internet Standard Multicast* [ISM]), the problem has also been alleviated by the use of GLOP^[61], which allocates sections of the address space by mapping Autonomous System numbers of a provider into Class D prefixes. This is potentially inefficient, but solves the contention, collision, revocation, or resolution problem that *Multicast Address Set Claim* (MASC) and *Multicast Address Allocation* (MALLOC)^[60] attempt to do in a distributed dynamic manner.

In the longer term this address allocation, as well as scalable solutions to many-to-many multicast in the local domain and interdomain, await further development on bidirectional trees [“Bi-dir PIM” and the *Border Gateway Multicast Protocol* (BGMP)], which we discuss next. It is likely that these will need IPv6 to scale to serious usage.

Bidirectional PIM-SM

The PIM-SM multicast routing protocol builds both source and shared trees for the distribution of multicast packets. PIM-SM shared trees are rooted at special routers called *Rendezvous Points* and are unidirectional in nature. Shared tree traffic always flows from the *Rendezvous Point* down to the leaf routers. In some types of multicast applications, namely many-to-many type applications, a unidirectional tree may be inefficient.

Other multicast protocols such as *Core Based Trees* (CBT) and BGMP provide bidirectional shared trees. Bidirectional trees^[40] do not have these inefficiencies in many-to-many applications. In a bidirectional tree, traffic from a source is forwarded directly onto the shared tree at the closest point; the traffic is then forwarded both “up” and “down” the tree to all receivers. This is in contrast to a unidirectional tree when the source packets are sent first to the Rendezvous Point (or root) and then down the tree. Recently, two proposals have been submitted that add bidirectional tree capabilities to PIM-SM^[40].

BGMP

BGMP^[33] is a new inter-domain multicast routing protocol that addresses many of the scaling problems of earlier protocols. BGMP attempts to bring together many of the ideas of previous protocols and adds features that make it more service provider friendly. BGMP is designed to be a unified inter-domain multicast protocol in much the same way that the *Border Gateway Protocol* (BGP) is used for unicast routing.

BGMP is an inter-domain protocol in that it adopts particular design features of BGP familiar to providers. Two of these features follow: it uses TCP connections for the transfer of routing information and it has a state machine (with error notifications) similar to BGP.

In order to accommodate different applications and backward compatibility, BGMP can build three types of multicast trees, both unidirectional source and shared trees and bidirectional shared trees. Unidirectional trees are useful for single-source applications and for backward compatibility with other multicast routing protocols. Shared trees are useful for many-to-many applications (for example, multi-player gaming, videoconferencing) and multicast forwarding state to scale for these types of applications.

One of the unique properties of BGMP is that its shared trees are rooted at an Autonomous System that is associated with the multicast group address of the tree. Having the root of the tree at the Autonomous System that is associated with the address is logical because there are likely members in that domain. Rooting the trees at an Autonomous System level also provides stability and inherent fault tolerance.

BGMP requires a way to discover which Autonomous Systems “own” which multicast addresses; this can be accomplished through the use of the MASC protocol or through globally assignable multicast addresses (for example, IPv6 multicast). The MASC protocol allocates temporary assignments from the IPv4 group D address space; it then distributes these assignments into *Multiprotocol BGP* (MBGP) so that BGMP will know which Autonomous System is associated with which group and, therefore, where to send join messages.

If globally assignable addresses are available, then BGMP can use any static address architecture for obtaining an Autonomous System from a multicast group address.

The combination of BGMP and a large multicast address space (for example, IPv6 address space) provide the best scaling for all types of multicast applications.

Transport and Congestion Control: Calling Down Traffic on a Site

Multicast is a multiplier. It gives an advantage to senders, but without their knowledge. Multicast (and its application level cousin, the CU-SeeMe reflector) can “attract” more traffic to a site than it can cope with on its Internet access link. (CU-SeeMe is a popular Macintosh- and PC-based Internet videoconferencing package that currently does not directly use IP multicast.) A user can do this by inadvertently joining a group for which there is a high-bandwidth sender, and then “going for a cup of tea.” This problem will be averted through access control, or through mechanisms such as charging^[58], which may result from the deployment of real-time traffic support.

The problem is seen as critical by ISPs who have a shared bottleneck in their access technology—this is the case for cable modem and in some cases for *Asymmetric Digital Subscriber Line* (ADSL), where a large number of fast lines converge on a slower interface to the backbone. Here, a single user may attract more traffic than this link can handle, without seeing a problem that he or she causes for other users (unicast or other multicast lower-capacity separate sessions using the same shared bottleneck). The use of IGMPv3 with authenticated join and configuration management would appear to be a possible solution to these woes. Alternatively, the use of TCP-friendly multicast congestion control (as envisaged for reliable multicast, but also as emerging in some *Real-Time Transport Protocol* (RTP)^[4] applications), would also solve this problem.

Congestion Control

One of the critical areas to clarify is the role of congestion control in multicast transport protocols^[1]. From an early stage, it was established that coexistence with TCP was a critical design goal for protocols that would operate in the wider Internet. Thus systems such as *TCP Friendly (Reliable) Multicast Congestion Control* (TFMCC)^[8], *Pragmatic General Multicast Congestion Control* (PGMCC)^[53], and receiver-driven congestion control^[54] all extend the classic work by Raj Jain^[15] and Van Jacobson^[17] and subsequent evolution^[16] on TCP congestion avoidance and control.

Recently, this line of thinking has even been extended back into the unicast world in the application of such control schemes to *User Datagram Protocol* (UDP)-like flows in the work on the *Datagram Congestion Control Protocol* (DCCP)^[62], suitable for adaptive multimedia flows on RTP, for example.

Reliable Multicast

There is a clear requirement for some sort of analog to TCP for multicast applications that need a level of reliability. The *Internet Research Task Force's* (IRTF's) *Reliable Multicast Research Group* (RMRG) group^[3] has developed numerous prototypical solutions to the problem, which turns out to be quite a large design space (not “one size fits all”).

The IETF *Reliable Multicast Transport* (RMT) working group has now been chartered to develop single-source reliable multicast transport solutions that meet the current Internet constraints^[1]. That group has developed a building block approach^[12], which is based partly on abstracting components from existing work such as *Reliable Multicast Transport Protocol* (RMTP) II^[18], *Receiver Driven Layered Congestion Control* (RLC)^[7], *Multicast File Transfer Protocol* (MFTP)^[28], *Pragmatic General Multicast* (PGM)^[41], and many other protocols.

Some applications of RMT products are likely to be infrastructural rather than of direct use to the ISPs' customers—for example, distributing software to mirror sites seems to be one popular compelling use.

However, reliable multicast is sometimes regarded as something of an oxymoron. When people talk about “Reliable Multicast,” they usually mean a single protocol at a single “layer” of a protocol stack, typically the transport layer (although we have seen people propose it in the network and even link [ATM!] layers too), that can act as any layered protocol can—to provide common functionality for applications (higher layers) that need it.

So what is wrong with that? Well, possibly three things (or more):

- *Fate sharing*: Fate sharing in unicast applications means that as long as there is a path that IP can find between two applications, then TCP can hang on to the connection as long as the parties like. However, if either party fails, the connection certainly fails.
Fate sharing between multicast end points is a more subtle idea. Should “reliability” extend to supporting the connection fork recipients failing? Clearly this will be application specific (just as timing out on not getting liveness out of a unicast connection is for TCP—we must permit per-recipient timeouts and failures).
- *Performance*: When A talks to B, the performance is limited by one path. Whatever can be done to improve the throughput (or delay bound) is done by IP (for example, load sharing the traffic over multiple paths). When A talks to B, C, D, E, or F, should the throughput or delay be that sustainable by the slowest or average?
- *Semantics*: As well as performance and failure modes, N-way reliable protocols can have different service models. We could support reliable one-to-n, reliable n-to-one, and reliable n-to-m.

Applications such as software distribution are cited as classic one-to-n requirements. Telemetry is given as an n-to-one reliable protocol. Shared whiteboards are cited as examples of n-to-m applications.

It is interesting to look at the reliability functions needed in these. The one-to-n and n-to-one protocols are effectively *simplex* bulk transfer applications. In other words, the service is one where reliability can be dealt with by “rounding up” the missing bits at the end of the transfer. Because this does not need to be especially timely, there is no need for this to be other than end to end, and application based. (Yes, we know telemetry could be time sensitive, but we are trying to illustrate major differences clearly for now.)

On the other hand, n-to-m processes such as whiteboards need timely recovery from outages. The implication is that the “service” is best done somewhat like the effect of having $n \times (m-1)/2$ TCP connections. If used in the WAN, the recovery may best be distributed, because requests for recovery will implode down the very links that are congested or error prone and cause the need for recovery.

Now there are different schemes for creating distributed recovery. If the application semantics are that operations (application data unit packets worth) are sequenced in a way that the application can index them, then any member of a multicast session can efficiently help any other member to recover (examples of this include Mark Handley’s Network Text tool^[16].) On the other hand, packet-based recovery can be done from data within the queues between network or transport and application, if they are kept at all members in much the same way as a sender in a unicast connection keeps a copy of all unacknowledged data.

The problem with this is that *because* it is multicast, we do not have a positive acknowledgement system. Therefore, there is no way to inform *all* end points when they can safely discard the data in the “retransmit” queue. Only the application really knows this!

Well, this is not to say that there is not an obvious toolkit for reliable multicast support—it would certainly be good to have RTP-style media timestamps (determined by the application, but filled in by the system). It would be good to have easy access to a timestamp-based receive queue so applications could use this to do all functions discussed previously. It might be advantageous to have virtual Token Ring, expanding ring search, token tree, and other toolkits to support retransmit “helper” selection.

Table 1 illustrates this in terms of where functions might be put to provide reliability (retransmit), sequencing, and performance (adaptive playout, say, versus end to end, versus hop-by-hop delay constraint).

Table 1: Reliable Multicast Semantics

	Recovery	Sequency	Dalliance
<i>Network</i>	not in our internet	ditto	int-serv
<i>Transport</i>	one-many	yes	adaptive
<i>Application</i>	many-many	operation semantics	adaptive

Router Assist for Reliable Multicast

As mentioned in previous sections, one of the difficulties in end-to-end multicast signaling is the “implosion” of signaling at a source from many receivers. This problem has been addressed in numerous ways, including the use of timers, the use of servers to aggregate signaling, and the use of router-assisted mechanisms. We now discuss three protocols that make use of router assistance in order to better scale end-to-end multicast protocols.

PGM^[41] is a *negative acknowledgement* (NAK)-based router-assisted reliable multicast protocol. PGM uses routers to aggregate receiver-to-source signals (for example, the NAKs) as they flow toward the source. PGM router support also includes a subcasting ability whereby repairs will flow down only to receivers who have requested them.

Extending the ideas of router assist in PGM is the *Generic Multicast Transport Service* (GMTS). GMTS provides generic, fixed, simple services for any end-to-end multicast transport protocol. These services include such features as signal aggregation with predicates and sophisticated subcasting ability. GMTS was used as a basis for *Generic Router Assist* (GRA)^[34], which is similar, IETF standards oriented, and a bit more streamlined.

Securing Multicast

Multicast security is more difficult than unicast security in several areas. The key exchange protocols used between unicast hosts do not scale to groups. Rekeying is required more often to maintain confidentiality as group membership changes. And the efficient authentication transforms used between two unicast hosts cannot protect traffic between mutually distrustful members of a group.

These problems are being worked on by the IETF *Multicast Security* (msec) and IRTF *Group Security* (gsec) working groups. Because of the wide range of application requirements in group communication, their work is based upon a building block approach similar to that of the RMT group.

The blocks being developed are data security transforms, group key management and group security association, and group policy management^[49]. An application may use different blocks together to create a protocol that meets its specific requirements.

Data Security Transforms

A data security transforms block provides confidentiality and authentication services for data being transported between group members. Confidentiality is reasonably easy to provide using standard encryption algorithms. Authentication is more difficult, because the algorithms used in unicast protocols such as *IP Security* (IPSec) would not allow a group member to authenticate data as being from another specific group member. This is because the secret used to authenticate the traffic must be shared between all sending and receiving parties. Public-key signatures would solve this problem, but are an order of magnitude slower than symmetric authentication algorithms and hence especially unsuitable for real-time traffic and low-powered communications devices.

Instead, blocks such as the *Timed Efficient Stream Loss-tolerant Authentication Protocol* (TESLA)^[55] are being developed that trade off small amounts of functionality (such as immediate rather than slightly delayed authentication) to retain the efficiency benefits of symmetric algorithms. TESLA senders use a hash chain of keys $k_{n,\dots,1}$ to sign data, where: $k_n = \text{hash}(k_{n-1})$

They release each key in the chain a short interval after the data the key has signed. As long as other group members received the data during that interval, they can be confident that the signature was made by the sender. If keys are lost during transmission, receivers can recompute any key earlier in the sequence simply by repeatedly applying the hash function used to any later key received. Finally, they can be sure that keys are coming from the sender because the first key in the sequence is digitally signed, while only the sender can know the later keys in the sequence (because by definition, a hash function must not be reversible).

Group Key Management and Group Security Association

To use data security transforms, group members need to possess the cryptographic keys necessary to encrypt or decrypt and sign or authenticate data. They also need to agree on parameters such as specific encryption algorithms. This building block allows this information to be shared between group members.

The Group Key Management architecture^[47] provides a unified model for key management blocks. A central *Group Controller/Key Server* (GCKS) provides *Traffic Encrypting Keys* (TEKs) or *Key Encrypting Keys* (KEKs) to new group members after authenticating them with a unicast protocol. The GCKS may also delegate some of its functions to other entities, improving scalability.

In groups with simple security requirements, this may be the only communication required between a group member and GCKS. But if group changes need to be cryptographically enforced, further TEKS, encrypted using a KEK, may be provided to members by multicast or a more scalable protocol such as the *Logical Hierarchy of Keys* (LHK)^[56] that does not require every rekey message to be sent to every group member. Alternatively, noninteractive mechanisms such as hash trees may be used to update keys^[48]. Finally, group members may explicitly de-register with the GCKS using a one- or two-step message.

Three key management building blocks are being developed. The *Group Domain of Interpretation* (GDOI) builds on the *Internet Security Association Key Management Protocol* (ISAKMP)^[52] to allow the creation and management of security associations for IPSec and other network or application layer protocols^[46]. *Multimedia Internet Keying* (MIKEY) is targeted at real-time multimedia communications, particularly those using the Secure RTP, and can be tunneled over the *Session Initiation Protocol* (SIP)^[45]. And a *Group Secure Association Key Management Protocol* (GSAKMP), along with a GSAKMP-Light profile, have also been developed^[51].

Group Policy Management

The final building block defines policies such as which roles various entities may play in the group; who may hold group information such as cryptographic keys; the cryptographic algorithms used to protect group data; and proof that the creator of a given policy is authorized to do so. A group policy token is used to hold all of this information^[50]. All or part of tokens can be made available to users in policy repositories or by using other out-of-band mechanisms.

Operational Deployment of Multicast

As mentioned previously, multicast seems to be difficult to deploy. One problem is that it has only recently moved from the research community (and typically implemented using tunnels) into the service community (running native IP multicast routing).

This means that debugging multicast sessions, applications, and routing is a common activity. However, because of the dynamic nature of multicast addresses and the anonymous nature of the multicast service model, debugging is somewhat more difficult than for the equivalent unicast case.

Fortunately, all current native multicast paths are at least computed from underlying unicast ones, and it is possible to use tools such as *mtrace* and *mrm* to query the underlying router system to try to figure out where things are going on. Of course, the relevant *Management Information Bases* (MIBs) need to be designed, but mere *Simple Network Management Protocol* (SNMP) access to the variables defined in these may not be enough.

Many multicast sessions are global, and not surprisingly, someone, somewhere, sometime in the session will have a problem. In a way, you only have to look at multicast as a way of sampling large pieces of the Internet at one time to see why it is difficult to understand. In fact, a research project called *Multicast-Based Inference of Network-Internal Characteristics* (MINC)^[9, 57] is using that very observation to build tools of more general use.

MRM

One recent tool that has been developed to facilitate multicast monitoring and debugging is the *Multicast Reachability Monitor* (MRM)^[32]. MRM consists of two parts; a MRM management station configures test senders and test receivers in multicast networks. A multicast test sender or test receiver is any server or router that supports the MRM protocol and can source or sink multicast traffic. MRM provides the ability to dynamically test particular multicast scenarios; this capability can be used for fault isolation and general monitoring of sessions.

MRM is typically used to configure MRM-capable routers as test senders and test receivers from a management station. Routers configured as test senders send multicast packets periodically to a configured multicast group at a configured rate. Routers configured as test receivers monitor traffic to a group and keep statistics that can be reported back via *RTP Control Protocol* (RTCP) packets. Test receivers can be configured to send RTCP reports when a given condition has been reached or when polled by a management station. Although the MRM protocol is simple itself, it provides powerful capabilities that can be used by future multicast debugging applications.

Research Ideas in Multicast Routing and Addressing

The seeming complexity exhibited by the full panoply of multicast protocols has led some people to develop doubts as to the eventual deployment of multicast. It is far too early to say whether these doubts are well founded. The slow pace of deployment is a symptom not just of this complexity, but also of the underlying complexity of handling growth and evolution of *any* type in such a large system as the Global Internet.

Having said that, it is worth mentioning four of the approaches that have been discussed in the Internet community recently:

- *Addressable Internet Multicast* (AIM), by Brian Levine, et al., attempts to provide explicit addressing of the multicast tree. The routers run a tree-walking algorithm to label all the branch points uniquely, and then make these labels available to end systems. This allows numerous interesting services or refinement of multicast services to be built. Of some particular interest would be the ability this service gives to end systems to do subcasting, which would be useful for some classes of reliable transport protocols.

- *Explicitly Requested Single-Source* (Express), by Hugh Holbrook et al., is aimed at optimizing multicast for a single source. The proposal includes additional features such as authentication and counting of receivers, which could be added to many other multicast protocols usefully. It is motivated by a perceived requirement from some ISPs for these additional features. Express makes use of an extended address (channel + group) to provide routing without global agreement on address assignment. A possible source of problem for AIM is the potential for unbounded growth in the size of identifiers for labeling subtree branch points.
- *Root Addressed Multicast Architecture* (RAMA), by Radia Perlman et al., is in some senses a generalization of Express type addressing, but it also requires bidirectional trees (CBT like, rather than current PIM-SM, although work on bidirectional PIM is under way too). The goal is to offer a single routing protocol for both intra- and interdomain. In fact, RAMA can be implemented by combining the address extensions proposed for Express, and two-level bidirectional PIM as an implementation of BGMP. RAMA and Express (and bidirectional PIM) require a mechanism for carrying additional information in multicast IP data packets.

There are two critical problems for carrying this identifier that are difficult to solve in general: first, it takes new space in the IP packet, and this has to be accessed by both hosts and routers—that represents a deployment problem; secondly, in the general case, the extra field must be examined on the “fast path,” in routers that have such a concept, and this takes valuable processing resources that may have to be taken away from some other forwarding task.

- *Connectionless Multicast* (CM) by Dirk Ooms, et al., is a proposal for small, very sparse groups to be implemented by carrying lists of IP unicast addresses in packets. The scheme is not simply a form of loose source routing, because it would make use of packet replication at appropriate branch points in the network. It may be well suited to IP telephony applications where a user starts with a unicast call, but then adds a third or fourth participant.
- The *L'Ecole Polytechnique Fédérale de Lausanne* (EPFL) work on *Distributed Core Multicast* (DCM) aims to address very large numbers of very small groups with mobile users, typical characteristics of mobile IP telephony users making conference or group calls.
- MIT has done some work on the use of wide-area “anycast” addresses for the core and Rendezvous Point. This results in a potential improvement in the availability of trees (and subtrees) for multicast delivery in the event of router or link outage. More importantly, it may be possible for a multicast group to survive network partitions (or lack of core reachability), a possibility that would make this an invaluable improvement to the service. It depends on the scalability of the wide-area anycast solution, which the MIT work shows is at least viable, and certainly worth more attention.

- *Yet Another Multicast* (YAM) routing protocol^[30] was devised by Ken Carlberg of SAIC to address the possibility of forming different multicast trees based on some QoS metric—the idea is that IGMP is modified to provide a “one-to-many” join, and a receiver sends this with required performance parameters. Routers receiving the request over links that can provide this service respond. The receiver (sender of the one-to-many IGMP) selects the one to then commit the join to.
- *Quality of Service Sensitive Multicast Internet protoCol* (QoSMIC) is a development from YAM by Faloutsos^[29] at Toronto, and slightly modifies the tree-building exercise.
- When multicast and *Multiprotocol Label Switching* (MPLS) are mentioned together, there is both confusion and surprise. MPLS can be used with multicast in two very different ways. The first method is by building multicast trees over MPLS traffic-engineered paths. Some multicast routing protocols already make use of unicast forwarding information for the construction of multicast trees. Using multicast traffic-engineered paths is simply an extension of this concept—with one caveat. Some multicast routing protocols use *Reverse Path Forwarding* (RPF) checks on incoming packets to prevent looping; this is accomplished by checking to see if the incoming interface is the “closest” to the source. With MPLS traffic engineering, RPF checks are difficult. A solution has not been presented at this time that addresses this problem.

The second method for using multicast with MPLS is through the use of point-to-multipoint virtual circuits in much the same way as ATM point-to-multipoint virtual circuits. These are useful in cases where receivers are statically configured to a multicast address or multicast traffic is always to be delivered to a destination. Mapping dynamic memberships into a multipoint circuit has proven difficult, for example, with ATM. There are currently several Internet drafts that propose various solutions for MPLS and multicast^[31].

- Several groups have been working on end system-only multicast schemes, probably most notably Carnegie-Mellon University^[59].

Summary and Conclusions

In this article, we have looked at some of the newer ideas in the research and development community in the area of multicast. There is still a lot to be done to close the loop between network services, transport, and applications, but present research indicates that we will eventually achieve this goal.

References

- [0] M. Handley and J. Crowcroft, "Internet Multicast Today," *The Internet Protocol Journal*, Vol. 2, No. 4, December 1999.
- [1] A. Mankin, A. Romanow, S. Bradner, and V. Paxson, "IETF Criteria for Evaluating Reliable Multicast Transport and Application Protocols," RFC 2357, June 1998.
- [2] J. W. Byers, M. Luby, M. Mitzenmacher, and A. Rege, "A Digital Fountain Approach to Reliable Distribution of Bulk Data," Proceedings of SIGCOMM '98, September 1998.
- [3] Reliable Multicast Research Group:
<http://www.east.isi.edu/RMRG/>
- [4] H. Schulzrinne, S. Casner, R. Frederick, and V. Jacobson, "RTP: A Transport Protocol for Real-Time Applications," RFC 1889, January 1996.
- [5] S. Floyd, V. Jacobson, C. Liu, S. McCanne, and L. Zhang, "A Reliable Multicast Framework for Light-weight Sessions and Application Level Framing, Scalable Reliable Multicast (SRM)," Proceedings of ACM SIGCOMM '95.
- [6] M. Handley and J. Crowcroft, "Network Text Editor (NTE): A scalable shared text editor for the Mbone," Proceedings of ACM SIGCOMM '97, September 1997.
- [7] L. Vicisano, L. Rizzo, and J. Crowcroft, "TCP-like Congestion Control for Layered Multicast Data Transfer," Proceedings of INFOCOM '98.
- [8] M. Handley, S. Floyd, and B. Whetten, "Strawman specification for TCP friendly (reliable) multicast congestion control (TFMCC)," work in progress.
- [9] S. R. Caceres, N. Duffield, J. Horowitz, D. Towsley, and T. Bu, "Multicast-Based Inference of Network-Internal Characteristics: Accuracy of Packet Loss Estimation," Proceedings of IEEE Infocom '99, March 1999.
- [10] S. J. Cowley, "Of Timing, Turn-taking, and Conversations," *Journal of Psycholinguistic Research*, 1998, Vol. 27, No. 5, pp. 541–571.
- [11] Jonathan Rosenberg and Henning Schulzrinne, "Timer Reconsideration for Enhanced RTP Scalability," Proceedings of the Conference on Computer Communications (IEEE Infocom), March/April 1998.
- [12] B. Whetten, L. Vicisano, R. Kermode, M. Handley, S. Floyd, and M. Luby, "Reliable Multicast Transport Building Blocks for One-to-Many Bulk-Data Transfer," RFC 3048, January 2001.
- [13] Handley, M. et al., "Rate Adjustment Protocol," Proceedings of Infocom 1999.
- [14] Kouvelas, I. et al., "Self Organising Transcoders," Proceedings of NOSSDAV 1998.

- [15] D-M. Chiu and R. Jain, "Analysis of the Increase and Decrease Algorithms for Congestion Avoidance," *Computer Networks and ISDN Systems*, Vol. 17, pp. 1–14, 1989.
- [16] S. Floyd and K. Fall, "Router Mechanisms to Support End-to-End Congestion Control," Technical report, <ftp://ftp.ee.lbl.gov/papers/collapse.ps>
- [17] V. Jacobson, "Congestion Avoidance and Control," Proceedings of ACM SIGCOMM '88, August 1988, pp. 314–329.
- [18] J. C. Lin and S. Paul, "RMTP: A Reliable Multicast Transport Protocol," Proceedings of IEEE INFOCOM '96, March 1996, pp. 1414–1424.
- [19] M. Mathis, J. Semke, J. Mahdavi, and T. Ott, "The Macroscopic Behaviour of the TCP Congestion Avoidance Algorithm," *ACM Computer Communication Review*, Vol. 27 No. 3, July 1997.
- [20] S. McCanne, V. Jacobson, and M. Vetterli, "Receiver-driven Layered Multicast," Proceedings of SIGCOMM '96, August 1996, pp. 1–14.
- [21] J. Padhye, V. Firoiu, D. Towsley, and J. Kurose, "Modelling TCP Throughput: A Simple Model and Its Empirical Validation," Proceedings of SIGCOMM '98, September 1998.
- [22] L. Rizzo and L. Vicisano, "A Reliable Multicast Data Distribution Protocol Based on Software FEC Techniques," The Fourth IEEE Workshop on the Architecture and Implementation of High Performance Communication Systems (HPCS '97), June 1997.
- [23] Dan Rubenstein, Jim Kurose, and Don Towsley, "The Impact of Multicast Layering on Network Fairness," Proceedings of ACM SIGCOMM '99, August 1999.
- [24] N. Shacham, "Multipoint Communication by Hierarchically Encoded Data," Proceedings of IEEE Infocom '92, 1992, pp. 2107–2114.
- [25] Chris Greenhalgh, Steve Benford, Adrian Bullock, Nico Kuijpers, and Kurt Donkers, "Predicting Network Traffic for Collaborative Virtual Environments," *Computer Networks and ISDN Systems*, Vol. 30, 1998, pp. 1677–1685.
- [26] Steve Deering, "Host Extensions for IP Multicasting," RFC 1112, August 1989.
- [27] S. Deering, C. Partridge, and D. Waitzman, "Distance Vector Multicast Routing Protocol," RFC 1075, November 1988.
- [28] Ken Miller, "Multicast File Transfer Protocol," White Paper, Starburst Technologies.
- [29] Michalis Faloutsos, Anindo Banerjee, and Rajesh Pankaj, "QoS-MIC: Quality of Service Sensitive Multicast Internet Protocol," *ACM Computer Communication Review*, Vol. 28, pp. 144–153, September 1998.
- [30] K. Carlberg and J. Crowcroft, "Building Shared Trees Using a One-To-Many Joining Mechanism," *ACM Computer Communication Review*, Vol. 27, pp. 5–11, January 1997.

- [31] D. Ooms, B. Sales, W. Livens, A. Acharya, F. Griffoul, and F. Ansari, "Framework for IP Multicast in MPLS," work in progress.
- [32] K. Almeroth, K. Sarac, and L. Wei, "Supporting Multicast Management Using the Multicast Reachability Monitor (MRM) Protocol," UCSB CS Technical Report, May 2000.
- [33] D. Thaler, D. Estrin, D. Meyer, et al., "Border Gateway Multicast Protocol (BGMP)," Proceedings of ACM SIGCOMM '98, 1998.
- [34] B. Cain, T. Speakman, and D. Towsley, "Generic Router Assist Building Block," work in progress.
- [35] B. Cain and D. Towsley, "Generic Multicast Transport Services (GMTS)," Proceedings of Networking 2000, Paris, France, May 2000.
- [36] B. Fenner, "Domain Wide Multicast Group Membership Reports," work in progress.
- [37] D. Farinacci et al., "Multicast Source Discovery Protocol," Internet Draft, January 2000, work in progress.
- [38] B. Cain, "Connecting Multicast Domains," Internet Draft, work in progress, October 1999.
- [39] D. Thaler, "Interoperability Rules for Multicast Routing Protocols," RFC 2715, October 1999.
- [40] D. Estrin and D. Farinacci, "Bi-directional Shared Trees in PIM-SM," work in progress.
- [41] T. Speakman et al., "PGM Reliable Transport Protocol Specification," RFC 3208, December 2001.
- [42] B. Cain, S. Deering, and A. Thyagarajan, "Internet Group Key Management Protocol, Version 3," work in progress.
- [43] H. Holbrook and D. Cheriton, "IP Multicast Channels: Express Support for Large-scale Single-source Applications," Proceedings of SIGCOMM '99, September 1999.
- [44] C. Diot, B. Levine, B. Lyles, H. Kassem, and D. Balensiefen, "Deployment Issues for the IP Multicast Service and Architecture," IEEE Network Magazine, Special Issue on Multicasting, January/February 2000.
- [45] J. Arkko, E. Carrera, F. Lindholm, M. Naslund, and K. Norrman, "MIKEY: Multimedia Internet KEYing," Internet Draft, work in progress, February 2002.
- [46] M. Baugher, T. Hardjano, H. Harney, and B. Weis, "The Group Domain of Interpretation," Internet Draft, work in progress, February 2002.
- [47] M. Baugher, R. Canetti, L. Dondeti, and F. Lindholm, "Group Key Management Architecture," Internet Draft, work in progress, February 2002.

- [48] B. Briscoe, "MARKS: Zero Side Effect Multicast Key Management Using Arbitrarily Revealed Key Sequences," Proceedings of Networked Group Communication, November 1999.
- [49] T. Hardjano, R. Canetti, M. Baugher, and P. Dinsmore, "Secure IP Multicast: Problem Areas, Framework, and Building Blocks," Internet Draft, work in progress, September 2000.
- [50] T. Hardjano, H. Harney, P. McDaniel, A. Colgrove, and P. Dilmore, "Group Security Policy Token," Internet Draft, work in progress, November 2001.
- [51] H. Harney, A. Schuett, and A. Colegrove, "GSAKMP Light," Internet Draft, work in progress, July 2001.
- [52] D. Maughan, M. Schertler, M. Schneider, and J. Turner, "Internet Security Association and Key Management Protocol (ISAKMP)," RFC 2408, November 1998.
- [53] Luigi Rizzo, "pgmcc: A TCP-friendly Single-Rate Multicast Congestion Control Scheme," Proceedings of ACM SIGCOMM '2000, August 2000.
- [54] Luby et al., "Wave and Equation Based Rate Control Using Multicast Round Trip Time," Proceedings of ACM SIGCOMM '2002, September 2002.
- [55] A. Perrig, R. Canetti, B. Briscoe, D. Tygar, and D. Song, "TESLA: Multicast Source Authentication Transform," Internet Draft, work in progress, November 2000.
- [56] D. M. Wallner, E. Harder, and R. C. Agee, "Key Management for Multicast: Issues and Architectures," RFC 2627, September 1998.
- [57] F. Lo Presti, N.G. Duffield, J. Horowitz, and D. Towsley, "Multicast-Based Inference of Network-Internal Delay Distributions," <http://www.cs.umass.edu/pub/Lopr99TR9955.ps.Z>
- [58] T. Henderson and S. Bhatti, "Protocol Independent Multicast Pricing," Proceedings of NOSSDAV 2001.
- [59] Yang-hua Chu, Sanjay G. Rao, and Hui Zhang, "A Case for End System Multicast," Proceedings of ACM SIGMETRICS, June 2000, pp. 1–12.
- [60] Multicast Address Allocation Working Group, <http://www.icir.org/malloc/>
- [61] D. Meyer and P. Lothberg, "GLOP Addressing in 233/8," RFC 3180, September 2001.
- [62] <http://www.icir.org/dccp/>

IAN BROWN holds a BSc from The University of Newcastle upon Tyne and a PhD from University College London. His research has focused on network security and active networking. He is a member of the ACM, IEEE, and is a contributor to the Internet Engineering Task Force, particularly in the area of authorized emergency communications. He has also worked extensively on the social implications of technology, and is a trustee of Privacy International and advisory board member of the Foundation for Information Policy Research. His e-mail address is:

I.Brown@cs.ucl.ac.uk

BRAD CAIN is a Senior Consulting Engineer at Storigen Systems, where he contributes to product development in the areas of networking and storage technology. Prior to joining Storigen, Cain was chief scientist at Cereva Networks, where he worked on system architecture and new product development. Cain also worked at Mirror Image Internet, one of the first commercial Content Delivery Networks (CDNs), where he helped architect their content distribution system. Cain is a contributor in the IETF and IRTF in the areas of IP multicast, IP routing, MPLS, and content networking. He has published numerous papers in the areas of routing and multicast and has more than 40 patents pending in the areas of multicast, security, routing, and router architecture. Cain holds a masters and bachelors in electrical engineering from the University of Delaware. E-mail: **Brad.Cain@storigen.com**

JON CROWCROFT is the Marconi Professor of Networked Systems at the University of Cambridge. Prior to that he was professor of networked systems at University College London (UCL) in the Computer Science Department. He is a member of the ACM, a Fellow of the British Computer Society, a Fellow of the IEE, and a Fellow of the Royal Academy of Engineering, as well as a senior member of the IEEE. He is a member of the IAB, and was general chair for ACM SIGCOMM from 1995 to 1999. He is on the editorial team for the ACM/IEEE *Transactions on Networks and Computer Communications*, as well as on the program committee for ACM SIGCOMM and IEEE Infocomm. He has published five books—the latest is *Linux TCP/IP Implementation*, published by Wiley in 2001.

E-mail: **Jon.Crowcroft@cl.cam.ac.uk**

MARK HANDLEY received his BSc in Computer Science with Electronic Engineering from University College London in 1988 and his PhD from UCL in 1997. For his PhD he studied multicast-based multimedia conferencing systems, and was technical director of the European Union funded “MICE” and “MERCI” multimedia conferencing projects. After two years working for the University of Southern California’s Information Sciences Institute (ISI), he moved to Berkeley to join the new ICSI Center for Internet Research (formerly known as ACIRI). Most of his work is in the areas of scalable multimedia conferencing systems, reliable multicast protocols, multicast routing and address allocation, and network simulation and visualization. He is co-chair of the IRTF Reliable Multicast Research Group, and he previously chaired the IETF Multiparty Multimedia Session Control working group. E-mail: **mjh@icir.org**

Zero Configuration Networking

by Edgar Danielyan, Danielyan Consulting

Zero configuration networking may sound like an oxymoron to many who spend most of their time setting up and mending networks. But don't decide on a career change yet—although zero configuration networks exist and work, they don't work always and everywhere. In this article I describe the current state of the affairs in zero configuration IP networking, introduce Zeroconf, the suite of zero configuration IP protocols, and tell what they do and how they work. This article is only a brief introduction to zero configuration networking and Zeroconf, so if you are really interested in all the details, refer to the sources listed in the References section at the end of this article.

The best introduction to Zeroconf is the one from the *Zeroconf Working Group* of the *Internet Engineering Task Force* (IETF)^[1]:

“The goal of the Zero Configuration Networking (Zeroconf) is to enable networking in the absence of configuration and administration. Zero configuration networking is required for environments where administration is impractical or impossible, such as in the home or small office, embedded systems ‘plugged together’ as in an automobile, or to allow impromptu networks as between the devices of strangers on a train.”

Essentially, to reduce network configuration to zero (or near zero) in *Internet Protocol* (IP) networks, it is necessary, inter alia, to:

- Distribute IP addresses (without a *Dynamic Host Configuration Protocol* [DHCP] server),
- Provide name resolution (without a *Domain Name System* [DNS] server),
- Find and list services (without a directory service), and
- Distribute multicast IP addresses, if necessary (without a multicast server).

These and other requirements are defined in an Internet Draft titled “Requirements for Automatic Configuration of IP Hosts” by Aidan Williams^[2]. This document does not define Zeroconf protocols themselves but instead spells out the requirements that should be met to achieve effective and useful zero configuration IP networking. One of the most important requirements for any Zeroconf protocol is that it should not interfere with other protocols and it must be able to exist on the same network with other non-Zeroconf protocols and devices. Another requirement is “no less” security—Zeroconf protocols should not be less secure than existing non-Zeroconf protocols—more on this later. Although IPv6 addresses some of the requirements of zero configuration networking (such as automatic allocation of link-local addresses), other requirements have yet to be met for both IPv4 and IPv6.

Zeroconf IETF Working Group

The Zeroconf Working Group of the IETF is chaired by Erik Guttman of Sun Microsystems and Stuart Cheshire from Apple Computer, with Thomas Narten (IBM) and Erik Nordmark (Sun) serving as area directors. It was chartered in September 1999 and had its first meeting at the 46th IETF in Washington, D.C., in November 1999. Those interested in the work of Zeroconf WG may find the mailing list archive of the working group at:

<http://www.merit.edu/mail.archives/zeroconf/>

Where and When to Use Zeroconf

For a correct understanding of the applicability and usefulness of Zeroconf it is necessary to keep in mind that it is a *link-local* technology. Link-local addressing and naming are meaningful only in a particular network; link-local addresses and names are not global and are not unique globally. In this case it means that Zeroconf is intended for use in small wired or wireless local-area networks in situations and places where zero configuration is necessary. It is appropriate to use Zeroconf in such networks when there is no possibility (or it is inappropriate) to set up a working IP network using the traditional technologies such as DNS and DHCP. Zeroconf is not appropriate and should not be used in many cases, for example in:

- Medium or large networks
- Networks where a high degree of security and control is required
- Large public access networks
- Networks with low bandwidth and high latency (such as some wireless networks)

When inappropriately used, Zeroconf may bring more problems and headaches than it solves. In contrast, examples of correct and appropriate use would include:

- Home and small office networks
- Ad hoc networks at meetings and conferences (especially wireless networks)
- Two devices needing to spontaneously share or exchange information

Likewise, Zeroconf advantages from one viewpoint may become annoying problems from another. Consider, for instance, the automatic distribution and configuration of link-local IP addresses. For a home network user this is a blessing—no longer do you have to spend time creating an addressing scheme and setting the IP addresses and netmasks on devices that should just work. But for an enterprise network (especially an incorrectly configured one), sudden appearance of nodes with (yet) unfamiliar and strange (this is not your regular **10.*** or **192.168.***) IP addresses may result in more than surprise and added workload for the network administrator.

Continuing in this manner, Multicast DNS (mDNS) that ends the misery of having to remember and type `ftp 10.20.30.1` every time you need to transfer files from or to your PC named Bobo and replaces it with just `ftp bobo` may result in strange behavior on some networks. The bottom line? Zeroconf is not a one-size-fits-all solution; it wasn't designed to be one, and will not work as one.

Zeroconf and Security

Security should occupy an important place in the minds of all networking professionals, so an introduction to zero configuration networking would be incomplete without a mention of its security position. Security goals of Zeroconf are defined in section 4, Security Considerations, of “Requirements for Automatic Configuration of IP Hosts”^[2]:

“Zeroconf protocols are intended to operate in a local scope, in networks containing one or more IP subnets, and potentially in parallel with standard configured network protocols. Application protocols running on networks employing zeroconf protocols will be subject to the same sets of security issues identified for standard configured networks. Examples are: denial of service due to the unauthenticated nature of IPv4 ARP and lack of confidentiality unless IPSec-ESP, TLS, or similar is used. However, networks employing zeroconf protocols do have different security characteristics, and the subsequent sections attempt to draw out some of the implications.

Security schemes usually rely on some sort of configuration. Security mechanisms for zeroconf network protocols should be designed in keeping with the spirit of zeroconf, thus making it easy for the user to exchange keys, set policy, etc. It is preferable that a single security mechanism be employed that will allow simple configuration of all the various security parameters that may be required. Generally speaking, security mechanisms in IETF protocols are mandatory to implement. A particular implementation might permit a network administrator to turn off a particular security mechanism operationally. However, implementations should be “secure out of the box” and have a safe default configuration.

Zeroconf protocols MUST NOT be any less secure than related current IETF-Standard protocols. This consideration overrides the goal of allowing systems to obtain configuration automatically. Security threats to be considered include both active attacks (e.g. denial of service) and passive attacks (e.g. eavesdropping). Protocols that require confidentiality and/or integrity should include integrated confidentiality and/or integrity mechanisms or should specify the use of existing standards-track security mechanisms (e.g. TLS (RFC 2246), ESP (RFC 1827), AH (RFC 2402) appropriate to the threat.”

Although this document does not address each and every aspect of security issues with Zeroconf, it sets requirements for Zeroconf protocols. As is the case with traditional IPv4 and IPv6, use of such techniques as *IP Security Architecture* (IPSec) or *Transport Layer Security* (TLS) may be appropriate in some cases. However, the nonstatic (or one may say non-durable) nature of both IP addresses and names in Zeroconf environment may pose a problem for IPSec and TLS deployment.

Dynamic Configuration of IPv4 Link-Local Addresses

Generally speaking, the first requirement that should be fulfilled before any useful IP communication can occur are the IP addresses of sender and recipient. The IP addresses are usually either assigned and set manually or provided by some other means such as DHCP or the *Point-to-Point Protocol* (PPP). However, neither of these is possible in zero configuration networks. Therefore, an automatic mechanism for dynamic configuration of IP addresses without any manual intervention or dependence on third-party service (that is, DHCP) is necessary. This mechanism already exists in IPv6 but not in IPv4. In “Dynamic Configuration of IPv4 Link-Local Addresses”^[3], Stuart Cheshire, Bernard Aboba, and Erik Guttman describe a method that may be used in IPv4 networks to automatically assign IPv4 addresses valid for local communication on a particular interface. A special network **169.254/16** is reserved with the *Internet Assigned Numbers Authority* (IANA) for this purpose. It is necessary to highlight that **169.254/16** addresses are reserved for link-local use only. The document also addresses such issues as support for multiple addresses and multiple interfaces, continuous address conflict detection, effects of joining previously not interconnected networks, and other considerations.

IPv4 Address Conflict Detection

Address conflicts in IP networks are annoying problems that (needlessly) take time and effort to detect and rectify, so a separate document on address conflict detection was deemed necessary. “IPv4 Address Conflict Detection”^[4] by Stuart Cheshire presents two things: first, a way to prevent this unfortunate situation of conflicting IP addresses from happening, and second, a way to detect address conflicts if they do happen even after all the precautions. Both of these are accomplished using the *Address Resolution Protocol* (ARP). Interestingly, in the Security Considerations section of the document the author states:

“The ARP protocol [RFC 826] is insecure. A malicious host may send fraudulent ARP packets on the network, interfering with the correct operation of other hosts. For example, it is easy for a host to answer all ARP requests with responses giving its own hardware address, thereby claiming ownership of every address on the network.

This specification makes this existing ARP vulnerability no worse, and in some ways makes it better: Instead of failing silently with no indication why, hosts implementing this specification are required to either attempt to reconfigure automatically, or if not that, at least inform the human user of what is happening.”

Although some may argue about the question of whether or not it is effective, appropriate, and useful to “inform the human user” in this case, this solution nevertheless follows the principle of at least not worsening the current security situation of an existing protocol.

Zeroconf Multicast Address Allocation Protocol

The *Zeroconf Multicast Address Allocation Protocol* (ZMAAP) defined in^[5] specifies a method for peer-to-peer allocation of *multicast addresses without a multicast* (MADCAP) server in small zero configuration networks. The word “small” is important here because ZMAAP is not scalable beyond small networks (and is not designed to be).

Multicast DNS

“Performing DNS queries via IP Multicast”^[6] by Stuart Cheshire suggests some very useful ideas on how to use mDNS with maximum benefit and minimum hassle in zero configuration networks. In my opinion, the best thing about this proposal is that it does not require any changes to the DNS protocol (messages, resource record types, etc.) itself. Instead it concentrates on the use of multicast for name resolution in environments where no DNS servers exist (and where one would not reasonably expect them to). The goal is to have a working name resolution service without name servers. The document proposes to use **local.arpa** (although the exact choice of this special domain is not the goal of this document) as the link-local domain (like the **169.254/16** network for dynamic allocation of IPv4 link-local addresses described earlier in this article). For reverse address resolution, **254.169.in-addr.arpa** is also link-local. The multicast address **224.0.0.251** that is used for mDNS queries is registered by the IANA for this purpose. No delegation is performed within mDNS domain **local.arpa**. There is also no *Start of Authority* (SOA) record for the mDNS domain because of the nature of zero configuration networks where it is intended to be used—in particular, there is no mailbox responsible for the zone. Likewise, zone transfers are not applicable with mDNS zones. To summarize, any local link has its own local and private **local.arpa** and **254.169.in-addr.arpa** zones, which have only link-local significance in the particular Zeroconf network.

DNS Service Discovery

Like the multicast DNS solution described previously, the *DNS Service Discovery* (DNS-SD)^[7] does not require any changes to the existing DNS protocol; thus it is completely compatible with the existing DNS server and client software.

What DNS-SD proposes is a naming scheme for *DNS Resource Records* (RRs) to allow for service discovery using the existing DNS—either the traditional or multicast DNS described in the previous paragraph. DNS-SD uses the SRV and PTR resource records to provide the required functionality. To cite from [7]:

“Service discovery requires a central aggregation server. DNS already has one: It’s called a DNS server.

Service discovery requires a service registration protocol. DNS already has one: It’s called DNS Dynamic Update.

Service discovery requires a security model. DNS already has one: It’s called DNSSEC.

Service discovery requires a query protocol. DNS already has one: It’s called DNS.”

It is necessary to note that DNS-SD is compatible with mDNS and vice versa, but neither requires the other one to function. However, it is practical to use mDNS for service discovery (using DNS-SD) to have a single protocol and interface and not have to implement another protocol just for service discovery.

Industry Support

Any new technology needs industry support to succeed, and Zeroconf is no exception. Several major vendors have announced plans to support or already support Zeroconf in their products, including Apple, Epson, Hewlett-Packard, Lexmark, Philips, Canon, Xerox, Sybase, and WorldBook. One can expect that more companies will Zeroconf-enable their products as the technology itself matures and hopefully becomes standardized and widespread.

Rendezvous

Rendezvous is Apple Computer’s implementation of Zeroconf in its Darwin 6 and Mac OS X 10.2 (“Jaguar”) operating systems. Apple has stated its full support for the Zeroconf and intent to completely replace the aging AppleTalk with Zeroconf-enabled Macs, without sacrificing the ease of use and transparency to end users provided by AppleTalk networks. A good example of Zeroconf’s use in OS X would be the iChat instant messaging (IM) client, which comes with the Version 10.2 of Mac OS X. It works not only with AOL *Instant Messenger* (AIM) and Mac networks but may also be used between Zeroconf-enabled Macs in a Zeroconf network.

Coupled with Apple’s implementation of IEEE 802.11b (“WiFi”) in ad hoc mode, it permits a wireless zero configuration network that just works without any configuration or additional hardware or software.

Apple has also made the source code for the mDNS Responder, a part of Rendezvous implementing mDNS, freely available through the Darwin Open Source Project. Mac OS X software developers are encouraged to use Zeroconf, and there are documentation and application examples to facilitate this. More information about Rendezvous and Zeroconf on Macs is available from Apple’s Web sites^[9].

Summary

With computers and computer networks becoming more and more complex and sophisticated, some people (including the author of this article) believe that care should be taken by those in the know not to create more problems than we solve using these computers and networks. Yes, we want more features—but we also need to remember that most users of these features do not have doctorates in computer science and (surprise, surprise) don't even wish to. Zero configuration networking would probably help in this regard, minimizing and even eliminating in some cases the need to configure and administer small networks. Let me conclude by quoting once more from the Zeroconf Working Group:

“It is important to understand that the purpose of Zeroconf is not solely to make current personal computer networking easier to use, though this is certainly a useful benefit. The long-term goal of Zeroconf is to enable the creation of entirely new kinds of networked products, products that today would simply not be commercially viable because of the inconvenience and support costs involved in setting up, configuring, and maintaining a network to allow them to operate.”

References

- [1] Zeroconf Working Group, Internet Engineering Task Force (IETF): <http://www.ietf.org/html.charters/zeroconf-charter.html>
- [2] Aidan Williams, “Requirements for Automatic Configuration of IP Hosts,” **draft-ietf-zeroconf-reqts-12.txt**
- [3] Stuart Cheshire, Bernard Aboba, and Erik Guttman, “Dynamic Configuration of IPv4 Link-Local Addresses,” **draft-ietf-zeroconf-ipv4-linklocal-07.txt**
- [4] Stuart Cheshire, “IPv4 Address Conflict Detection,” **draft-cheshire-ipv4-acd-02.txt**
- [5] Octavian Catrina, Dave Thaler, Bernard Aboba, and Erik Guttman, “Zeroconf Multicast Address Allocation Protocol (ZMAAP),” **draft-ietf-zeroconf-zmaap-02.txt**
- [6] Stuart Cheshire, “Performing DNS Queries via IP Multicast,” **draft-cheshire-dnsext-multicastdns-00.txt**
- [7] Stuart Cheshire, “Discovering Named Instances of Abstract Services Using DNS,” **draft-cheshire-dnsext-nias-00.txt**
- [8] Zeroconf: <http://www.zeroconf.org>
- [9] Rendezvous: <http://developer.apple.com/macosx/rendezvous/>
<http://www.apple.com/macosx/jaguar/rendezvous.html>
- [10] Erik Guttman, “Autoconfiguration for IP Networking: Enabling Local Communication,” *IEEE Internet Computing*, June 2001.

EDGAR DANIELYAN is a self-employed consultant, author, and editor specialising in UNIX, networking, and information security. In previous life he has been a cofounder of a national ISP and manager of a country TLD. He is currently working on his next book (*WLAN Security*) which is due to be published in 2003. His previous book, *Solaris 8 Security*, was published by New Riders Publishing in 2001. He is also a member of IEEE, IEEE Standards Association, IEEE Computer Society, ACM, USENIX, and the SAGE. He is online at <http://www.danielyan.com> and can be reached by e-mail at edd@danielyan.com

Book Reviews

Ruling the Root *Ruling the Root: Internet Governance and the Taming of Cyberspace*, by Milton L. Mueller, ISBN 0-262-13412-8, The MIT Press, 2002, <http://mitpress.mit.edu>

“WASHINGTON, Apr. 1 /Governance Newswire/ — The organizations that create street names, assign addresses, and assign telephone numbers have issued a joint announcement: Henceforth any conversation not conducted in Bahasa Malayu will result in termination of the relevant address or telephone number assignment.”

The above bit of fiction is not pure silliness. Fear of equivalent, Internet-related excesses is the essence of Milton Mueller’s book, *Ruling the Root*. The Syracuse University professor believes that administration of Internet addresses and domain names provides a fulcrum for overall Internet governance. He says they create a “political economy” vulnerable to serious abuse. Domain name administration is equated with control over Internet content, because, “a domain name record [is] very much like an Internet driver’s license” as if it provides permission to use the Net, and even authorizes the locations one may visit.

Organization

The book covers both IP address and domain name administration. The material on IP addresses is thin, perhaps because it is a well-managed area without significant controversy. This is in marked contrast to the recent history of debate on *Domain Name System* (DNS) oversight. So it might have been instructive to see a comparison between the two administrative models, beyond simply noting that domain names can be interesting.

Discussion covers Internet technology, the history and politics of DNS and IP administrative management structure, and the intellectual property aspects of name assignment conflicts. Mueller suggests a three-layer hierarchy: technical, economic, and policy. What is missing from this “architecture” and from the entire book is any concern for the pragmatic details of administration and operation of these global, mission-critical services. Yet such tasks are difficult to perform well, as Network Solutions repeatedly demonstrated over the years, by losing registrations and corrupting critical data files; and the effects of problems are large.

When *Star Trek*’s Captain Picard commands, “make it so,” we know that he fully appreciates the challenges in implementing his directive. However, for *Ruling the Root*, policy development is not concerned with the operational complexities.

Not surprisingly, the book often demonstrates a misunderstanding of constraints inherent in DNS technology, although the tutorial on basic Internet technology is adequate, in spite of making the common error about the “T” in TCP/IP.^[1]

Differing Opinions

Other reviewers of the book have called it well written, insightful, and nuanced. Indeed the discussion of history that is fully documented and involves simple, clear, objective facts is quite good. The rest of the time Mueller presents biased and unfounded descriptions of Internet governance, motives, and decisions, while failing to distinguish between what is fact and what is his opinion.

Ruling the Root sees adversaries, conspiracies, and threats, and permits no balancing sense of diverse collaboration, constructive criticism, or productive compromise. The technical community is somewhat less suspect, but is deprecated with the usual cliché about its naivete. So Mueller misses the essential point that techies designed, built, operated, and grew this robust, survivable, equitable system for global operations and service governance.

Professor Mueller’s treatment of the dominant DNS registry, *Network Solutions* (NSI), now VeriSign, is curiously superficial and soft. NSI benefited spectacularly from the National Science Foundation’s decision to permit charging for domain names, and from the policies and delays in the formation of the *Internet Corporation for Assigned Names and Numbers* (ICANN), as well as ICANN’s distraction away from its intended registry oversight function and toward abstract debates about Internet governance. Yet the book does not consider NSI’s role in ICANN-related political processes.

Mueller fails to understand the history of the organization that managed the DNS from its inception, the *Internet Assigned Numbers Authority* (IANA) and Jon Postel’s role in running it. IANA is incorrectly represented as a simple operations arm of the U.S. Government. The grass-roots basis for its real legitimacy is missed. Its policy role is missed. Its collaborative processes are denied. For example, Mueller tells us that the description of IANA in RFC 1083, published in 1988 meant, “a new world was being defined by the RFC.” In reality it was simply documenting established practice, as is typical for operations RFCs.

Validation

Mueller’s substantiation of his analyses is also problematic. The book must be read with careful attention to the actual authority of each source. Goals and agendas are often misstated. For example, he characterizes the pre-ICANN *International Forum for the White Paper* (IFWP) as “the real arena for arriving at a decision [about the details of the new organization].” Its actual goal was simply to be a forum for discussion. Discussion, not decision-making.^[2]

The book claims that the pre-ICANN *International Ad Hoc Committee* (IAHC) was formed “to develop and implement a blueprint for a global governance structure for the domain name system.” In fact, the IAHC was formed for “specifying and implementing policies and procedures relating to iTLDs (international top-level domains, now called ‘generic’ TLDs, or gTLDs).”^[3] He claims, “They had asserted that the root was theirs to dispose of.” To the contrary, the IAHC was explicitly subordinate to IANA, and had nothing at all to do with management of the DNS root or any non-gTLD part of the DNS. Interestingly, the endnote Mueller offers as substantiation disproves his characterization.

Ruling the Root is loaded with endnotes—27 pages of small print. However, even the formal citations are problematic. Note #55 cites a newspaper article as a primary source, as if it were definitive proof the person discussed in the article held a specific opinion. Mueller’s Note #45 claims to substantiate that, “Postel himself... admitted...it is unclear who actually controls the name space.” Yet the note is for *Internet Architecture Board* (IAB) minutes. Attributing it to Postel was a fabrication.

Back-room, deal-making, conspiracy explanations are offered without substantiation. Of changes to *Internet Engineering Task Force* (IETF) management, Mueller states: “The most important reason the IETF didn’t institute voting was that Jon Postel and several other senior figures vowed that they would refuse to run for office.” Postel never made such a vow, and the process to effect these IETF changes did not experience any such attempts at influence. Of Postel’s instructing some root servers to retrieve copies of the DNS root from a non-NSI master, Mueller claims that Postel was “apparently concerned about the direction U.S. policy was taking.”

No substantiation is offered, because the claim is false. Postel and others were concerned about NSI’s reaction to its own loss of control. The switch was intended to see what it would take to move NSI out of the hierarchy. These are not small matters of nuance. They show a pattern of misrepresentation.

The Author

Professor Mueller’s credibility would have been aided by disclosing his own affiliations. The only ICANN constituency (the Non Commercial Domain Name Holders Constituency) claiming to represent the non-commercial world focuses on the civil society concerns that dominate the public debate about ICANN. Professor Mueller’s discussion of the group is quite thin and does not disclose the fact that he held a dominant management position in it. In his criticism of dispute-resolution activities, he neglects to mention that he is a paid arbitration panelist.

An important book should be read because it has factual detail and thoughtful insight. *Ruling the Root* is, instead, important because it so thoroughly embodies the difficulties that have emerged in discussing Internet policy. Because so many people take *Ruling the Root* seriously, it should be read. However, the serious problems of the book encourage borrowing it, rather than buying a copy. Based on the pattern noted in this review, a thorough audit of those problems would be appropriate for the relevant Syracuse University academic ethics committee.

—Dave Crocker^[4], *Brandenburg Internet Working*
dcrocker@brandenburg.com

References

- [1] The “T” stands for transmission, not transport or transfer.
- [2] <http://web.archive.org/web/19981206105122/http://www.ifwp.org/>
- [3] <http://www.iahc.org/iahc-charter.html>
- [4] Factual claims in the review that do not have citations are based on the reviewer’s direct experience. Dave Crocker wrote the first Internet standard for domain name syntax (RFC 822). He also was the IETF area director for initial work on DNS security. More recently he was one of Jon Postel’s appointees to the IAHC. He naively thought that its work should be conducted in the manner that had been typical for Internet administration. So the last few years of charged, global politicization have been an education. He must also note that he was once Jon Postel’s officemate.

High-Speed Networks and Internets

High-Speed Networks and Internets: Performance and Quality of Service, 2nd ed., by William Stallings, ISBN 0-13-032221-0, Prentice Hall, 2002. <http://www.prenhall.com/stallings>

This thoroughly updated classic covers topics of traffic engineering, queuing, and traffic modeling. The book gives a complete look around the protocols of the next generation: *Resource Reservation Protocol* (RSVP), *Multiprotocol Label Switching* (MPLS), and *Real-Time Transport Protocol* (RTP). It gives the keys to understand the way Frame Relay, TCP, and ATM react to congestion and flow control. The book also deals with new trends and standards that will lead the telecommunications industry in the following years. A very useful book, from the same author of traditional titles such as: *Data Communications*, *Cryptography*, *Computer Architecture*, and many more.

Organization

High-Speed Networks is divided into seven parts. The first one discusses the basic background needed to understand the rest of the book. Following the introduction, the second chapter goes on with the classical: the *Open System Interconnection* (OSI) model and the TCP/IP suite.

Part II explains packet-switching technologies in detail. The forth chapter explains the architecture of Frame Relay, and the next one focuses on ATM, including its operation and the adaptation layers. Chapter 6 works on high speed LANs, covering Fast Ethernet and Gigabit Ethernet, with the different media supported by each.

The third part is one of the most important; chapter 7 presents an overview of probability and stochastic processes. Although it is a brief one, it is useful to make revision of some concepts. The next chapter works on queuing analysis, introducing the basic elements of a queuing model. It explains the topics with plenty of examples: M/M/1, multiserver queues, and networks of queues, presenting all the formulas. Chapter 9 is dedicated to self-similar traffic. As recent studies indicate, traffic on high speed networks does not have the characteristics needed for the queuing theory. It introduces and explains the concept of self-similarity. Then the author applies this concept to data traffic analysis and examines performance implications. Based on papers on this subject, Stallings explains this new approach to traffic modeling not analyzed before.

The forth part focuses on another main topic: congestion and traffic management. Chapter 10 explains the effects of congestion and the different ways to control and avoid it. In the following chapter the author discusses control mechanisms at the link level. He examines different ways used by protocols to handle flow control: *Stop and Wait*, *Sliding Window*, and *Go back N-ARQ*. An analysis of the performance gained by using *Automatic Repeat Request* (ARQ) techniques follows.

These chapters give a detailed description of the different ways that communications can be handled. Chapter 12 focuses on transport-level traffic management. It explains TCP flow control in detail, including the retransmission strategy. The way TCP avoids congestion is discussed thoroughly. The next chapter continues with congestion control in ATM networks. The framework for traffic control is explained in detail, with sections dedicated to *Available-Bit-Rate* (ABR) and *Guaranteed-Frame-Rate* (GFR) traffic management.

The next part of the book is about Internet routing. Chapter 14 presents the algorithms used to compute the minimum path, and introduces some elementary concepts in graph theory. Later the author concentrates on Interior routing protocols, analyzing the *Routing Information Protocol* (RIP) and *Open Shortest Path First* (OSPF), the most important ones. Next the book discusses exterior routing protocols and multicast. The author describes in a simple way these addressing schemes and the related protocols.

The following section is dedicated to *Quality of Service* (QoS) in IP networks. The first chapter discusses integrated services, with coverage of queuing disciplines such as *Weighted Fair Queuing* (WFQ). A review of the Differentiated Services architecture follows.

After discussing the concepts, the author examines the protocols that support QoS: RSVP, MPLS, and RTP. He explains the philosophy behind each protocol, its characteristics, and its implementation.

In the final part of the book, the author changes the subject to compression. In Chapter 19 he presents an overview of information theory, discussing typical areas such as entropy. The next chapter continues with loss-less compression, facsimile compression, and others. It discusses the Lempel-Ziv algorithm used in PKZIP. The final chapter reviews lossy compression, explaining the discrete cosine transform, a key component of the *Joint Photographics Expert Group* (JPEG) and *Motion Picture Experts Group* (MPEG) standards.

Two very interesting appendices end the book: one for Internet standards and the standardization process and the other one dedicated to sockets, containing source code. Although the book is not dedicated to programming, the inclusion of TCP sockets can be useful to understand its implementation.

A book worth reading

We are facing an essential book for networking professionals, designers, and engineers. It covers unusual topics such as self-similar traffic and data compression. It is the basement for the design of any high speed network. As Internet traffic continues to grow, the optimization of network resources becomes a critical topic. Also, more and more voice traffic is carried over packet networks, congestion being one of its worst enemies. The time-sensitive traffic needs attention, and this book provides the tools to manage it.

In addition to its solid coverage of topics, the book has plenty of bibliography and many links to the principal sites for each chapter. With no doubt this is a very useful book, from the well-known technical author William Stallings.

—Rodrigo J. Plaza, *Iplan Networks, Argentina*
rplaza@iplan.com.ar

Would You Like to Review a Book for IPJ?

We receive numerous books on computer networking from all the major publishers. If you've got a specific book you are interested in reviewing, please contact us and we will make sure a copy is mailed to you. The book is yours to keep if you send us a review. We accept reviews of new titles, as well as some of the "networking classics." Contact us at **ipj@cisco.com** for more information.

Letters to the Editor

ENUM Ole,

As the co-chair of the ENUM work group in the IETF, I was delighted with Geoff Huston's article. (*The Internet Protocol Journal*, Volume 5, No. 2, June 2002, page 13).

I would like to point out and clarify several other issues raised by the Letters to the Editor published in the subsequent issue.

First, as a practical matter though the North American Numbering Plan uses a single country code "1," there will not be a single administration of ENUM within "1." The agreements between the IAB and the ITU on the administration of **e164.arpa** clearly indicate that these resources will be administered on a nation-state basis.

www.iab.org/DOCUMENTS/enum-pr.html

www.iab.org/DOCUMENTS/sg2-liaison-e164-sep-02.html

The United States, Canada, Bermuda, and the 18 countries of the NANP will be free to administer their numbering resources as they so choose through the use of 1 + NPA (area codes) zones within the root of **e164.arpa**.

Dr. Deleuze writes, "E.164 numbers are really telephone addresses. They are tied to telephone network topology and are surely not user friendly. There are no user-friendly names in the telephone system."

In fact, this is not exactly correct either. Since the advent of Number Portability by several national telephone administrations, including the United States, telephone numbers are no longer tied to the underlying network or routing structure of the PSTN. Actual routing of phone calls in the United States is done on Local Routing Numbers for all landline calls and, beginning in November of 2003, for wireless calls as well.

Phone numbers even now are essentially names, much like domain names in the Internet. In the United States, phone numbers can be taken or "ported" to any wireline service provider within proscribed geographic boundaries, in 2003 between wireless service providers and from wireline to wireless providers as well.

I partially take issue with Dr. Deleuze's thought that telephone numbers are not "user-friendly." Phone numbers are readily identifiable, easy to use, and are not tied to culture or language, problems we have not yet solved with domain names.

—Richard Shockey, NeuStar Inc.
rich.shockey@NeuStar.com

Visitor Networks Dear Editor,

The September 2002 issue of IPJ featured a very interesting, comprehensive article on visitor networks. One aspect I found not mentioned, however, is the danger of users in such scenarios falling victim to fake visitor gateways. In public wireless hot spots, as they are increasingly being setup at numerous locations these days, attackers could employ their own mobile WLAN device to direct visitors trying to log on to the hot spot to their own fake login page, enabling them to easily collect their login details such as credit card information. Using encryption does not help here as long as the gateway does not need to authenticate itself to the customer's mobile device. The average user should not have a chance to realize whether he or she is connected to a legitimate or a fake login page—if he or she is aware of that potential danger at all. Given the fact that all such an attack would need, apart from readily available equipment such as a portable computer with a WLAN card, is some small piece of appropriate software and that it would be quite difficult to detect, that kind of threat unfortunately should be quite realistic in such environments.

—Dr. Georg Schwarz
Detecon International GmbH, Berlin, Germany
Georg.Schwarz@detecon.com

The author responds:

This is a good point that was not discussed in the article. There are actually at least three cases that visitors need to worry about. The first is, as you mentioned, that the service provider is not who they say they are. This can be dealt with by using SSL certificates assuming the visitor is conscious of the URL that he/she is being directed to and knows that it belongs to the real service provider. If the visitor has no idea who is a reasonable service provider, this is a different class of problem, very similar to what has happened with public telephones that accept standard calling and credit cards—someone makes a call, receives the service but then gets charged an outrageous rate. The third case is a man-in-the-middle attack or passive snooping where someone with a laptop as you describe is able to grab traffic and gather passwords.

Some basic advice to visitors is for services that require subscription, although possibly inconvenient, never subscribe on a potentially compromised connection. That way, only the service provider-assigned username and password is compromised, instead of more sensitive personal information related to the account. Connections using 802.1x authentication with EAP-TLS provide mutual authentication and are in the long run, a better solution than redirection of web pages. No matter what kind of security one has, inevitably there will be legally legitimate providers that will take advantage of visitors and in that case it's just "buyer beware."

—Dory Leifer
leifer@del.com

Again, I found the latest issue of IPJ quite enlightening and useful. However, I do have one comment regarding the article by Greg Scholz on “An Architecture for Securing Wireless Networks.” Although the use of source IP addresses to provide policy group membership on the firewall works in most cases, some client OSs and some IPsec VPN boxes allow the source address (even if it is the endpoint address of the tunnel, not the “real” address of the host) to be changed, provided the source address of the enciphered traffic does not change. This would allow users to change the policy group they belong to. A better solution is to use a VPN box that can associate groups of IPsec tunnels to VLANs. Then the firewall could be configured to allow policy group membership based on VLANs. This takes all determination of policy group membership off the client host and places it in the domain of trust of the VPN and firewall boxes.

—Chris Liljenstolpe
Cable and Wireless
chris@cw.net

Fragments

Upcoming Events

The IETF will meet in San Francisco, California, USA March 16–21, 2003. The IETF will also meet in Vienna, Austria, July 13–18, 2003 and in Minneapolis, Minnesota November 9–14, 2003.

See <http://www.ietf.org/meetings>

The next APRICOT (*Asia and Pacific Regional Internet Conference on Operational Technologies*) will be held in Taipei, Taiwan, February 19–28. See <http://www.apricot2003.net/>

The *Internet Corporation for Assigned Names and Numbers* (ICANN) will meet in Rio de Janeiro, Brazil, March 23–27, 2003, in Montreal, Canada, June 22–26, 2003, and in Carthage, Tunisia, December 1–5, 2003. See <http://www.icann.org>

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Technology
WorldCom, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 2002 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRST STD U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol Journal

March 2003

Volume 6, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Measuring IP Networks.....	2
Session Initiation Protocol	20
Letters to the Editor.....	31
Book Review.....	36
Call for Papers	39

FROM THE EDITOR

Even the most carefully designed and operated IP network is subject to any number of performance problems ranging from overloaded links and mis-configured routers to server failures. For these situations, the network manager has several diagnostic tools as options. Geoff Huston gives us an overview in an article entitled “Measuring IP Network Performance.”

Voice over IP (VoIP) is an emerging application, as well as a rapidly growing market. Use of the corporate network or the Internet at large to carry telephone traffic has many advantages, not the least economic ones. A successful VoIP network must not only support IP-based telephones, but also provide a means of seamlessly integrating the IP-based network with traditional telephone networks. At the core of VoIP lies the *Session Initiation Protocol* (SIP) and a few related protocols. Bill Stallings describes SIP in our second article.

Book reviews published in *The Internet Protocol Journal* can rarely be characterized as “controversial.” However, when the book in question deals with ICANN, it is perhaps not surprising that strong opinions emerge. Thus, following the review of *Ruling the Root* in our last issue, we received a letter from the author that is included in our “Letters to the Editor” section (along with a response from the book reviewer). I would like to take this opportunity to remind our readers that book reviews do represent the *opinion* of the reviewer and should be read in that light.

Our online subscription system has been up and running for a couple of months. Please give it a try at: www.cisco.com/ipj.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Measuring IP Network Performance

by Geoff Huston, Telstra

If you are involved in the operation of an IP network, a question you may hear is: “How *good* is your network?” Or, to put it another way, how can you measure and monitor the quality of the service that you are offering to your customers? And how can your customers monitor the quality of the service you provide to them?

These questions have been lurking behind many public and enterprise IP networks for many years now. With the increasing levels of deployment of various forms of high-speed (or broadband) services within today’s Internet there is new impetus to find some usable answers that allow both providers and users to place some objective benchmarks against the service offerings. With the lift in access speed with broadband services, there is an associated expectation on the part of the end user or service customer about the performance of the Internet service. It should be “better” in some fashion, where “better” relates to the performance of the network and the service profile that is offered to network applications. And not only is there an expectation of “better” performance, it should be measurable. This article looks at network performance and explores its definition and measurement.

A Functional Definition of Network Performance

An informal functional approach to a definition of network performance is measuring the speed of the network. How fast is the network? Or, what is the elapsed time for a particular network transaction? Or, how quickly can I download a data file? This measurement of time for a network transaction to complete certainly relates to the speed of the network, and speed is a good network performance benchmark, but is speed everything?

When looking at the broad spectrum of performance, the answer is that speed is not everything. The ability of a network to support transactions that include the transfer of large volumes of data, as well as supporting a large number of simultaneous transactions, is also part of the overall picture of network load and hence of network performance. But large data sets is not everything in performance. Consideration should also be given to the class of network applications where the data is implicitly clocked according to some external clock source. Such real-time applications include interactive voice and video, and their performance requirements include the total delay between the end points, or latency, as well as the small-scale variation of this latency, or *jitter*. Such performance measurements also include the ratio of discarded packets to the total number of packets sent, or loss rate, as well as the extent to which a sequence of packets is reordered within the network, or even duplicated by the network. Taken together, this set of performance factors can be considered as a form of the amount of distortion of the original real-time signal.

Accordingly, a functional description of network performance encompasses a description of speed, capacity, and distortion of transactions that are carried across the network. This informal description of what

constitutes network performance certainly feels to be on the correct path, given that if one knew the latency, available bandwidth, loss, and jitter rates and packet reorder probability as a profile of network performance between two network end points, as well as the characteristics of the network transaction, it is possible to make a reasonable prediction relating to the performance of the transaction.

Taking this informal definition, the next step is to create a more rigorous framework for measuring performance. For any single network path between an entry and egress point, it is possible to measure the path latency, available peak bandwidth, loss rates, jitter profile, and reorder probability. But there is a difference between a description of the performance of a particular path across a network and the performance of the network as an aggregate entity. Given a set of per-path performance measurements, how can you construct a view of the performance of the network? A common methodology is to take a relatively complete set of path measurements across a network and then combine them to create an average metric. Although this accomplishes a useful reduction in the size of the data, there is also a loss of information. The average network performance measurements have little relationship to the performance of any individual path.

There are various ways to improve this loss of information, including weighting the individual path measurements by the amount of traffic passed along the path. Such techniques are indeed to ensure that paths that use far-flung network outliers that carry relatively low volumes of traffic have a much lower impact on the overall network performance metric than the major network transit paths.

Measuring Network Performance

Given these performance indicators, the next step is to determine how these indicators may be measured, and how the resulting measurements can be meaningfully interpreted. At this point it is useful to look at numerous popular network management and measurement tools and examine their ability to provide useful measurements. There are two basic approaches to this task; one is to collect management information from the active elements of the network using a management protocol, and from this information make some inferences about network performance. This can be termed a *passive approach* to performance measurement, in that the approach attempts to measure the performance of the network without disturbing its operation. The second approach is to use an active approach and inject test traffic into the network and measure its performance in some fashion, and relate the performance of the test traffic to the performance of the network in carrying the normal payload.

Measuring Performance with SNMP

In IP networks the ubiquitous network management tool is the *Simple Network Management Protocol* (SNMP). There is no doubt that SNMP can provide a wealth of data about the operational status of each management network element, but can it tell you anything about the overall network performance?

The operation of SNMP is a *polling* operation, where a management station directs periodic polls to various managed elements and collects the responses. These responses are used to update a view of the operating status of the network.

The most basic tool for measuring network performance is the periodic measurement of the interface byte counters. Such measurements can provide a picture of the current traffic levels on the network link, and when related to the total capacity of the link, the relative link loading level can be provided. As a performance indicator this relative link loading level can provide some indication of link performance, in that a relatively lightly loaded link (such as a load of 5 to 10 percent of total available capacity) would normally indicate a link that has no significant performance implications, whereas a link operating at 100 percent of total available capacity would likely be experiencing high levels of packet drop, queuing delay, and potentially a high jitter level. (Figure 1) In between these two extremes there are performance implications of increasing the load. Of course it should be noted that the characteristics of the link have a bearing on the interpretation of the load levels, and a low-latency 10-Gbps link operating at 90-percent load will have very significantly lower levels of performance degradation than a 2-Mbps high-latency link under the same 90-percent load. (Figure 2)

Figure 1a: Relative Link Loading – An Optimally Loaded Link

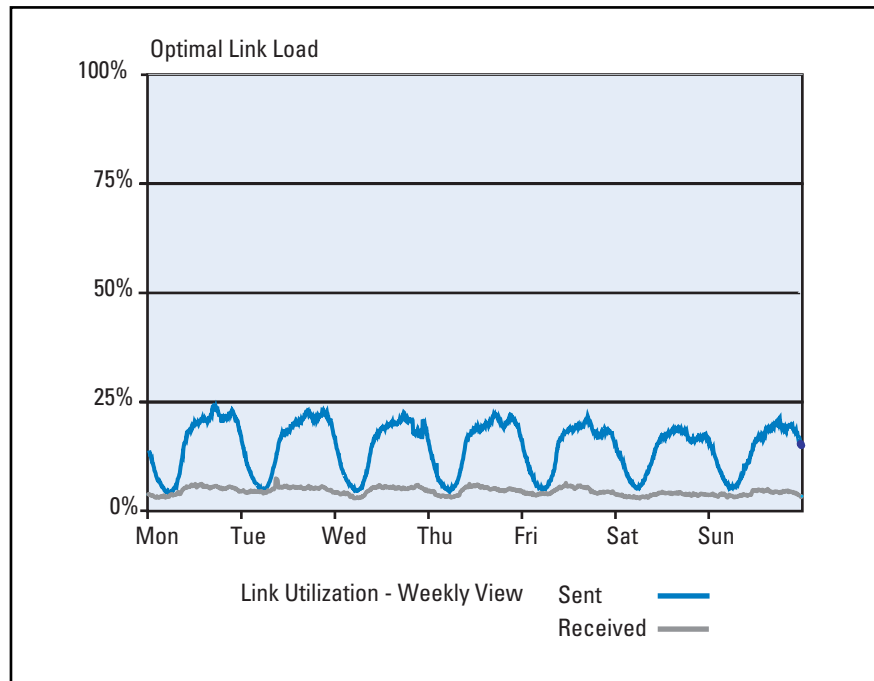


Figure 1b: Relative Link Loading – A Maximally Loaded Link

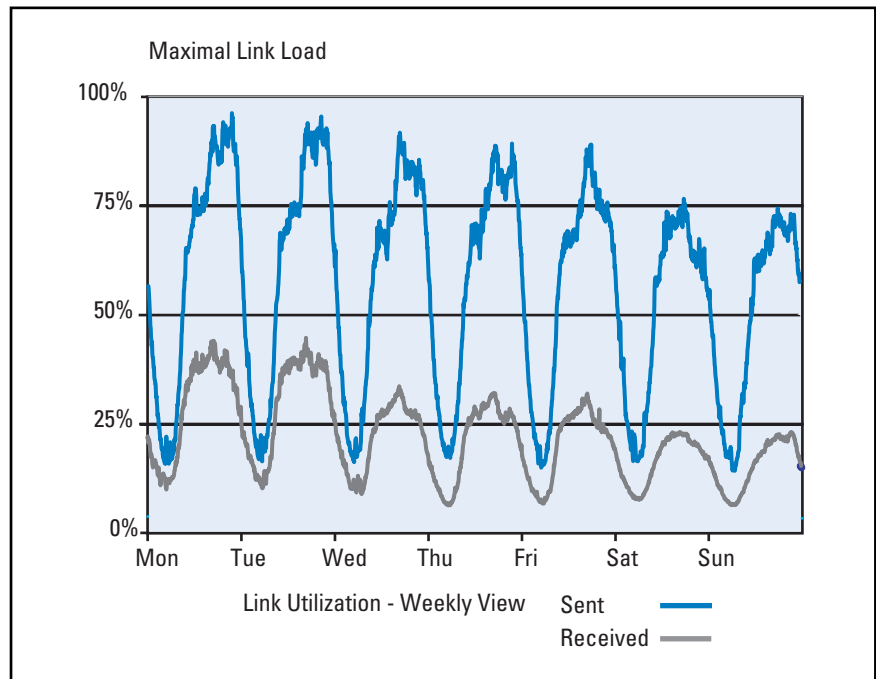


Figure 1c: Relative Link Loading – Highly Degraded Link

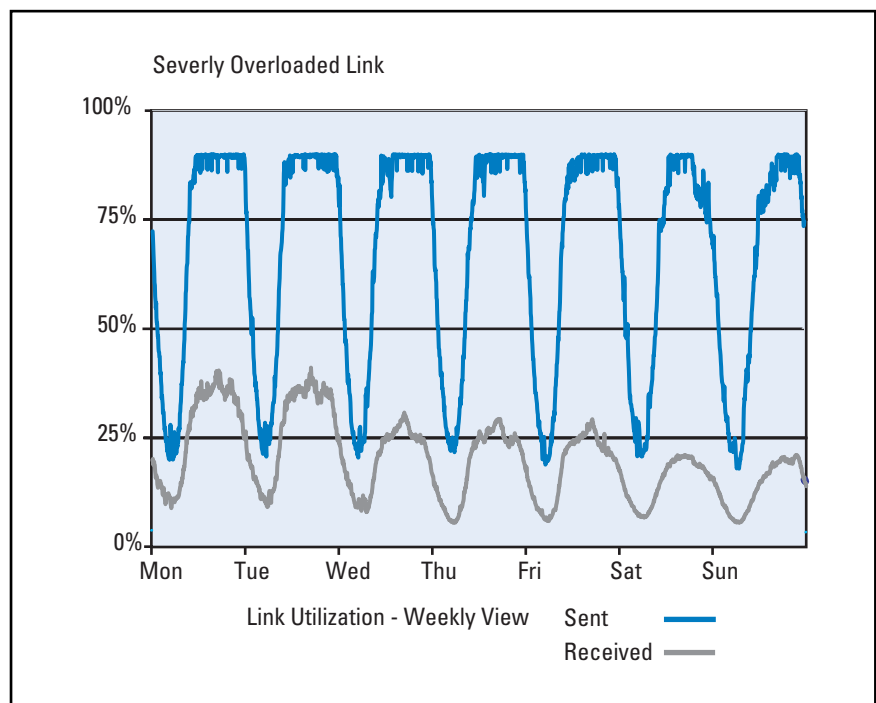
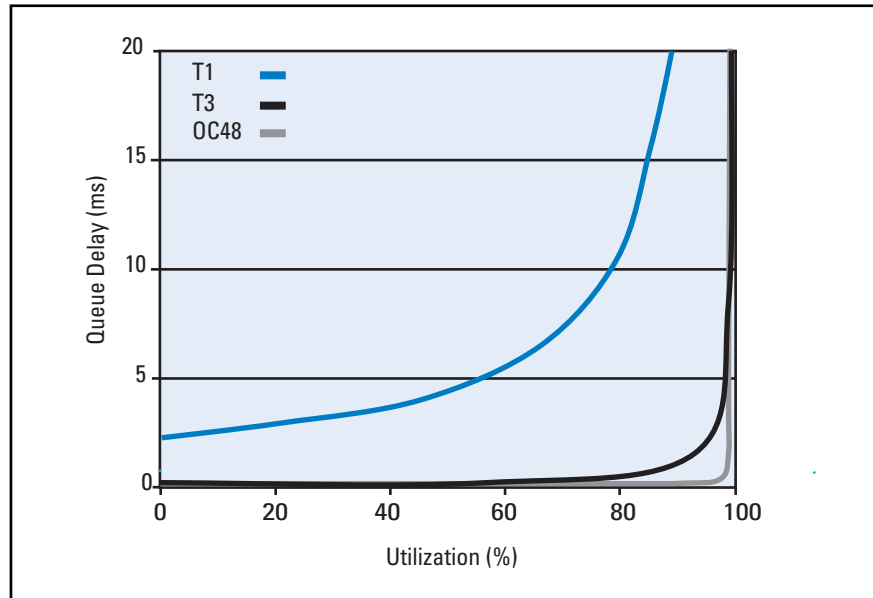


Figure 2: Queuing Delay Simulation
(David Meyer, Sprint, November 2002)



Relative traffic load on each link can be complemented by measurement of performance-related SNMP counters. A management system can poll each active network element to retrieve the number of packets dropped for each interface, and the number of packets successfully forwarded. From these two data items, the relative drop proportion of packets can be calculated on an element-by-element and potentially a link-by-link basis, and a series of element measures can provide a per-path drop proportion by combining the individual packet-forwarding measurements for the interfaces on the path.

Because some count of relative packet drop rate can be gathered from each network element, with the additional input of the current forwarding state of the network it is possible to predict the path a packet will take through the network, and hence estimate the path probability of drop. However, this information is still well short of being a reliable measurement of service performance.

Queuing delay is somewhat more challenging to measure on an element-by-element basis using element polling with SNMP. In theory, the polling system could use a rapid sequence of polling the output queue length of a router and estimating the queuing delay based on an average packet size estimate, together with the knowledge of the available output capacity. Of course, such a measurement methodology assumes a simple *first-in, first-out* (FIFO) queuing discipline, a queue size that varies slowly over time, and slow link speeds. Such assumptions are rarely valid in today's IP networks. As the link speed increases, the queue size may oscillate with a relatively high frequency as a function of both the number and capacity of the input systems and of the capacity of the output system. In general, queuing delay is not easily measured using network element polling.

There is no ready way for a polling mechanism to detect and count the incidence of reordered packets. Packet reordering occurs in many situations, including the use of parallel switching fabrics within a single network element and the use of parallel links between routers.

IP routers are not typically designed to detect, let alone correct, packet reordering and because they do not detect this condition, they cannot report on the incidence of reordering via SNMP polling.

The generic approach of network management polling systems is that the polling agent, the network management station, is configured with an internal model of the network; status information, gathered through element polling, is integrated to the network model. The correlation of the status of the model to the status of the network itself is intended to be accurate enough to allow operational anomalies in the network to be recognized and flagged. The challenge is that a sequence of snapshots of element status values cannot readily be reconstructed into a comprehensive view of the performance of the network as an entire system, or even as a collection of edge-to-edge paths. Measurement techniques using polling and modeling can track the performance of the individual elements of the network, but they cannot track per-path service levels across the network. The network-element polling approach can indicate whether or not each network element is operating within the configured operational parameters, and alert the network operator when there are local anomalies to this condition. But such a view is best described as *network centric*, rather than service centric. An implicit assumption is that if the network is operating within the configured parameters, then all service-level commitments are being met. This assumption may not be well founded.

The complementary approach to performance instrumentation of network elements is active network probing. This requires the injection of marked packets into the data stream; collection of the packets at a later time; and correlation of the entry and exit packets to infer some information regarding delay, drop, and fragmentation conditions for the path traversed by the packet. The most common probe tools in the network today are *ping* and *traceroute*.

Measuring Performance with Ping

The best known, and most widely used active measurement tool is *ping*. Ping is a very simple tool: a sender generates an *Internet Control Message Protocol* (ICMP) echo request packet, and directs it to a target system. As the packet is sent, the sender starts a timer. The target system simply reverses the ICMP headers and sends the packet back to the sender as an ICMP echo reply. When the packet arrives at the original sender's system, the timer is halted and the elapsed time is reported. An example ping output is shown in Figure 3.

Figure 3: Example Ping Report

```
% ping www.iab.org
PING www.iab.org (132.151.6.25): 56 data bytes
64 bytes from 132.151.6.25: icmp_seq=0 ttl=44 time=254.409 ms
64 bytes from 132.151.6.25: icmp_seq=1 ttl=44 time=254.197 ms
64 bytes from 132.151.6.25: icmp_seq=2 ttl=44 time=255.238 ms
64 bytes from 132.151.6.25: icmp_seq=3 ttl=44 time=255.874 ms
--- www.iab.org ping statistics ---
4 packets transmitted, 4 packets received, 0% packet loss
round-trip min/avg/max/stddev = 254.197/254.930/255.874/0.670 ms
```

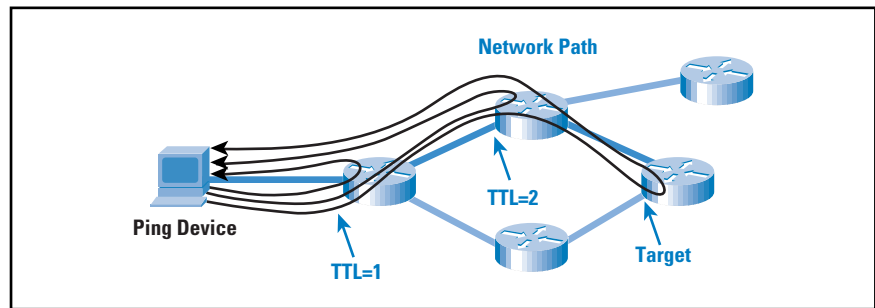
This simple active sampling technique can reveal a wealth of information. A ping response indicates that the target host is connected to the network, is reachable from the query agent, and is in a sufficiently functional state to respond to the ping packet. In itself, this response is useful information, indicating that a functional network path to the target host exists. Failure to respond is not so informative because it cannot be reliably inferred that the target host is not available. The ping packet, or perhaps its response, may have been discarded within the network because of transient congestion, or the network may not have a path to the target host, or the network may not have a path back to the ping sending host, or there may be some form of firewall in the end-to-end path that blocks the ICMP packet from being delivered.

However, if you can ping a remote IP address, then you can obtain numerous performance metrics. Beyond simple reachability, further information can be inferred by the ping approach with some basic extensions to our simple ping model. If a sequence of labeled ping packets is generated, the elapsed time for a response to be received for each packet can be recorded, along with the count of dropped packets, duplicated packets, and packets that have been reordered by the network. Careful interpretation of the response times and their variance can provide an indication of the load being experienced on the network path between the query agent and the target. Load will manifest a condition of increased delay and increased variance, due to the interaction of the router buffers with the traffic flows along the path elements as load increases. When a router buffer overflows, the router is forced to discard packets; and under such conditions, increased ping loss is observed. In addition to indications of network load, high erratic delay and loss within a sequence of ping packets may be symptomatic of routing instability with the network path oscillating between many path states.

A typical use of ping is to regularly test numerous paths to establish a baseline of path metrics. This enables a comparison of a specific ping result to these base metrics to give an indication of current path load within the network.

Of course, it is possible to interpret too much from ping results, particularly when pinging routers within a network. Many router architectures use fast switching paths for data packets, whereas the central processing unit of the router may be used to process ping requests. The ping response process may be given a low scheduling priority because router operations represent a more critical router function. It is possible that extended delays and loss, as reported by a ping test, may be related to the processor load or scheduling algorithm of the target router processor rather than to the condition of the network path. (Figure 4)

Figure 4: Ping Path



Ping sequences do not necessarily mimic packet flow behavior of applications. Typical TCP flow behavior is prone to cluster into bursts of packet transmissions on each epoch of the round-trip time. Routers may optimize their cache management, switching behavior, and queue management to take advantage of this behavior. Ping packets may not be clustered; instead, an evenly spaced pacing is used, meaning that the observed metrics of a sequence of ping packets may not exercise such router optimizations. Accordingly, the ping results may not necessarily reflect an anticipation of application performance along the same path. Also a ping test does not measure a simple path between two points. The ping test measures the time to send a packet to a target system and for the target to respond back to the sender. Ping is measuring a loop rather than a simple path.

With these caveats in mind, monitoring a network through regular ping tests along the major network paths can yield useful information regarding the status of the network service performance.

Many refinements to ping can extend its utility. Ping can use *loose source routing* to test the reachability of one host to another, directing the packet from the query host to the loose source routed host, then to the target host and back via the same path through the specified approach. However, many networks disable support for loose source routing, given that it can be exploited in some forms of security attacks. Consequently, the failure of a loose source routed ping may not be a conclusive indication of a network fault.

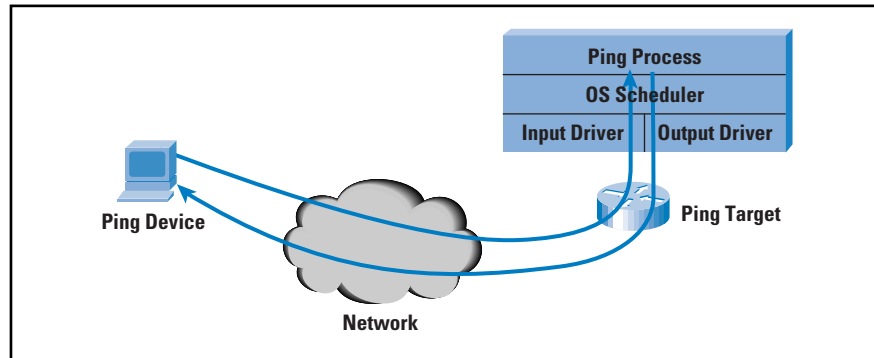
Ping also can be used in a rudimentary way to discover the provisioned capacity of network links. By varying the packet length and comparing the ping times of one router to the next-hop router on a path, the bandwidth of the link can be deduced with some degree of approximation because of a background queue-induced level of network jitter.

A more sophisticated variation of ping is to pace the transmission of packets from the received packets, mimicking the behavior of the TCP flow control algorithms with *Slow Start* and subsequent congestion avoidance. *Treno* is such a tool. In *Treno*, the transmission of ping packets is managed by the TCP Reno flow-control algorithm, such that further ping packets are triggered by the reception of responses to earlier packets, and the triggering of further packets is managed by an implementation of the TCP control function. Such a tool can indicate available flow rate-managed capacity on a chosen path.

Path Discovery Using Traceroute

The second common ICMP-based network management tool, *traceroute*, devised by Van Jacobson, is based on the ICMP *Time Exceeded* message. Here, a sequence of *User Datagram Protocol* (UDP) packets are generated to the target host, each with an increased value of the *Time To Live* (TTL) field in the IP header. This generates a sequence of ICMP Time Exceeded messages sourced from the router where the TTL expired. These source addresses are those of the routers, in turn, on the path from the source to the destination. (Figure 5)

Figure 5: Traceroute Path



Like ping, traceroute measures the elapsed time between the packet transmission and the reception of the corresponding ICMP packet. In this way, the complete output of a traceroute execution exposes not only the elements of the path to the destination, but also the delay and loss characteristics of each partial path element. Traceroute also can be used with loose source route options to uncover the path between two remote hosts. The same caveats mentioned in the ping description relating to the relative paucity in deployment of support for loose source routing apply. An example of a traceroute report is shown in Figure 6.

Figure 6. Traceroute report

```

% traceroute www.cisco.com
traceroute to www.cisco.com (198.133.219.25), 64 hops max, 40 byte packets
 1 dickson-gw1.Canberra.telstra.net (203.50.0.1) 0.272 ms 0.265 ms 0.270 ms
 2 GigabitEthernet4-1.civ12.Canberra.telstra.net (203.50.8.1) 0.402 ms 0.272 ms 0.259 ms
 3 GigabitEthernet3-1.civ-core2.Canberra.telstra.net (203.50.7.5) 0.214 ms 0.227 ms 0.193 ms
 4 GigabitEthernet2-2.dkn-core1.Canberra.telstra.net (203.50.6.126) 0.459 ms 0.394 ms 0.385 ms
 5 Pos4-0.ken-core4.Sydney.telstra.net (203.50.6.121) 3.806 ms 3.762 ms 3.770 ms
 6 Pos2-0.pad-core4.Sydney.telstra.net (203.50.6.22) 3.907 ms 3.959 ms 3.913 ms
 7 GigabitEthernet0-1.syd-core01.Sydney.net.reach.com (203.50.13.246) 3.898 ms 3.866 ms 3.977 ms
 8 i-13-2.sjc-core01.net.reach.com (202.84.143.41) 191.361 ms 191.365 ms 191.341 ms
 9 sl-st21-sj-6-1.sprintlink.net (144.223.242.1) 186.955 ms 186.851 ms 187.010 ms
10 sl-bb25-sj-5-1.sprintlink.net (144.232.20.73) 187.241 ms 187.337 ms 187.055 ms
11 sl-gw11-sj-10-0.sprintlink.net (144.232.3.134) 187.279 ms 186.898 ms 186.821 ms
12 sl-ciscopsn2-11-0-0.sprintlink.net (144.228.44.14) 187.572 ms 187.495 ms 187.620 ms
13 sjck-dirty-gw1.cisco.com (128.107.239.5) 184.533 ms 184.686 ms 184.694 ms
14 sjck-sdf-ci0d-gw1.cisco.com (128.107.239.106) 184.676 ms 184.686 ms 184.644 ms
15 www.cisco.com (198.133.219.25) 185.017 ms 185.122 ms 185.019 ms
  
```

Notes:

- 1) There are interprovider handovers at hops 7, 9, and 13.
- 2) There is a sudden jump in response times at hop 8. The additional 182 ms of round-trip latency corresponds to a 36,000-km submarine cable path. This can be explained by the hop-7 to hop-8 segment, including a submarine cable path between Australia and the United States.

Traceroute is an excellent tool for reporting on the state of the routing system. It operates as an excellent “sanity check” of the match between the design intent of the routing system and the operational behavior of the network.

The caveat to keep in mind when interpreting traceroute output has to do with asymmetric routes within the network. Whereas the per-hop responses expose the routing path taken in the forward direction to the target host, the delay and loss metrics are measured across the forward and reverse paths for each step in the forward path. The reverse path is not explicitly visible to traceroute.

One-Way Measurements

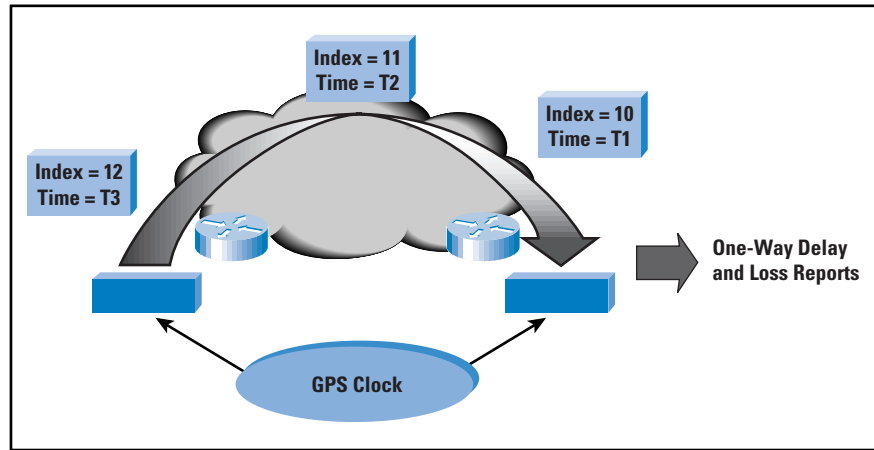
Round-trip probes, such as ping and traceroute, are suited to measuring the total network path between two ends of a transaction, but how can a network provider measure the characteristics of a component of the total end-to-end path? In such a case the network provider is interested in the performance of a set of unidirectional transit paths from an network ingress point to an egress point. There are now some techniques that perform a one-way delay and loss measurement, and they are suited to measuring the service parameters of individual transit paths across a network. A one-way approach does not use a single network management system, but relies on the deployment of probe senders and receivers using synchronized clocks.

The one-way methodology is relatively straightforward. The sender records the precise time a certain bit of the probe packet was transmitted into the network; the receiver records the precise time that same bit arrived at the receiver. Precisely synchronizing the clocks of the two systems is an interesting problem, and initial implementations of this approach have used *Global Positioning System* (GPS) satellite receivers as a synchronized clock source.

One of the noted problems with the use of GPS was that computers are generally located within machine rooms and a clear GPS signal is normally available only on a rooftop. Later implementations of this approach have used the clock associated with the *Code Division Multiple Access* (CDMA) mobile telephone network as a highly accurate, synchronized, distributed clock source, with the advantage that the time signal is usually available close to the measurement unit.

Consequent correlation of the sender’s and receiver’s data from repeated probes can reveal the one-way delay and loss patterns between sender and receiver. To correlate this to a service level requires the packets to travel along the same path as the service flow and with the same scheduling response from the network.

Figure 7: One-Way Measurements



Ping and traceroute are ubiquitous tools. Almost every device can support sending ping and traceroute probes, and, by default almost every device, including network routers, will respond to a ping or traceroute probe. One-way measurements are a different matter, and such measurements normally require the use of dedicated devices in order to undertake the clocking of the probes with the required level of precision (Figure 7).

Choosing the Right Time Base

Whether it is an active or passive measurement regime, the next basic decision is the time base to use for the measurements. Many applications are very sensitive to short-lived transient network conditions. This may take the form of a burst of packet loss, or a period of packet reordering, or a switch to a longer round trip time. TCP may react by halving its sending rate, or by entering an extended wait state while awaiting the retransmission timer to expire. In either case it will take numerous round trip time intervals for the transport session to recover, and this may impact the behavior of the application. On the other hand, a periodic network probe may miss the transient event altogether and report no abnormalities whatsoever.

IP networks have bursty traffic sources, and there is a marked self-similarity in the traffic patterns. This appears to be consistent over a wide range of networks, where large-capacity systems tend to observe large burst patterns and smaller systems also see bursts of a similar proportionate size. So the question is, what time interval for measurements can provide meaningful aggregation of information, while at the same time be sensitive enough to report on the outcomes of transient bursts within the network? Intuitively a measurement time base of hourly measurements is very insensitive to capturing transient bursts, whereas a time base of a millisecond would generate a massive amount of data, a scenario that would tend to smother the identification of abnormalities. Interestingly enough, the choice of a measurement base has little to do with the capacity of the links within a network, but it has a close relationship to the average routing trip time of the individual transport sessions that are active within the network.

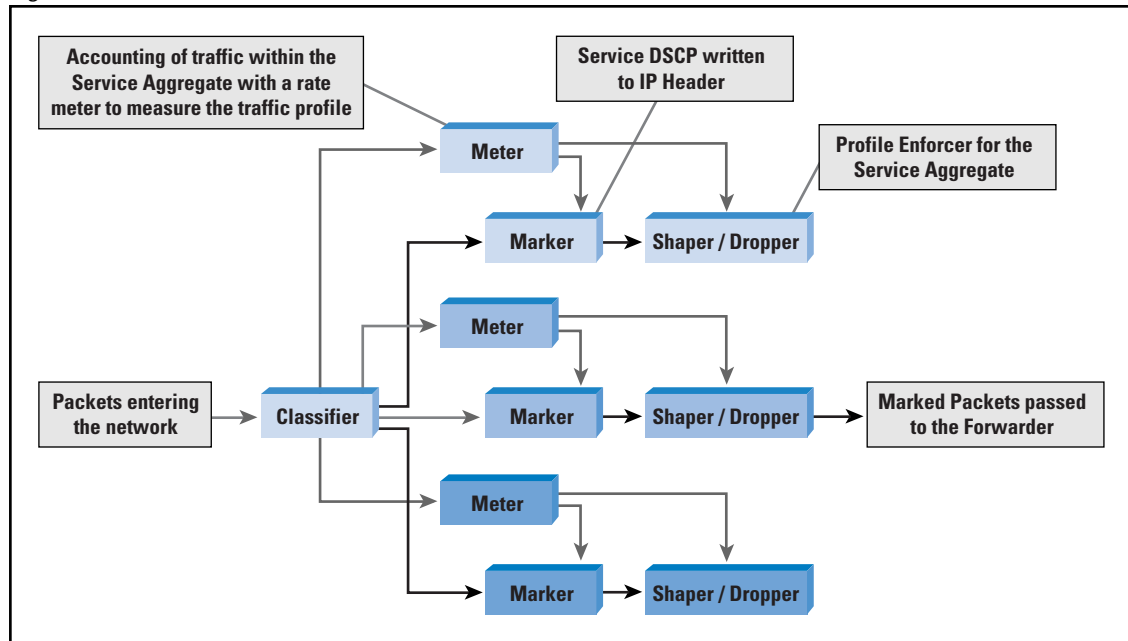
The profile of IP networks is one that is dominated by TCP traffic, and TCP traffic uses a transport control mechanism where the returning stream of *acknowledgement* (ACK) packets governs the actions of the sender. This implies that network-based distortion in the forward data path will not be signaled back to the sender for one complete round-trip time interval, and the consequent adaptation of the sender to the conditions of the network will take numerous additional round-trip times. The implication is that in order to capture a comprehensive view of network performance, a time base of 1 to 2 seconds is appropriate. However, for large networks, such a view generates a massive amount of data. It appears that many networks use a measurement time base of about 60 to 300 seconds, representing an acceptable compromise between sensitivity of the measurement system and the consequent volume of measurement data to analyze.

What About QoS Networks?

So far the assumption has been that the network operates with a single service level and that probes of the network operate at the same service level as the network payload. This is certainly a common situation, but the total picture is slightly broader. When the network provider attempts to create a premium response for certain classes of traffic, and where the customer is paying a premium tariff to use such a premium service, the question of performance becomes a matter of deep concern to both the provider and the customer. After all, the customer is now paying a premium for improved performance, so it would help all concerned if this could be clearly defined and measured.

Solutions exist in both the passive and active polling domains. In the case of SNMP there is a monitoring framework (or Management Information Base, MIB) relating to the *Differentiated Services* (DiffServ) model of *Quality of Service* (QoS), and also MIBs relating to the *Integrated Services* (IntServ) QoS model. For the DiffServ MIB, it is first necessary to define an abstract model of the operation of a DiffServ admission router, by looking at the major functional blocks of the router. The first of these blocks is the definition of the supported behavior aggregates provided by the network. Within the network path, the initial active path element is the traffic classification module, which can be modeled as a set of filters and an associated set of output streams. The output stream is passed to the traffic-conditioning elements, which are the traffic meters and the associated action elements. Many meter profiles can be used in the model: an average data rate, an exponential weighted moving average of one of numerous various traffic profiles that can be expressed by a set of token-bucket parameters using an average rate, a peak rate, and a burst size. More elaborate meter specifications can be constructed using a multilevel token-bucket specification. From the meter, the traffic is passed through an action filter, which may mark the packets and shape the traffic profile through queues or discard operations. Together, this sequence of components forms a *traffic conditioning block*. The traffic is then passed into a queue through the use of a queuing discipline that applies the desired service behavior. (Figure 8)

Figure 8: DiffServ Control Architecture



From this generic model it is possible to define instrumentation for SNMP polling, where each of these five components—the behavior aggregate, the classifier, the meter, profile actions, and the queuing discipline—correspond to a MIB table. With this structure it is possible to parameterize both the specific configuration of the DiffServ network element and its dynamic state. This MIB is intended to describe the configuration and operation of both edge and interior DiffServ network elements, the difference being that interior elements use just a behavior aggregate classifier and a queue manager within the management model, whereas the edge elements use all components of the model.

A comparable MIB is defined for the IntServ architecture and an additional MIB for the operation of guaranteed services. The IntServ MIB defines the per-element reservation table used to determine the current reservation state, an indication of whether or not the router can accept further flow reservations, and the reservation characteristics of each current flow. No performance polling parameters or accounting parameters are included in the MIB. The guaranteed services MIB adds to this definition with a per-interface definition of a backlog. This is a means of expressing *packet quantization delay*, a delay term, which is the packet propagation delay over the interface, and a slack term, which is the amount of slack in the reservation that can be used without redefining the reservation. Again, these are per-element status definitions, and they do not include performance or accounting data items.

The IntServ MIB is being further defined as a *Resource Reservation Protocol (RSVP) MIB* for the operation of IntServ network elements^[14]. There are a larger number of objects within the MIB, including General Objects, Session Statistics Table, Session Sender Table, Reservation Requests Received Table, Reservation Requests Forwarded Table, RSVP Interface Attributes Table, and an RSVP Neighbor Table.

Interestingly, the MIB proposes a writeable RSVP reservation table to allow the network manager to manually create a reservation state that can be removed only through a comparable manual operation. The MIB enables a management system to poll the IntServ network element to retrieve the status of every active IntServ reserved flow and the operational characteristics of the flow, as seen by the network element.

In a QoS DiffServ environment, ping and traceroute pose some interesting engineering issues. Ping sends an ICMP packet. The network QoS admission filters may choose a different classification for these packets from that chosen for normal data-flow TCP or UDP protocol packets; as a result, the probe packet may be scheduled differently or even take a completely different path to the network. In an IntServ QoS network, the common classification condition for a flow is a combination of the IP header source and destination addresses and the TCP or UDP header source and destination port addresses. The ping probe packet cannot reproduce this complete flow description, and therefore cannot, by default, be inserted into the flow path that it is attempting to measure. With traceroute, the packet does have a UDP protocol address, but it uses a constant port address by default, causing a similar problem of attempting to be inserted to an IntServ flow. DiffServ encounters similar problems when attempting to pass the probe packet into the network via the DiffServ admission classification systems. Inside the network, it is possible to insert the probe packet into the network with the IP *Differentiated Services Code Point* (DSCP) field set to the DiffServ behavior aggregate that is being measured.

The measurement of delay and loss taken by ping and traceroute is a cumulative value of both the forward and return path delay and loss. When attempting to measure unidirectional flow-path behavior, such as an IntServ flow path, this measurement is of dubious value, given the level of uncertainty as to which part of the path, forward or reverse, contributed to the ping or traceroute delay and loss reports.

For one-way delay measurements, in DiffServ networks, this can be done within the network, setting the DSCP field to the value of the service aggregate being monitored. Of course, from the customer's perspective, the DiffServ network service profile includes the admission traffic-conditioning block, and the interior one-way measurements are only part of the delivered service. In the IntServ network, the packets have to be structured to take the same path as the elevated service flows; they are classified by each element as part of the collection of such elevated service flows for the purposes of scheduling.

Measuring Performance—The Client Perspective

From the client's perspective, the measurement choices are more limited. A client does not normally enjoy the ability to poll network elements within a provider's network. One way for a client to measure service quality is to instigate probing of the network path, whereby a sender can pass a probe packet into the network and measure the characteristics of the response. Of course, the problems of inserting probe packets into the service flow remain, as do the issues of unidirectional elevated service flows with bidirectional probes.

However, the client does have the advantage of being able to monitor and manipulate the characteristics of the service flow itself. For TCP sessions, the client can monitor the packet retransmission rate, the maximum burst capacity, the average throughput, the *round-trip time* (RTT), RTT variance, and misordered packets, by monitoring the state of the outbound data flow and relating it to the inbound ACK flow. For UDP sessions, there is no corresponding transport-level feedback information flow to the sender as a part of the transport protocol itself. The receiver can measure the service quality of the received datastream using information provided in the *Real-Time Protocol* (RTP) information feedback fields—if RTP is being used for real-time data or as an application-related tool for other application types. If sender and receiver work in concert, the receiver can generate periodic quality reports and pass these summaries back to the sender. Such applications can confirm whether an application is receiving a specified level of service. This approach treats the network like a black box; no attempt is made to identify the precise nature or source of events that disrupt the delivered service quality. There are no standardized approaches to this activity, but numerous analysis tools are available for host platforms that perform these measurements.

Though the client can measure and conform service quality on a per-application level of granularity, the second part of the client's motivation in measuring service quality is more difficult to address. The basic question is whether the service delivered in response to a premium service request is sufficiently differentiated from a best-effort service transaction. Without necessarily conducting the transaction a second time, the best approach is to use either one-way delay probes, for unidirectional traffic, or a bulk TCP capacity probe, to establish some indication of the relativity in performance. From a client perspective none of these are simple to set up, and the dilemma that the customer often faces is the basic question of whether the cost of operating the measurement setup is adequately offset by the value of the resulting answers.

Measuring Networks—Looking for Problems

So far we have been looking at the ways of measuring network performance as a general task. Of course degraded performance does not happen by accident (well, sometimes accidents do happen), and it makes the measurement task easier if you can identify precisely what it is that you are looking for. This approach requires identification of the various situations that can impact network performance and then set up network measurement and monitoring systems that are tuned to identify these situations.

Within this approach, the motives for network measurement are concerned with identification of traffic load patterns that cause uneven network load, monitoring, and verification of service-level agreements, detection of abnormal network load that may be a signature of an attack, forecasting and capacity planning, and routing stability.

The objective here is to create a stable and well-understood model of the operational characteristics of the network, and then analyze the situations that could disrupt this stable state and the implications in terms of delivered performance under such conditions.

Such an approach could be described in terms of opposites—instead of measuring network performance, the approach is measuring the network to identify the conditions that cause nonperformance at particular times within particular network paths. As a performance management technique, this approach has been very effective—rather than taking a larger amount of performance data and merging and averaging it into a relatively meaningless index, the approach is to isolate those circumstances where performance is compromised and report on these exceptions rather than on the remainder of the time.

Of course measuring what is “normal” may involve more than assembling a benchmark set of SNMP-derived polling data and a collection of latency, loss, and jitter profiles obtained from analysis of large volumes of ping data. One additional tool is the router itself. Because the router uses many IP packet header fields to switch each packet, one approach is to get the router to assemble and aggregate information about the characteristics of traffic that has been passed through the router, and send these aggregated reports to a network management station for further analysis. *NetFlow* is the most common tool to undertake this form of reporting. Like SNMP, NetFlow can report on the characteristics of traffic as it passes a point in the network. For measuring end-to-end performance of individual applications, NetFlow has the same limitations as SNMP. The analogy is one of standing on a street corner counting cars that go past and from that measurement attempting to derive the average time for a commuter to drive to or from work. However, the value of NetFlow is that in this context of performance measurement, it can be used to derive a picture of the baseline characteristics of the network, including identification of the endpoints of the traffic flows. Extending the car analogy further, NetFlow can provide an indication of the origins and ultimate destinations of the cars as they pass the monitoring point. This information is useful in terms of designing networks that are adequately configured to handle the transit traffic load. In addition, with careful analysis, NetFlow can be used to identify exceptional traffic conditions. The advantage here is that NetFlow data can be used to identify both the abnormal traffic load and also provide some indication of the endpoints of the abnormal flows. In this way, NetFlow can be deployed as both a baseline network traffic profile benchmarking tool and a performance exception diagnosis tool.

This approach of capturing the packet header information as the traffic passes a monitoring point in the network has been implemented in numerous ways, and NetFlow is not the only data-collection tool in this space. One interesting approach has been used by NeTraMet, an implementation of the *Internet Engineering Task Force's* (IETF's) *Realtime Traffic Flow Measurement* architecture for traffic flow measurement.

The feature here is a powerful ruleset within the tool that allows the flow collector to be configured to collect information about particular traffic flows and their characteristics. In the context of measuring performance, one of the abilities of the tool is to match the outbound data flow with the inbound acknowledgement stream, allowing an analyzer some ability to infer end-to-end performance of the application based on the collected information.

Where to Go from Here

It is clear that the picture is so far very incomplete. The active probe measurements require either some latitude of interpretation or dedicated instrumentation to take measurements with some necessary level of frequency and precision. The passive approach of probing the active switching elements of the network is constrained by a very basic model of the switching system, so that the collectable values provide only a very indirect relationship to the manner in which the switching element is generating queuing delays and traffic flow instability.

Perhaps what is also increasingly unclear is the relationship between performance and networks in any case. The last few years have seen a massive swing in public Internet platforms away from networks where some level of congestion and contention was anticipated to networks that are extensively overprovisioned, and where packet jitter and loss are simply not encountered. With the ever-decreasing cost of transmission bandwidth in many markets, this environment of abundant network capacity is now also finding its way into various enterprise network sectors. In such worlds of abundant supply and overengineering of networks, there is really little left to measure within the network. The entire question of performance then becomes a question phrased much closer to home: how well is your system tuned to make the most of its resources and those of the server? Often the entire issue with performance is a situation of abundant network resources, abundant local memory and processing resources, and poor tuning of the transport protocol stack. That is, of course, quite properly the subject of another article.

Further Reading

The Internet offers a wealth of material on the topic of network measurement, and the major exercise is undertaking some filtering to get a broad collection of material that encompasses a range of perspectives on this topic. The following sources were used to prepare this article, and are recommended as starting points for further exploration of this topic.

- [1] *Internet Performance Survival Guide*, Geoff Huston, Wiley Computer Publishing, 2000.
- [2] "IPPM Metrics for Measuring Connectivity," J. Mahdavi, V. Paxson, RFC 2678, September 1999.
- [3] "A One-way Delay Metric for IPPM," G. Almes, S. Kalidinki, M. Zeukuaskas, RFC 2679, September 1999.

- [4] "A One-way Packet Loss Metric for IPPM," G. Almes, S. Kalidinki, M. Zeukuaskas, RFC 2680, September 1999.
 - [5] The RIPE Test Traffic Measurement service at:
<http://www.ripe.net/ripence/mem-services/ttm/>
 - [6] Treno, online at:
http://www.psc.edu/networking/treno_info.html
 - [7] "Trends in Measurement and Monitoring of Internet Backbones," session at the 26th North American Network Operators Group, hosted by D. Meyer,
<http://www.nanog.org/mtg-0210/measurement.html>,
October 2002.
 - [8] "Some thoughts on CoS and Backbone Networks," D. Meyer, presentation to the IEPREP Working Group, IETF-55,
<http://www.maoz.com/~dmm/IETF55/ieprep/>, November 2002.
 - [9] NetFlow resource page:
http://www.cisco.com/warp/public/732/Tech/nmp/netflow/netflow_techdoc.shtml
 - [10] Netramet, and many other interesting measurement tools are referenced in a resource page at: <http://www.caida.org/tools>
- This area of research is active, and numerous activities are ongoing in the area of research group activities and workshops.
- [11] The Internet Research Task Force has an Internet Measurement Research Group. Further details can be found at:
<http://www.irtf.org/charters/imrg.html>
 - [12] ACM SIGCOMM, the ACM Special Interest Group on Data Communications, sponsors an Internet Measurement Workshop. Proceeding of the November 2002 workshop can be found at:
<http://www.acm.org/sigcomm/imw2002/>
 - [13] The details of the 2003 Passive and Active Measurement Workshop can be found at: <http://www.pam2003.org>
 - [14] "RSVP Management Information Base using SMIv2," F. Baker, J. Krawczyk, A. Sastry, RFC 2206, September 1997.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. Huston is currently the Chief Scientist in the Internet area for Telstra. He is also a member of the Internet Architecture Board, and is the Secretary of the APNIC Executive Committee. He is author of *The ISP Survival Guide*, ISBN 0-471-31499-4, *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, and coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. All three books are published by John Wiley & Sons. E-mail: gih@telstra.net

The Session Initiation Protocol

by William Stallings

The *Session Initiation Protocol* (SIP), defined in RFC 3261^[6], is an application level signaling protocol for setting up, modifying, and terminating real-time sessions between participants over an IP data network. SIP can support any type of single-media or multi-media session, including teleconferencing.

SIP is just one component in the set of protocols and services needed to support multimedia exchanges over the Internet. SIP is the signaling protocol that enables one party to place a call to another party and to negotiate the parameters of a multimedia session. The actual audio, video, or other multimedia content is exchanged between session participants using an appropriate transport protocol. In many cases, the transport protocol to use is the *Real-Time Transport Protocol* (RTP). Directory access and lookup protocols are also needed.

The key driving force behind SIP is to enable Internet telephony, also referred to as *Voice over IP* (VoIP). There is wide industry acceptance that SIP will be the standard IP signaling mechanism for voice and multimedia calling services. Further, as older *Private Branch Exchanges* (PBXs) and network switches are phased out, industry is moving toward a voice networking model that is SIP signaled, IP based, and packet switched, not only in the wide area but also on the customer premises^[2, 3].

SIP supports five facets of establishing and terminating multimedia communications:

- *User location*: Users can move to other locations and access their telephony or other application features from remote locations.
- *User availability*: This step involves determination of the willingness of the called party to engage in communications.
- *User capabilities*: In this step, the media and media parameters to be used are determined.
- *Session setup*: Point-to-point and multiparty calls are set up, with agreed session parameters.
- *Session management*: This step includes transfer and termination of sessions, modifying session parameters, and invoking services.

SIP employs design elements developed for earlier protocols. SIP is based on an HTTP-like request/response transaction model. Each transaction consists of a client request that invokes a particular method, or function, on the server and at least one response. SIP uses most of the header fields, encoding rules, and status codes of HTTP. This provides a readable text-based format for displaying information. SIP incorporates the use of a *Session Description Protocol* (SDP), which defines session content using a set of types similar to those used in *Multipurpose Internet Mail Extensions* (MIME).

SIP Components and Protocols

A system using SIP can be viewed as consisting of components defined on two dimensions: client/server and individual network elements. RFC 3261 defines client and server as follows:

- *Client*: A client is any network element that sends SIP requests and receives SIP responses. Clients may or may not interact directly with a human user. User agent clients and proxies are clients.
- *Server*: A server is a network element that receives requests in order to service them and sends back responses to those requests. Examples of servers are proxies, user agent servers, redirect servers, and registrars.

The individual elements of a standard SIP configuration include the following:

- *User Agent*: The user agent resides in every SIP end station. It acts in two roles:
 - User Agent Client (UAC): Issues SIP requests
 - User Agent Server (UAS): Receives SIP requests and generates a response that accepts, rejects, or redirects the request
- *Redirect Server*: The redirect server is used during session initiation to determine the address of the called device. The redirect server returns this information to the calling device, directing the UAC to contact an alternate *Universal Resource Identifier* (URI). A URI is a generic identifier used to name any resource on the Internet. The URL used for Web addresses is a type of URI. See RFC 2396^[1] for more detail.
- *Proxy Server*: The proxy server is an intermediary entity that acts as both a server and a client for the purpose of making requests on behalf of other clients. A proxy server primarily plays the role of routing, meaning that its job is to ensure that a request is sent to another entity closer to the targeted user. Proxies are also useful for enforcing policy (for example, making sure a user is allowed to make a call). A proxy interprets, and, if necessary, rewrites specific parts of a request message before forwarding it.
- *Registrar*: A registrar is a server that accepts REGISTER requests and places the information it receives (the SIP address and associated IP address of the registering device) in those requests into the location service for the domain it handles.
- *Location Service*: A location service is used by a SIP redirect or proxy server to obtain information about a callee's possible location(s). For this purpose, the location service maintains a database of SIP-address/IP-address mappings.

The various servers are defined in RFC 3261 as logical devices. They may be implemented as separate servers configured on the Internet or they may be combined into a single application that resides in a physical server.

Figure 1: SIP
Components and
Protocols

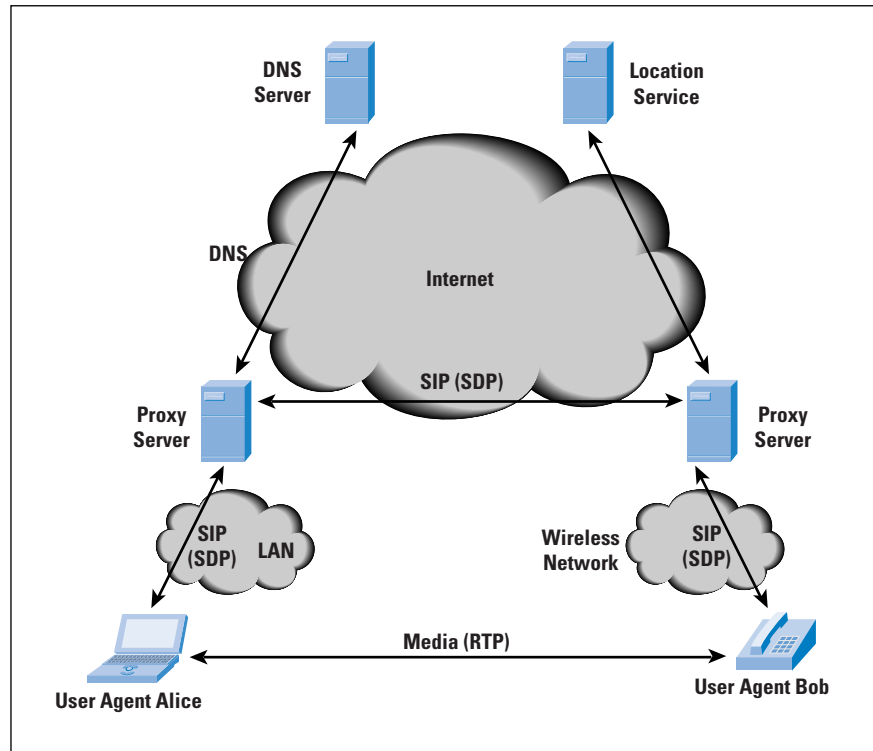


Figure 1 shows how some of the SIP components relate to one another and the protocols that are employed. A user agent acting as a client (in this case UAC Alice) uses SIP to set up a session with a user agent that acts as a server (in this case UAS Bob). The session initiation dialogue uses SIP and involves one or more proxy servers to forward requests and responses between the two user agents. The user agents also make use of the SDP, which is used to describe the media session.

The proxy servers may also act as redirect servers as needed. If redirection is done, a proxy server needs to consult the location service database, which may or may not be colocated with a proxy server. The communication between the proxy server and the location service is beyond the scope of the SIP standard. The *Domain Name System* (DNS) is also an important part of SIP operation. Typically, a UAC makes a request using the domain name of the UAS, rather than an IP address. A proxy server needs to consult a DNS server to find a proxy server for the target domain.

SIP often runs on top of the *User Datagram Protocol* (UDP) for performance reasons, and provides its own reliability mechanisms, but may also use TCP. If a secure, encrypted transport mechanism is desired, SIP messages may alternatively be carried over the *Transport Layer Security* (TLS) protocol.

Associated with SIP is the SDP, defined in RFC 2327^[4]. SIP is used to invite one or more participants to a session, while the SDP-encoded body of the SIP message contains information about what media encodings (for example, voice, video) the parties can and will use. After this information is exchanged and acknowledged, all participants are aware of the participants' IP addresses, available transmission capacity, and media type. Then, data transmission begins, using an appropriate transport protocol. Typically, the RTP is used. Throughout the session, participants can make changes to session parameters, such as new media types or new parties to the session, using SIP messages.

SIP Universal Resource Indicators

A resource within a SIP configuration is identified by a URI. Examples of communications resources include the following:

- A user of an online service
- An appearance on a multiline phone
- A mailbox on a messaging system
- A telephone number at a gateway service
- A group (such as “sales” or “help desk”) in an organization

SIP URIs have a format based on e-mail address formats, namely **user@domain**. There are two common schemes. An ordinary SIP URI is of the form:

sip:bob@biloxi.com

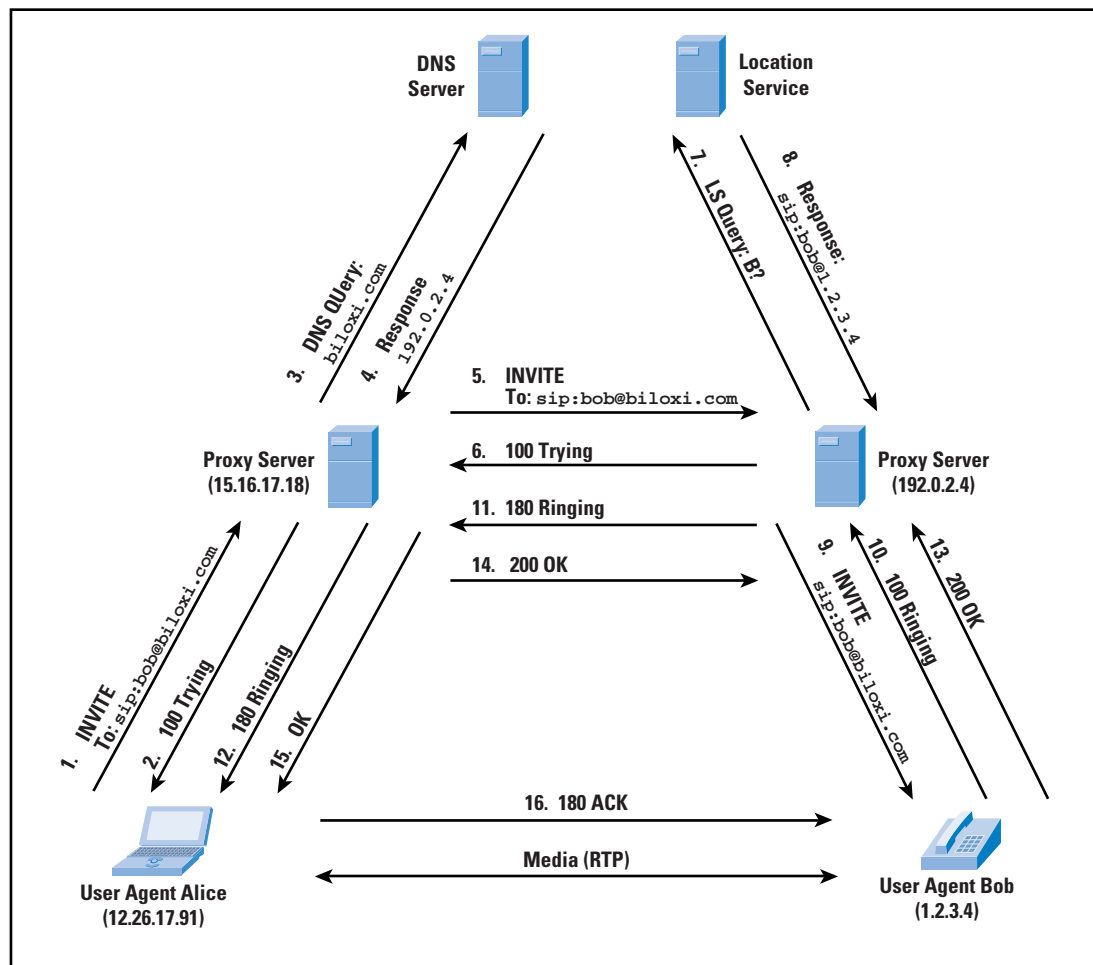
The URI may also include a password, port number, and related parameters. If secure transmission is required, “**sip:**” is replaced by “**sips:.**” In the latter case, SIP messages are transported over TLS.

Examples of Operation

The SIP specification is quite complex; the main document, RFC 3261, is 269 pages long. To give some feel for its operation, we present a few examples.

Figure 2 shows a successful attempt by user Alice to establish a session with user Bob, whose URI is **bob@biloxi.com**.^[9] Alice's UAC is configured to communicate with a proxy server (the outbound server) in its domain and begins by sending an INVITE message to the proxy server that indicates its desire to invite Bob's UAS into a session (1); the server acknowledges the request (2). Although Bob's UAS is identified by its URI, the outbound proxy server needs to account for the possibility that Bob is not currently available or that Bob has moved. Accordingly, the outbound proxy server should forward the INVITE request to the proxy server that is responsible for the domain **biloxi.com**. The outbound proxy thus consults a local DNS server to obtain the IP address of the **biloxi.com** proxy server (3), by asking for the DNS SRV resource record that contains information on the proxy server for **biloxi.com**.

Figure 2: SIP Successful Call Setup

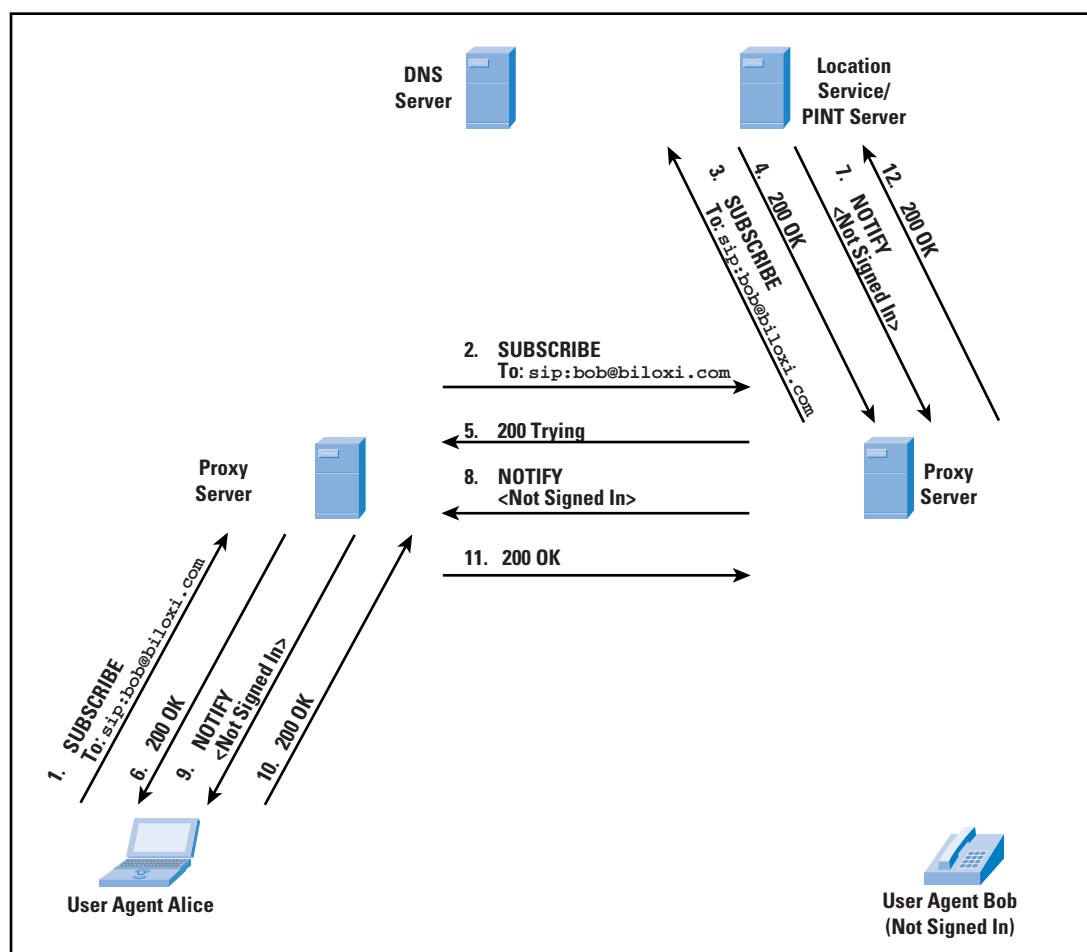


The DNS server responds (4) with the IP address of the **biloxi.com** proxy server (the inbound server). Alice's proxy server can now forward the INVITE message to the inbound proxy server (5), which acknowledges the message (6). The inbound proxy server now consults a location server to determine Bob's location (7), and the location server responds with Bob's location, indicating that Bob is signed in, and therefore available for SIP messages (8).

The proxy server can now send the INVITE message on to Bob (9). A ringing response is sent from Bob back to Alice (10, 11, 12) while the UAS at Bob is alerting the local media application (for example, telephony). When the media application accepts the call, Bob's UAS sends back an OK response to Alice (13, 14, 15).

Finally, Alice's UAC sends an acknowledgement message to Bob's UAS to confirm the reception of the final response (16). In this example, the ACK is sent directly from Alice to Bob, bypassing the two proxies. This occurs because the endpoints have learned each other's address from the INVITE/200 (OK) exchange, which was not known when the initial INVITE was sent. The media session has now begun, and Alice and Bob can exchange data over one or more RTP connections.

Figure 3: SIP Presence Example

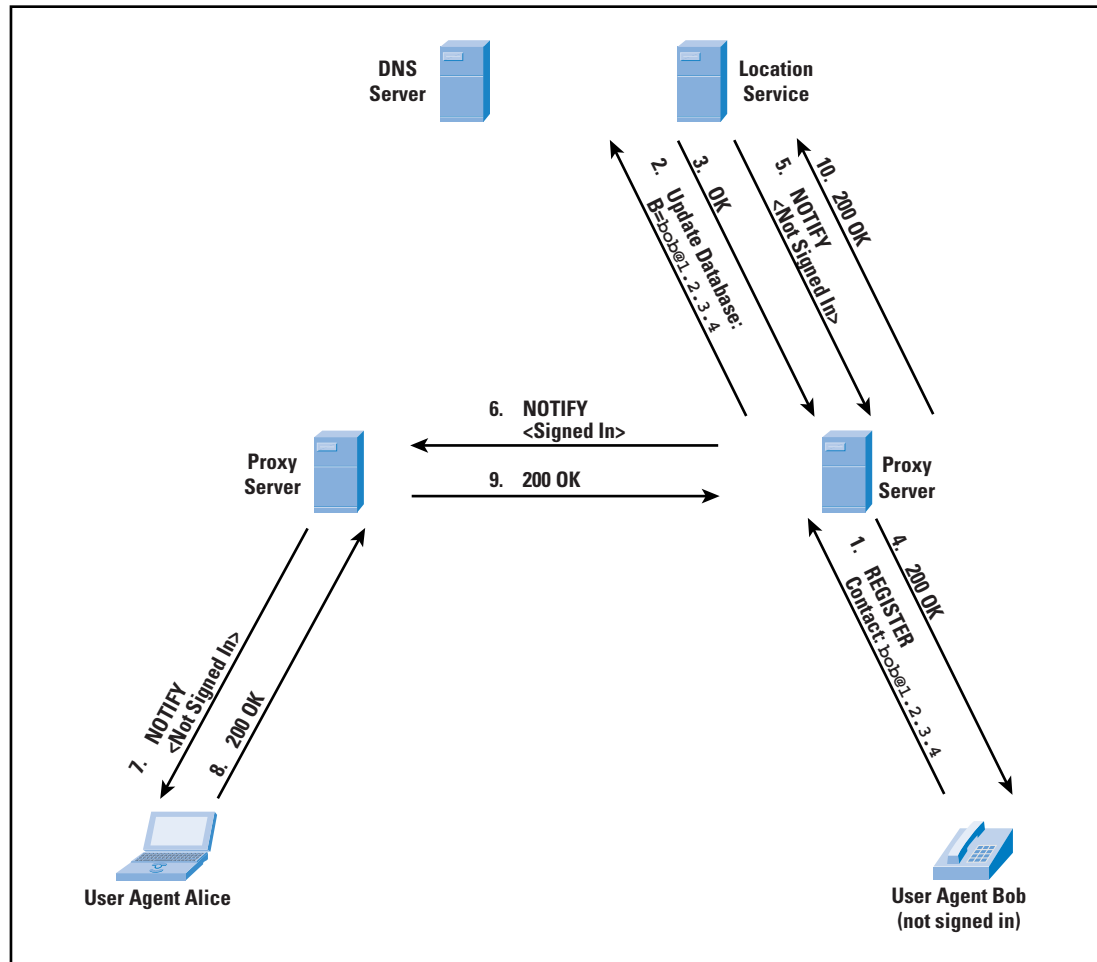


The next example (Figure 3) makes use of two message types that are not yet part of the SIP standard but that are documented in RFC 2848^[5] and are likely to be incorporated in a later revision of SIP. These message types support telephony applications. Suppose that in the preceding example, Alice was informed that Bob was not available. Alice's UAC can then issue a SUBSCRIBE message (1), indicating that it wants to be informed when Bob is available.

This request is forwarded through the two proxies in our example to a PINT (*Public Switched Telephone Network [PSTN]-Internet Networking*) server (2, 3). A PINT server acts as a gateway between an IP network from which comes a request to place a telephone call and a telephone network that executes the call by connecting to the destination telephone. In this example, we assume that the PINT server logic is colocated with the location service. It could also be the case that Bob is attached to the Internet rather than a PSTN, in which case the equivalent of PINT logic is needed to handle SUBSCRIBE requests. In this example, we assume the latter and assume that the PINT functionality is implemented in the location service. In any case, the location service authorizes subscription by returning an OK message (4), which is passed back to Alice (5, 6). The location service then immediately sends a NOTIFY message with Bob's current status of not signed in (7, 8, 9), which Alice's UAC acknowledges (10, 11, 12).

Figure 4 continues the example of Figure 3. Bob signs on by sending a REGISTER message to the proxy in its domain (1). The proxy updates the database at the location service to reflect registration (2). The update is confirmed to the proxy (3), which confirms the registration to Bob (4). The PINT functionality learns of Bob's new status from the location server (here we assume that they are colocated) and sends a NOTIFY message containing Bob's new status (5), which is forwarded to Alice (6, 7). Alice's UAC acknowledges receipt of the notification (8, 9, 10).

Figure 4: SIP Registration and Notification Example



SIP Messages

As was mentioned, SIP is a text-based protocol with a syntax similar to that of HTTP. There are two different types of SIP messages, *requests* and *responses*. The format difference between the two types of messages is seen in the first line. The first line of a request has a method, defining the nature of the request and a Request-URI, indicating where the request should be sent. The first line of a response has a response code. All messages include a header, consisting of a number of lines, each line beginning with a header label. A message can also contain a body such as an SDP media description.

For SIP requests, RFC 3261 defines the following methods:

- *REGISTER*: Used by a user agent to notify a SIP configuration of its current IP address and the URLs for which it would like to receive calls
- *INVITE*: Used to establish a media session between user agents
- *ACK*: Confirms reliable message exchanges
- *CANCEL*: Terminates a pending request, but does not undo a completed call
- *BYE*: Terminates a session between two users in a conference
- *OPTIONS*: Solicits information about the capabilities of the callee, but does not set up a call

For example, the header of message (1) in Figure 2 might look like the following:

```
INVITE sip:bob@biloxi.com SIP/2.0
Via: SIP/2.0/UDP 12.26.17.91:5060
Max-Forwards: 70
To: Bob <sip:bob@biloxi.com>
From: Alice <sip:alice@atlanta.com;tag=1928301774>
Call-ID: a84b4c76e66710@12.26.17.91
CSeq: 314159 INVITE
Contact: <sip:alice@atlanta.com>
Content-Type: application/sdp
Content-Length: 142
```

The first line contains the method name (**INVITE**), a SIP URI, and the version number of SIP that is used. The lines that follow are a list of header fields. This example contains the minimum required set.

The **via** headers show the path the request has taken in the SIP configuration (source and intervening proxies), and are used to route responses back along the same path. As the INVITE message leaves, there is only the header inserted by Alice. The line contains the IP address (**12.26.17.91**), port number (**5060**), and transport protocol (**UDP**) that Alice wants Bob to use in his response.

The **Max-Forwards** header limits the number of hops a request can make on the way to its destination. It consists of an integer that is decremented by one by each proxy that forwards the request. If the **Max-Forwards** value reaches 0 before the request reaches its destination, it is rejected with a 483 (**Too Many Hops**) error response.

The **To** header field contains a display name (Bob) and a SIP or SIPS URI (**sip:bob@biloxi.com**) toward which the request was originally directed. The **From** header field also contains a display name (Alice) and a SIP or SIPS URI (**sip:alice@atlanta.com**) that indicate the originator of the request. This header field also has a **tag** parameter that contains a random string (**1928301774**) that was added to the URI by the UAC. It is used to identify the session.

The **Call-ID** header field contains a globally unique identifier for this call, generated by the combination of a random string and the host name or IP address. The combination of the **To** tag, **From** tag, and **Call-ID** completely defines a peer-to-peer SIP relationship between Alice and Bob and is referred to as a dialog.

The **CSeq** or *Command Sequence* header field contains an integer and a method name. The CSeq number is initialized at the start of a call (314159 in this example), incremented for each new request within a dialog, and is a traditional sequence number. The CSeq is used to distinguish a retransmission from a new request.

The **Contact** header field contains a SIP URI for direct communication between user agents. Whereas the **Via** header field tells other elements where to send the response, the **Contact** header field tells other elements where to send future requests for this dialog.

The **Content-Type** header field indicates the type of the message body. The **Content-Length** header field gives the length in octets of the message body.

The SIP response types defined in RFC 3261 are in the following categories:

- *Provisional* (1xx): The request was received and is being processed.
- *Success* (2xx): The action was successfully received, understood, and accepted.
- *Redirection* (3xx): Further action needs to be taken in order to complete the request.
- *Client Error* (4xx): The request contains bad syntax or cannot be fulfilled at this server.
- *Server Error* (5xx): The server failed to fulfill an apparently valid request.
- *Global Failure* (6xx): The request cannot be fulfilled at any server.

For example, the header of message (13) in Figure 2 might look like the following:

```
SIP/2.0 200 OK
Via: SIP/2.0/UDP server10.biloxi.com
Via: SIP/2.0/UDP bigbox3.site3.atlanta.com
Via: SIP/2.0/UDP 12.26.17.91:5060
To: Bob <sip:bob@biloxi.com;tag=a6c85cf>
From: Alice <sip:alice@atlanta.com;tag=1928301774>
Call-ID: a84b4c76e66710@12.26.17.91
CSeq: 314159 INVITE
Contact: <sip:bob@biloxi.com>
Content-Type: application/sdp
Content-Length: 131
```

The first line contains the version number of SIP that is used and the response code and name. The lines that follow are a list of header fields. The **Via**, **To**, **From**, **Call-ID**, and **CSeq** header fields are copied from the INVITE request. (There are three **via** header field values—one added by Alice’s SIP UAC, one added by the **atlanta.com** proxy, and one added by the **biloxi.com** proxy.) Bob’s SIP phone has added a **tag** parameter to the **To** header field. This tag is incorporated by both endpoints into the dialog and is included in all future requests and responses in this call.

Session Description Protocol

The *Session Description Protocol* (SDP), defined in RFC 2327, describes the content of sessions, including telephony, Internet radio, and multimedia applications. SDP includes information about^[8]:

- *Media streams*: A session can include multiple streams of differing content. SDP currently defines audio, video, data, control, and application as stream types, similar to the MIME types used for Internet mail.
- *Addresses*: SDP indicates the destination addresses, which may be a multicast address, for a media stream.
- *Ports*: For each stream, the UDP port numbers for sending and receiving are specified.
- *Payload types*: For each media stream type in use (for example, telephony), the payload type indicates the media formats that can be used during the session.
- *Start and stop times*: These apply to broadcast sessions, for example, a television or radio program. The start, stop, and repeat times of the session are indicated.
- *Originator*: For broadcast sessions, the originator is specified, with contact information. This may be useful if a receiver encounters technical difficulties.

Although SDP provides the capability to describe multimedia content, it lacks the mechanisms by which two parties agree on the parameters to be used. RFC 3264^[7] remedies this lack by defining a simple offer/answer model, by which two parties exchange SDP messages to reach agreement on the nature of the multimedia content to be transmitted.

References

- [1] T. Berners-Lee, R. Fielding, and L. Masinter, “Uniform Resource Identifiers (URI): Generic Syntax,” RFC 2396, August 1998.
- [2] S. Borthick, “SIP Services: Slowly Rolling Forward,” *Business Communications Review*, June 2002.
- [3] S. Borthick, “SIP for the Enterprise: Work in Progress,” *Business Communications Review*, February 2003.

- [4] M. Handley and V. Jacobson, "SDP: Session Description Protocol," RFC 2327, April 1998.
- [5] S. Petrack and L. Conroy, "The PINT Service Protocol: Extensions to SIP and SDP for IP Access to Telephone Call Services," RFC 2848, June 2000.
- [6] J. Rosenberg, H. Schulzrinne, G. Camarillo, A. Johnston, J. Peterson, R. Sparks, M. Handley, and E. Schooler, "SIP: Session Initiation Protocol," RFC 3261, June 2002.
- [7] J. Rosenberg and H. Schulzrinne, "An Offer/Answer Model with the Session Description Protocol," RFC 3264, June 2002.
- [8] H. Schulzrinne and J. Rosenberg, "The Session Initiation Protocol: Providing Advanced Telephony Access Across the Internet," *Bell Labs Technical Journal*, October-December 1998.
- [9] Figures 2 through 4 are adapted from ones developed by Professor H. Charles Baker of Southern Methodist University.

WILLIAM STALLINGS is a consultant, lecturer, and author of over a dozen books on data communications and computer networking. He also maintains a computer science resource site for CS students and professionals at WilliamStallings.com/StudentSupport.html. He has a PhD in computer science from M.I.T. His latest book is *Computer Networks, with Internet Protocols and Technology* (Prentice Hall, 2003). His home in cyberspace is WilliamStallings.com and he can be reached at ws@shore.net

Letters to the Editor

Ruling the Root Ole,

As a matter of principle, I don't mind having my book *Ruling the Root* reviewed by David Crocker. Mr. Crocker was a significant figure in some of the key events covered in the book. His assessment and opinion of the book had the potential to be quite interesting.

One can only be disappointed with the results, however. The review reveals an inability to rise above partisan sniping and engage rationally with an different view. That, as a matter of policy, is why serious journals don't publish unsolicited reviews of books. Unsolicited reviewers tend to fall into one of two types: unabashed promoters with a personal interest in the success of the book, or people with an axe to grind trying to shoot down a perceived enemy.

I offer a rebuttal only because I think it is vital that the Internet technical community, the presumed readers of *The Internet Protocol Journal*, achieve a higher standard in their discussion of Internet-related policy issues.

Ruling the Root is a serious attempt to analyze the intersection of technology and policy. It offers a way of understanding that intersection based on theories of institutions and property rights. I know that this intersection irritates many engineers, who often harbor a wish that it would go away. By now we should know that it won't. Technical systems raise political issues. Technical people, economists, lawyers, and policy analysts, therefore, must be able to engage in rational dialogue about institutional issues, even when the discussion comes uncomfortably close to home. If we can't, the world is in big trouble.

The review completely misses this big picture. It begins with an attempt to belittle the policy significance of domain name management by inventing a mythical decree that all street names have to be in an obscure language. Crocker's attempt at humor falls flat, given today's headlines. Virtually the same day his review was published a German registrar was ordered to take a domain name away from a Web site with objectionable content. Not too long after, an ICANN Task Force published a WHOIS policy proposal that allows domain names to be shut down after 15 days if someone challenges the accuracy of the contact information, raising issues of privacy and harassment. ICANN regulates the prices of registries and entry into the market for domain name services. No one, not even ICANN itself these days, pretends that domain name administration is an exclusively technical matter.

Instead of engaging on those terms, the review concentrated on factual nitpicking. Take this one: "...the book does not consider NSI's role in ICANN-related political processes." This is an astoundingly inaccurate statement. The index of the book under "Network Solutions" contains 33 listings under 5 separate headings.

The book analyzes at length NSI's origins and ownership changes, its opposition to the IAHC and gTLD-MoU, its implied threat to establish a new root, and its policy conflicts with ICANN and the U.S. Department of Commerce.

Crocker claims that I “[characterize] the pre-ICANN *International Forum for the White Paper* (IFWP) as ‘the real arena for arriving at a decision [about the details of the new organization].’” His use of a sentence fragment covers up what appears to be a deliberate distortion. I really wrote that some people viewed the IFWP in that way, while others, notably Joe Sims, Jon Postel, and the Information Technology Association of America, did not; see pages 176–178. I wrote at length about how that basic lack of agreement between adherents of IFWP and followers of IANA over legitimacy led to lasting conflict over ICANN's formation.

Crocker was one of Jon Postel's appointees to the *International Ad Hoc Committee* (IAHC). The review takes issue with my characterization of the IAHC, but unfortunately only to maintain Crocker's fictional self-conceptions. He denies that the IAHC ever claimed that “the root was theirs to dispose of.” He also denies that IAHC was intended to be the seed of an alternative DNS governance structure. He's wrong on both counts. There is a voluminous record on this question, comprising contemporary news accounts, e-mail list archives, and my own recorded interviews with principal figures such as Don Heath.

Crocker's assertion that IAHC was “explicitly subordinate to IANA” is rather disingenuous, because IANA's U.S. government funding was ending and IAHC was explicitly perceived by Postel and ISOC as a mechanism for continuing its funding. So IAHC was intended to be the governance and support structure for IANA, just as ICANN now is. Indeed, today's ICANN has many features in common with the IAHC proposal, such as the shared registry concept, the slant toward intellectual property interests, the treatment of TLDs as “public resources,” and a compulsory and uniform dispute resolution procedure.

What is really at issue here? It is this: Crocker cannot accept the simple fact that a political battle was under way for control of the root, and Postel/IAHC, as well as NSI and the U.S. government, were contenders for that control. Crocker's review challenges the claim in the book that Postel's root redirection exercise in January 1998 was “apparently” based on “concerns about the direction U.S. policy was taking.” This judgment was based on interviews with people who were involved with Postel's effort. Of course I cannot read Postel's mind, but neither can David Crocker. My interpretation of why Postel acted is based on the timing and on evidence drawn from first-hand participants. Crocker offers an alternative interpretation, plausible but based on nothing but his own assertion. There is plenty of room for legitimate debate about historical interpretations. Such debate is useful, however, only if it is aimed at discovering the truth.

Regarding the status of IANA, I am sure we will never agree. I see it fundamentally as a DARPA contractor subject to U.S. governmental authority; Crocker views it in almost mystical terms as the embodiment of the Internet community. He says nothing about who paid the bills. Yet, we are not as far apart on the facts as he wants to make it seem. Contrary to the review, the book does document in great detail how a new community for Internet standards development grew up around the old DARPA-funded cadre of Postel, Cerf, and the IAB, and created its own standards of legitimacy and process. My book doesn't dispute Postel's tremendous respect and legitimacy among the technical community. But when it comes to institutionalizing control and ownership of the name and address roots of the Internet, whoever pays the piper calls the tune. And Postel's ability to perform the IANA functions was supported by U.S. government money from day one.

Hence, it was unrealistic to expect Postel to be exempt from governmental authority after domain names became resources of economic value and produced legal and political conflict over that value. Nor is it correct to imply, as Crocker does, that knowledge of the operational details of a technology automatically confers wisdom as to the correct public policies that should be adopted when that happens. Of course, policy decisions must respect technical facts and technical constraints. It is this relationship between technical system, technical community, and the worlds of business, law, and government that is central to the story told by *Ruling the Root*.

Crocker's final stab at discrediting the book involves some rather spurious charges of ethics problems. "In his criticism of dispute-resolution activities, he neglects to mention that he is a paid arbitration panelist," he writes. Crocker here refers to the fact that I was one of the few nonlawyers allowed by WIPO to serve as one of three judges in domain name—trademark disputes brought under ICANN's UDRP. The "pay" he refers to is a \$500 or \$750 honorarium for each case. I do about ten cases a year. I fail to see any conflict of interest or ethical problem here. Crocker implies that my meager remuneration for assuring that justice is done in UDRP cases somehow corrupts me, but he knows perfectly well that I am an opponent of the UDRP and would happily stop receiving those honoraria if the darn thing went away. Besides, no one is in a better position to understand what is right and what is wrong with UDRP than someone who is involved in the actual cases. I do not even understand what his concern is about the noncommercial DNSO constituency. I deal with it in one sentence in the book, and most of my activity in a "management capacity" (i.e., as an elective representative) came after the book manuscript was written.

—Milton Mueller, Syracuse University
Mueller@syr.edu

The author of the book review responds:

Professor Mueller's response discusses his goals of the book and his opinions of my review, to which he is, of course, entitled. He characterizes *Ruling the Root* as an academic consideration of the policy issues pertaining to the Domain Name Service, which he casts as global Internet administrative services. Note that the tag line to the title of his book, however, casts it more even more generally as "Internet governance." Academic and policy work need to be conducted carefully. Unfortunately, Professor Mueller confuses the issues, rather than elucidating them.

The opening, mythical decree of the review was carefully constructed to make the perspective of the book on communication system administrative policy clear: Professor Mueller confuses an administrative agency, such as ICANN or its telephonic equivalent, with a national government such as Germany. He also confuses control over administrative information, such as names and addresses associated with registrations, with primary content, such as a Web page.

Professor Mueller defends his writing about the IFWP as merely reporting the view of others, rather than being his own advocacy. However, his reporting is highly selective and results in his confusing the difference between tension that was *within* the IFWP process, versus *between* IFWP and IANA. His casting the issue as being with IANA is contrary to the formal documentation of IFWP, and contrary to the style and content of its process. IFWP was not designed, nor was it conducted, as a decision-making body.

Professor Mueller confuses the actions and intent of the IAHC with those of IANA (and ISOC). He claims to have extensive substantiation for his assessment of the IAHC. Yet none that is relevant to this confusion appears in his book or his letter. This omission is in spite of the fact that his view is at odds with the formal charter for the IAHC, the group's published report, and the direct record of the group's actions.

The review cites IANA's community-based authority. Professor Mueller confuses this with a rejection of the importance of funding, which it was not. He further confuses the IETF technical standards specification process with the operations administrative work of IANA. He continues to misunderstand the role of operational expertise in policy planning for critical infrastructure services, and he ignores the particular 15-year history of successful administrative policy activities provided by operations geeks, for DNS and IP addresses.

Lastly, given the minor points that Professor Mueller chose to address in his response, it is curious that he fails to respond to the primary ethics point raised in the review, namely his pattern of erroneous or absent citations that substantially undermine many of the assertions of his book.

—Dave Crocker, *Brandenburg Internet Working*
dcrocker@brandenburg.com

I recently read Edgar Danielyan's article on Zero Configuration Networking in the December 2002 issue of IPJ. As is always the case, as a journalist Edgar is entitled to hold and express his own opinions, so as I began the article I didn't know whether to expect glowing praise of Zeroconf, or a savage attack. Thankfully I needn't have worried. I found an excellent and well-balanced article.

I have two brief comments to make.

1. Since Edgar wrote his article, the old expired Internet Drafts have been updated. The drafts Edgar worked from discussed names ending in local.arpa. The actual shipping version of Mac OS X 10.2 ("Jaguar") uses names ending in just local. to designate link-local names, (link-local names are locally assigned, unique only within the local link, not required to be globally unique).
2. Edgar expressed the opinion that Zeroconf is only useful on small networks, not large networks.

While Edgar is correct that Zeroconf per se is aimed at solving the "small network" problem, discovering your local peers is useful no matter how big the network. At the recent IETF meeting in San Francisco, there was a large network with full connectivity to the Internet, including IPv6, yet the printers were still advertised using Rendezvous, and for Mac users those printers showed up automatically in the "Printer" popup menu in the print dialogs, with zero configuration.

There is also the issue that Rendezvous (the Apple product) will go beyond just what is required for Zeroconf (the IETF Working Group). Service Discovery, on which Rendezvous is based, doesn't have to be used only with link-local multicast DNS. It can also be used with conventional unicast DNS. For a preview of what the future might hold, you can browse to find an example list of printers at my house. Type: `nslookup -q=ptr _ipp._tcp.stuartcheshire.org`

Thanks for publishing a great article.

—Stuart Cheshire, Apple Computer, Inc.
cheshire@apple.com

List of Acronyms

DARPA	<i>Defense Advanced Projects Agency</i>
DNS	<i>Domain Name System</i>
DNSO	<i>Domain Name Supporting Organization</i>
gTLD-MoU	<i>generic Top Level Domain-Memorandum of Understanding</i>
IAHC	<i>International Ad Hoc Committee</i>
IANA	<i>Internet Assigned Numbers Authority</i>
ICANN	<i>Internet Corporation for Assigned Names and Numbers</i>
IETF	<i>Internet Engineering Task Force</i>
IFWP	<i>International Forum for the White Paper</i>
ISOC	<i>Internet Society</i>
UDRP	<i>Uniform Domain Name Dispute Resolution Policy</i>
WIPO	<i>World Intellectual Property Organization</i>

Book Review

Troubleshooting Campus Networks

Troubleshooting Campus Networks: Practical Analysis of Cisco and LAN Protocols, by Priscilla Oppenheimer and Joseph Bardwell, Wiley, 2002

It is perhaps rare that a book review would encompass the acknowledgements. A break from tradition here is warranted, though, because both authors reveal up front what every prospective reader should know when faced with a purchase decision: Is this work drawn merely from professional circumstance on the part of the author or does it embody a passion held by the author? Judge for yourself. How often do the words “love,” “wonderful,” and “protocol analysis” congregate?

Coauthors Priscilla Oppenheimer and Joseph Bardwell consider the spectrum of protocols and technologies likely to be encountered in a campus environment. A campus network, it is said by the authors, is any one that spans buildings (whether or not in an educational setting). Of course, bricks and mortar are functionally transparent to most modern technologies, and thus the definition of campus could easily be narrowed to any collection of departments or perhaps even any collection of LANs. A contrast is simply being made against the larger metropolitan or wide-area arena.

Although this book does include substantial theory and background for context, it is not yet another rehash of how things *ought* to behave in the vacuum of a lab environment (indeed, the authors occasionally express surprise at their own observations). Neither is it a step-by-step troubleshooting checklist for novice network administrators. To generalize the format, a thorough decomposition of the whole into its many parts follows an introductory discussion of the subject protocol or technology. It is next released into the wild and is quietly observed. Some conclusions are then drawn (some by the authors, some by the reader) regarding appropriate and inappropriate behavior. Lastly, possible courses of action in response to poor or abnormal performance or behavior are considered. This, again, is merely a generalization. The authors take great care to keep the discussion interesting and relevant, often doing so by sharing real-world experiences.

Organization

The six pages that comprise chapter 1 seek to set a stage, define a scope, and target an audience. The reviewer would add only that those of us who trade in wide-area networks also stand to gain a great deal from the experience.

If chapter 2 were packaged for individual sale, it would find its way under the Christmas tree of every colleague, customer, and boss this reviewer has ever encountered. Those readers familiar with Ms. Oppenheimer’s acclaimed *Top-Down Network Design*^[1] may be surprised to find the expression “bottom-up” in any of her work. It is, however, cornerstone not only to the chapter, but also to the remainder of the book.

This seemingly obvious approach to trouble-shooting and analysis could not possibly be emphasized enough according to this reviewer's professional observation.

Chapters 3, 5, and 6 delve into campus datalink layer technologies, protocols, and architectures, including Ethernet, *Spanning-Tree Protocol* (STP), and *Virtual Local-Area Networks* (VLANs). Yawn? The reviewer challenges the reader to finish these three chapters without learning something of considerable value. The Ethernet discussion, for example, breaks from the traditional approach where a cursory review of frame types, cable types, and topologies is deemed sufficient. Where Ethernet came from, where it is going, how it is encoded and presented to the physical layer (and why), and how to interpret frame size distribution using *Remote Monitoring* (RMON) or a protocol analyzer are but a few of the topics considered. Extensive use of protocol analyzer capture files casts new light on STP and VLANs.

Chapter 4 additionally addresses a Layer 2 technology (IEEE 802.11 wireless LANs) but warrants honorable mention. Rare is the *radio frequency* (RF) engineer who possesses a full appreciation for the heretofore all-digital, all-wired campus realm. Perhaps less common would be the network administrator with a capacity to do much other than tune in an FM radio station on a digital set. The authors masterfully string together all the relevant RF concepts, at exactly the right level of detail, to allow for a solid fundamental comprehension of 802.11 networks, technologies, architectures, and deployment. This chapter also would do superbly for anyone with a generic interest in RF units of measurement.

Chapter 7 advances the discussion up to the network layer. Although this may seem common knowledge for readers of a publication such as the IPJ, it is written from the perspective of seasoned protocol analysts. It is worth your time.

Chapter 8 persists at Layer 3 with a thorough discussion of relevant routing protocols. It is again worth noting the emphasis on analysis versus simple textbook theory. It, too, is worthy of your investment.

Chapter 9 rounds out the protocol stack, beginning with an emphasis on Layer 4 protocols *Transmission Control Protocol* (TCP) and *User Datagram Protocol* (UDP). One of the highlights found here is a thorough lesson on TCP window size analysis. Could there perhaps be a little more to this seemingly intuitive concept than you at first thought? The chapter closes following an in-depth consideration of application layer protocols such as the *File Transfer Protocol* (FTP), *Hypertext Transfer Protocol* (HTTP), and the *Domain Name System* (DNS). The fundamental mechanics of these protocols and how they interact with their lower-layer counterparts make for a good page-turner.

Chapters 10, 11, and 12 are dedicated to troubleshooting and analysis of *Internetwork Packet Exchange* (IPX), AppleTalk, and Windows networking, respectively. The latter is arguably the more relevant. The other two are nonetheless interesting and left the reviewer longing for a decent AppleTalk trace file with which to recreate.

Chapter 13, WAN Troubleshooting for LAN Engineers, covers the obvious wide-area technologies and architectures, such as *Integrated Services Digital Network* (ISDN), Frame Relay, and *Synchronous Optical Network* (SONET) in about as much detail as the typical LAN engineer or administrator is likely to tolerate. The subject of WAN analysis warrants a volume or two on its own in any case and thus would have been out of place if explored in much greater detail.

Conclusion

The reading of *Troubleshooting Campus Networks* is not to be approached as a spectator sport. Although the protocol analyzer screen captures are aplenty, and they suitably complement the lessons, merely thumbing the pages would be an opportunity missed. This reviewer chose a free, open-source protocol analyzer (readily available on the Internet) as a reading companion. Although likely far less capable, particularly in terms of graphing, than the oft-referenced Wildpackets EtherPeek product, it nevertheless affords the reader a Layer 2 through 7 window into a living, breathing network.

It bears mentioning that although “Cisco” appears in the subtitle, vendor neutrality is, on the whole, maintained. The Cisco sanctioned troubleshooting methodology is given brief mention in chapter 2. Coverage of the Cisco proprietary *Interior Gateway Routing Protocol* (IGRP), the *Enhanced IGRP* (EIGRP), and the Cisco Discovery Protocol is included, as is coverage of Cisco’s “enhancements” to STP. Lastly, where appropriate, Cisco IOS® “show” and “debug” output is included alongside protocol analyzer screen captures. None of this coverage appears to be included in the spirit of product promotion (bear in mind that this is not a Cisco Press title and that neither author is presently employed by Cisco Systems). Rather, it seems simply to be an acknowledgement that the target audience might very well include candidates for Cisco’s professional and expert-level certification programs (and rightly so).

It is probably anticlimactic that the reviewer would offer a strong buy recommendation for those with an interest in the fundamental interworkings of campus protocols and technologies. The authors’ enthusiasm for packet capture and analysis is infectious. Mr. Bardwell, in fact, is apparently so infatuated that he is at times moved to poetry. This could well be one for the ages.

—Scott Vermillion, IT Artisans Group
scott@itartisans-group.com

References

- [1] *Top-Down Network Design*, Priscilla Oppenheimer, ISBN 1578700698, Cisco Press, 1998.

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Internet Architecture and Technology
WorldCom, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.
Copyright © 2003 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSR STD U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol Journal

June 2003

Volume 6, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
BGP Communities	2
WAP	10
IPv6 Operations Group	20
The Myth of IPv6	23
Letters to the Editor.....	30
Book Review.....	35
Fragments	37
Call for Papers	39

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

FROM THE EDITOR

Articles in *The Internet Protocol Journal* broadly fall into three categories. First, we have articles that explain well-established technologies or operational practices. Second, we offer tutorials on new or emerging protocols and systems, not yet deployed but on the horizon. Finally, IPJ brings you insights, lessons learned and opinions on aspects of networking that have not completely lived up to their promises. In this issue, you will find a mixture of all three.

Our first article is an example from the “nuts-and-bolts” category. The *Border Gateway Protocol* (BGP) is one of the core routing protocols that is widely used in the Internet and has been around for a long time. Kris Foster explains how the *BGP Community* attribute can be used in service provider networks.

Efforts to provide cellular telephones with Internet access systems have produced mixed results. Japan has been leading the way in this area with widespread deployment of iMode devices or variants thereof. Having used such a system I must say I am both impressed and somewhat frustrated. It is wonderful to receive e-mail while on a busy Tokyo train, but accessing the Internet on a tiny screen (typically a 2-inch display with a resolution of 120 x 160 pixels) is not particularly rewarding. Not to mention the bandwidth limitations inherent with this technology. Another system, the *Wireless Application Protocol* (WAP) has been implemented in most countries that offer *Global System for Mobile Communications* (GSM) cell phone service. WAP is the subject of our second article. Edgar Danielyan describes the WAP architecture and looks at some of the lessons learned from its deployment.

The push for deployment of *IP Version 6* (IPv6) is taking place on several fronts and we cover some of them in this issue. In the IETF, a recently formed group has been chartered to help design transition strategies from IPv4 to IPv6. We have a short overview of this effort starting on page 20. Additionally, both the U.S. and Japanese governments are promoting the use of IPv6 in various ways. The U.S. Department of Defense has recently adopted IPv6 as one of its official protocols. In Japan the “IPv6 Appli-Contest 2003” is underway in an effort to encourage development of software and applications for IPv6. See “Fragments,” page 37–38 for further details.

Of course, not everyone is convinced that IPv6 is such a good idea, and with that in mind we bring you an opinion piece as well as a Letter to the Editor on this topic.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

Application of BGP Communities

by Kris Foster, TELUS

The *Border Gateway Protocol* (BGP) is the glue that binds networks and their individual policies together. Several attributes are passed along and possibly modified with each individual prefix, one of which is the *community* attribute. BGP communities are described poorly in most texts. The problem is not in explaining how they fit into the protocol, but in how to apply these to the real world. In this article I describe how they can be applied within a service provider network and between service provider networks. However, communities are not limited to service providers and can be applied creatively in enterprise networks.

The density of interconnection among service providers, and the various business agreements or political policies, means that controlling who can talk to whom over your network can become difficult. At a basic level there are two types of agreements between service providers: transit/customer and peers.

- Customers pay to receive every prefix from a transit provider.
- Customers advertise only the prefixes they own (along with their customers' prefixes) to the transit provider.
- Peers agree to send only their customers' prefixes to each other, and not other peers' prefixes.

Several methods are available to implement these policies. They can include prefix filters, *Autonomous System* (AS) path filters, and communities. With only prefix and AS path filters, service providers must ensure that as a new customer or peer is added, the prefixes and *AS Numbers* (ASNs) associated with the customer (and potentially *their* customers) are added to the filters on all of the BGP edge routers. This can be automated with scripts, possibly in combination with a route registry database. Very small service providers may be able to manage such a scheme, but as they grow and customer churn begins, this can quickly get out of control. The more time network operators spend in router configurations, the greater likelihood of human error. Communities provide an elegant solution for these problems.

The BGP Community Attribute

Within an AS, all BGP-speaking routers run *Internal BGP* (iBGP) in a full mesh to prevent routing loops (route reflectors can be used to relax this rule). This means that every BGP-speaking router passes its prefixes to each of its iBGP neighbors. ASs that are adjacent typically run eBGP on directly connected routers. All BGP routers share their prefixes—that is, the network number, network mask, and BGP attributes with each other—allowing each to run its own best-path selection algorithm. As a prefix is passed between ASs, an attribute called the AS-PATH is updated with the corresponding ASN. The AS-PATH is used to prevent routing loops between eBGP neighbors.

A community is a BGP attribute that may be added to each prefix. Communities are transitive optional attributes^[1], meaning BGP implementations do not have to recognize the attribute and at the network operator's discretion carry it through an AS or pass it on to another AS. The community attribute can be thought of as simply a flat, 32-bit value that can be applied to any set of prefixes. It can be read as a 32-bit value or split into two portions, the first 2 bytes representing an ASN and the last 2 bytes as a value with a predetermined meaning. The format of the community attribute is shown in Figure 1.

The values **0x00000000** through **0x0000FFFF** and **0xFFFF0000** through **0xFFFFFFFF** are reserved. Most modern router software displays communities as **ASN:VALUE**. In this format the communities **1:0** through **65534:65535** are available for use. The convention is to use the ASN of your own network as the leading 16 bits for your internal communities and communities that you accept from and send to your customers.

Three communities are defined in RFC 1997^[2] and are standard within BGP implementations: NO-EXPORT (**0xFFFFFFFF01**), NO-ADVERTISE (**0xFFFFFFFF02**), and NO-ADVERTISE-SUBCONFED (**0xFFFFFFFF03**). Additionally, NO-PEER (**0xFFFFFFFF04**) has been proposed in an Internet Draft^[3].

NO-EXPORT is commonly used within an AS to instruct routers not to export a prefix to eBGP neighbors. For instance, subnets of a larger block can be advertised to influence external AS best-path selection, and those not required for this traffic engineering purpose may be tagged NO-EXPORT to prevent them from being leaked to the Internet (and thus contributing to unnecessary global routing table growth). If a neighboring AS accepts this community, it can be used to selectively leak more specifics for traffic engineering but limit their propagation to just one AS.

NO-ADVERTISE instructs a BGP-speaking router not to send the tagged prefix to any other neighbor, including other iBGP routers.

NO-ADVERTISE-SUBCONFED is used to prevent a prefix from being advertised to other members within a *confederation*. A confederation can be thought of as a single AS, broken down into sub-ASs. The use of confederations within service provider networks is rare or nonexistent, so they are not considered here.

Finally, NO-PEER is used in situations where traffic engineering control over a more specific prefix is required, but to constrain its propagation only to transit providers and not peers. That is, the prefix is advertised from AS to AS provided there is a transit/customer relationship, unlike NO-EXPORT, which restricts propagation of the prefix to only the adjacent AS. Because peers of the various upstream providers will not see this prefix, the larger prefix encompassing the more specific one is used for routing, thereby conserving an extra entry for some in the global routing table. At this time the community is not recognized by major vendors and requires manual implementation.

Adding Depth: The Extended Community

The current community attribute is getting an upgrade with a new transitive-optional attribute (Type 16) called the *Extended Community*^[4]. Missing from regular communities was any real form of structure. The current Internet Draft defines the Extended Community as an 8-octet value as shown in Figure 1. The first octet specifies the type (and optionally the second value can specify a subtype). This value dictates the structure given to the remaining octets.

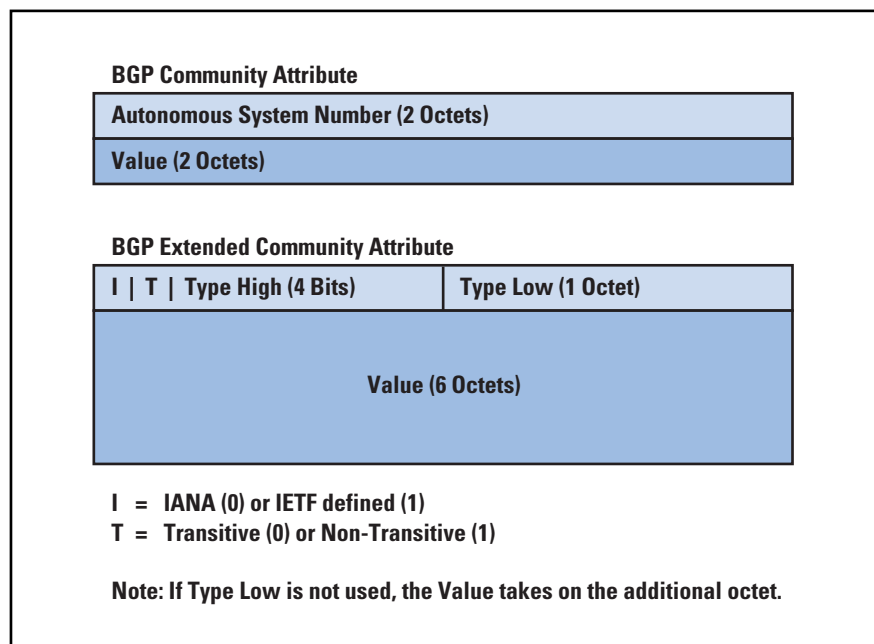
The Type field gives the community some immediate flexibility. The first is the use of bit 0 to represent whether the community is registered with the *Internet Assigned Numbers Authority* (IANA) or if it is specified by the *Internet Engineering Task Force* (IETF). The second bit gives the Extended Community a coarse scope, either *Transitive*, meaning it may be passed between ASs, or *Non-Transitive*, meaning it should be carried only within the local AS.

The Internet Draft also specifies numerous types available for use as templates.

The *Route Target Community* is already in popular use within *Multi-protocol Label Switching Virtual Private Networks* (MPLS VPNs). The Route Target Community identifies a set of routers that may receive this prefix. In the MPLS VPN context, this is necessary to limit the resources required to support individual VPN services; only routers that are part of the individual VPN need to hear about the routes within the VPN.

The *Link Bandwidth Community* gives the network operator additional control in influencing the best path selection. As prefixes are learned from eBGP neighbors, the local neighbor applies this community to specify in bytes per second the bandwidth of the link. It is a *Non-Transitive Community*, so its scope is limited to the local AS.

Figure 1: Community Formats

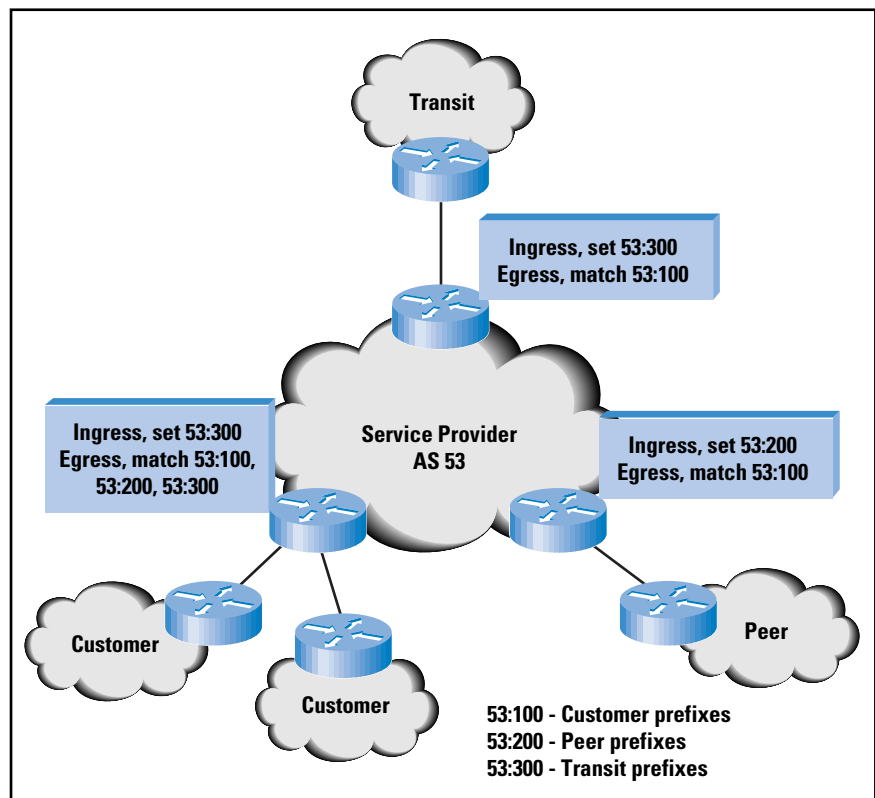


Intra-Autonomous System Communities

Policy control using communities within an AS can go farther than this, and their true value is evidenced when they are used to create new and complex policies. If we take our example of the three basic types of neighbor relationships, customers of a transit provider will want to send their customers' prefixes but not their peers' prefixes. To distinguish between a customer's prefix, a peer's prefix, and a transit provider's prefix, we can add a community to each as we learn it from the neighbor.

When advertising a prefix to a customer, peer, or transit provider, simply match all prefixes carrying the communities associated with the correct policy. As shown in Figure 2, all prefixes received from customers are tagged with **53:100**, peers are tagged with **53:200**, and transit is tagged with **53:300**. Our basic definition of a customer is someone who expects to receive all prefixes, so each customer-facing BGP session is preconfigured to send all prefixes matching **53:100**, **53:200**, and **53:300**. Again, from our definition of a peer being someone who wants to see only our customers, we would preconfigure all of our peers' BGP sessions to send only prefixes tagged with **53:100**.

Figure 2: Internal Use of Communities for Applying a Basic Service Provider Policy



We can extend this community coding and turn it into a useful troubleshooting tool by adding more information such as where the route was learned geographically. Codes could be assigned per continent, country, state/province, city, or central office.

During redistribution from an Interior Gateway Protocol, a community can be used to specify the original protocol (for example, *Intermediate System-to-Intermediate System* [IS-IS], *Open Shortest Path First*

[OPSF], or *Routing Information Protocol* [RIP]). These can be used to quickly determine where a prefix came from without tracing it back to the point of its origination.

It is possible to assign these additional properties in two different ways (or a combination). A single community value may represent a single meaning, such as **53:100**, meaning a customer-learned prefix. We could then add additional communities such as **53:1** to mean a prefix learned on the east coast, **53:2** to mean central, and **53:3** to mean west coast. Alternatively, a single community could represent both a customer and a prefix learned on the west coast by tagging with the single tag **53:103**. To support these complex values, most vendors allow for pattern matching of specific values, ranges of values, and logical operators such as OR and NOT, in the form of regular expressions. Using regular expressions and complex communities can help to make a router configuration more economical and easier to read.

Inter-Autonomous System Communities

We have some options for Inter-AS traffic engineering: we can prepend additional AS numbers onto a prefix path, use *Multi-Exit Discriminators* (if the provider supports this), announce more specific prefixes or not announce prefixes at all, modify the origin type, or use communities designed by the other service provider. Communities are clean and consistent with regard to the method of signaling to an adjacent AS how each prefix should be treated.

Of most concern to downstream customers is controlling their primary and backup circuits. Small service providers and enterprises may negotiate different rates on different circuits. Customers purchasing transit with a commitment to send a high amount of traffic with a lower cost per megabit on one circuit, and on a second circuit purchase transit with a very low commitment but at a higher cost per megabit can save some money, assuming they use only the second circuit during outages on the first. Two simple communities can be used to effectively influence a service provider into using the appropriate primary and backup circuits: one value to lower and another to raise the preference of specific prefixes during the transit provider's best-path selection.

An example of adjusting Local Preference with communities can be found in RFC 1998, "An Application of the BGP Community Attribute in Multi-home Routing"^[5].

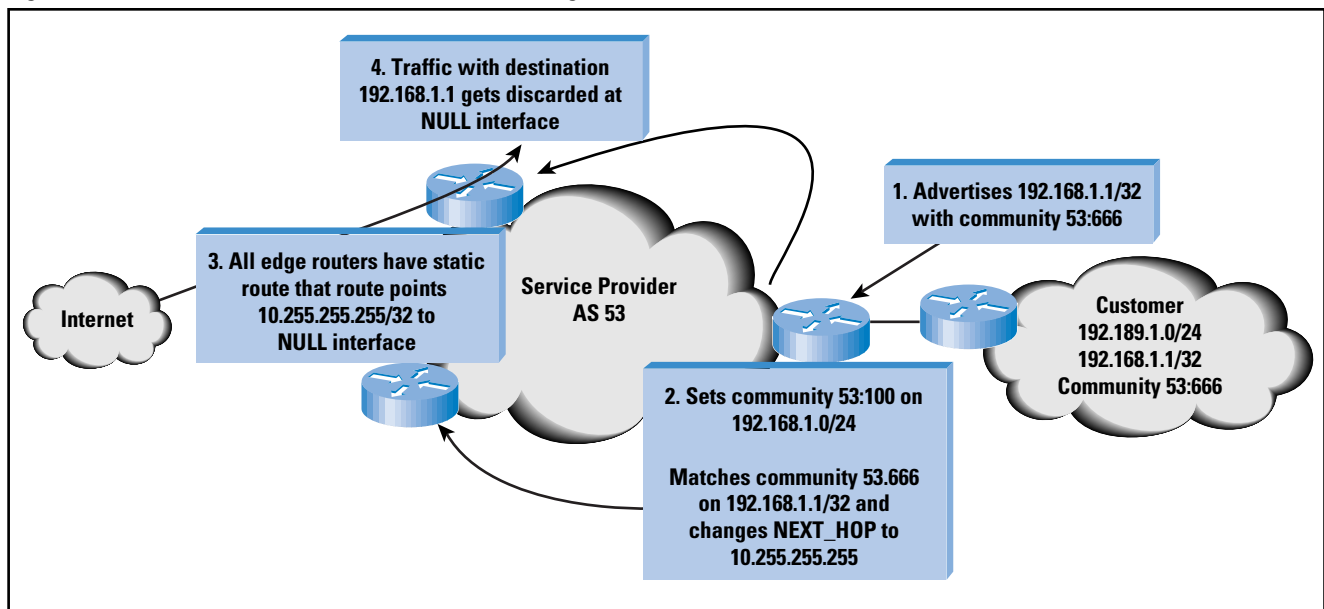
Some other traffic engineering signaling possibilities include:

- Force the adjacent AS to prepend its ASN a certain number of times to a prefix sent to customers or peers.
- Force the other side to selectively advertise a prefix to specific neighbors.
- Request that the neighbor drop all traffic to a prefix.

The last example may seem a little strange; if you are paying someone to deliver traffic, you expect to receive that traffic. Here is where communities can play a role in network security. *Denial-of-Service* (DoS) attacks may take out an entire customer's service, but the attack may be

focused on one or several hosts and not an entire network, as illustrated in Figure 3, allowing customers to tag individual host routes (a subnet consisting of a single address), the customer can signal to the provider to drop all traffic (black hole) for that specific address. To achieve this, the provider selects a single IP address and routes all traffic destined for it to the NULL interfaces on every BGP-speaking router. When a customer signals for a prefix to be blackholed, the service provider replaces the NEXT_HOP information in the BGP advertisement (which under normal circumstances is the edge router IP address) with the specific address that all other routers have statically routed to the NULL interface. When a packet arrives destined for the host under attack, the edge router performs a routing table lookup to find the BGP prefix; using the NEXT_HOP, it then performs a recursive lookup and ultimately sends the packet out the NULL interface. It is important to use other techniques such as prefix lists to prevent a third party from exploiting this technique to disrupt service for others in the Internet.

Figure 3: Customer-Initiated Black Hole to Defend Against a DoS Attack



A service provider may elect to send communities to its customers, leaving it up to the customers to decide for themselves which communities to act on. For a customer who is dual-homed to the same service provider in multiple states or countries, it may be helpful to know where a prefix was originated. A customer could use this community to prefer a connection in New York instead of a Los Angeles connection for European traffic. A single composite metric composed of all relevant geographical information is best, because this gives customers maximum flexibility in choosing the values that are meaningful to them.

Tagging the type of prefix may help other networks to selectively filter more specific addresses. Adding a community specifying if a block is a more specific part of a *Classless Inter-Domain Routing* (CIDR) block being advertised, the CIDR block itself, or if it is a more specific block but the CIDR block is not being advertised, can help the downstream network avoid incorrect filtering.

Example: Network A announces

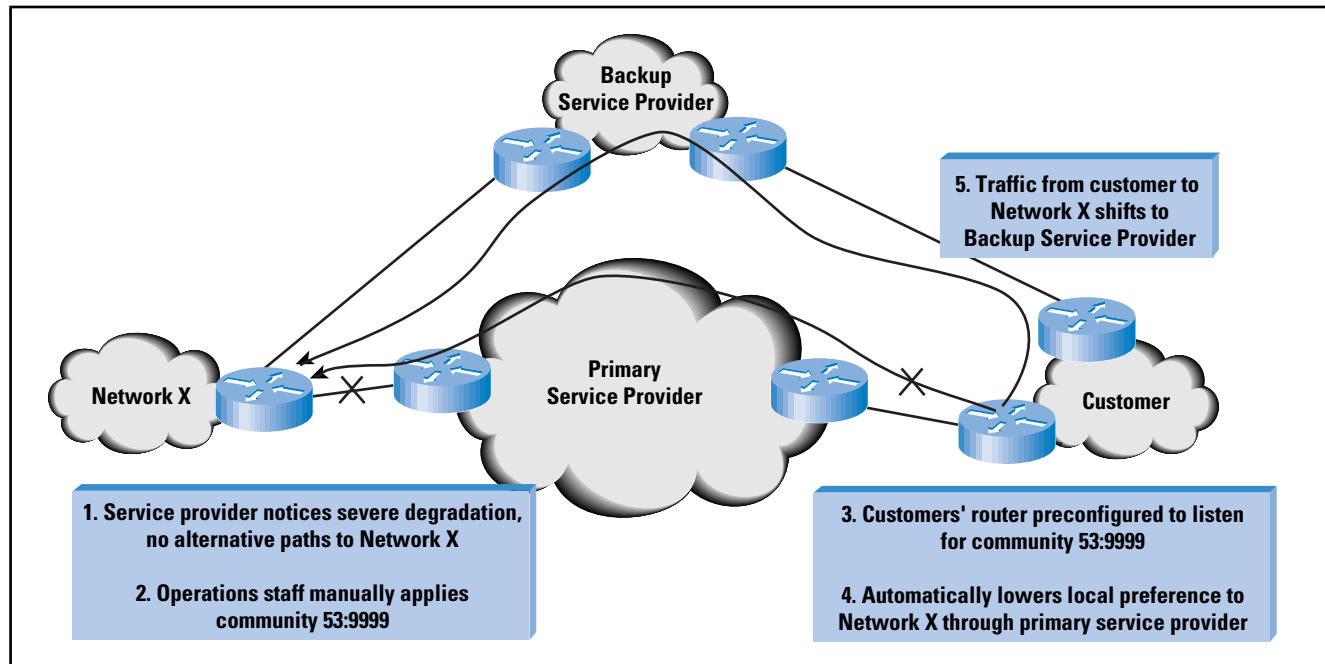
142.77.0.0/16 with a tag of 1:77
 142.77.1.0/24 with a tag of 1:88
 150.3.12.0/24 with a tag of 1:99

1:77 means it is a CIDR block
 1:88 means it is a more specific block within a CIDR block
 1:99 means that the full CIDR block is not being announced

Network B then has the option of accepting the more specific 142.77.1.0/24. It also knows that it must accept 150.3.12.0/24 because there is no other route to this network.

In extreme cases providers may find that a portion of their network has become severely degraded. Planned with customers in advance, the upstream provider manually sets a specific community on prefixes associated with the degradation to indicate that this path should be avoided. This could be helpful during natural disasters, fiber cuts, or other unanticipated network outages/degradation. The downstream customers' inbound filters would then match this community and lower the preference on the prefixes tagged with it, causing them to automatically shift traffic to an alternative source if it is available. The degradation signalling process can be seen in Figure 4.

Figure 4: Provided Initiated Signalling of Severe Route Degradation



Design Recommendations

The following are some suggestions if you are just starting out with using communities in your own network. Even the smallest network can benefit from starting early with a clean community design.

- Choose a set of internal communities that best reflects the topology and characteristics of your network. For external communities some service providers offer none, others offer only enough to allow for the tagging of primary and backup circuits, and others provide a seemingly endless list.
- Keep the set simple. Adding additional complexity typically requires changes to all the BGP-speaking edge routers. Router configurations can quickly grow to enormous proportions to accommodate the numerous community combinations. Troubleshooting a routing mess with a complex community structure can be difficult for those on the graveyard shift.
- Avoid transiting communities received from neighboring ASs blindly through your network. This could be abused intentionally or unintentionally to influence traffic to use your costly transit over settlement-free peering and revenue-generating customer circuits. Problems can be created farther out in the Internet and can be very difficult to locate. Depending on the support of your router software, you may be able to selectively add and remove communities, or failing that, you may need to remove all communities and re-add what is acceptable.
- Document your communities internally and externally. Your customers will appreciate the additional control, and your operations team will have an easier time troubleshooting.

Summary

Communities add power to BGP, changing it from a routing protocol to a tool for signaling and policy enforcement. If deployed correctly and consistently, communities can help make a network scale, easier to operate, easier to troubleshoot, and can give its customers what they want.

References

- [1] Y. Rekhter and T. Li, "A Border Gateway Protocol 4 (BGP-4)," RFC 1771, March 1995.
- [2] R. Chandra, P. Traina, and T. Li, "BGP Communities Attribute," RFC 1997, August 1996.
- [3] G. Huston, "NOPEER Community for BGP Route Scope Control," Internet Draft, May 2003.
- [4] S. Sangli, D. Tappan, and Y. Rekhter, "BGP Extended Communities Attribute," Internet Draft, May 2002.
- [5] E. Chen and T. Bates, "An Application of the BGP Community Attribute in Multi-home Routing," RFC 1998, August 1996.

KRIS FOSTER, CCIE® #7749, currently lives in Calgary, Alberta, and spends his time in TELUS' IP backbone. His industry affiliations include the Association for Computing Machinery (ACM), the Internet Society (ISOC), and the North American Network Operators Group (NANOG). He can be reached at kris.foster@telus.com

WAP: Broken Promises or Wrong Expectations?

by Edgar Danielyan, Danielyan Consulting LLP

The *Wireless Application Protocol* (WAP) was once hailed as the ultimate mobile Internet solution that would revolutionize how we use the Internet and mobile phones. As you may already know, it didn't. What is to blame? Is it bad technology, wrong time, or greedy network operators? Actually, is there a reason to blame anyone? This article introduces WAP with its related technologies and tries to answer these questions. Although WAP is available on a variety of wireless mobile networks, such as those employing *Code Division Multiple Access* (CDMA) IS-95, *Time Division Multiple Access* (TDMA) IS-136, *International Mobile Telecommunications* (IMT-2000), *Universal Mobile Telecommunication System* (UMTS), and *Wideband Code Division Multiple Access* (W-CDMA), in addition to GSM/GPRS this article covers WAP over GSM/GPRS networks only.

A Case for WAP

Before looking at WAP itself, let's first recall what sparked its idea and development. As we all know, most if not all *second-generation* (2G) mobile phones and networks suffer from numerous limitations that make it impossible or impractical to use standard Internet protocols and technologies on today's mobile phones. The most visible of these limitations include the following:

- Low bandwidth (usually 9.6 kbps)
- High network latency
- Small, mostly monochrome displays
- Numeric keypads
- Slow processors
- Limited memory

All these limitations meant that it was necessary to develop an alternative suite of protocols and technologies that would work on these mobile phones but still provide functionality comparable to the standard Internet technologies used on wired networks and desktops. WAP was developed to address these issues^[1].

WAP Forum and Open Mobile Alliance

The *WAP Forum* is the industry organization behind WAP and its associated protocols and technologies. In 2002, the WAP Forum and the Open Mobile Architecture Initiative merged, creating the *Open Mobile Alliance* (OMA), which will continue work on WAP 2 and develop new mobile and wireless solutions. Nearly 200 of the world's top network operators, vendors, and content providers are members of the Open Mobile Alliance^[2]. Other organizations such as the *Location Interoperability Forum* (LIF)^[3], *Multimedia Messaging (MMS) Interoperability Group* (MMS-IOP)^[4], *SyncML Initiative*^[5], and *Wireless Village Initiative*^[6] have announced their support for the new organization.

Global System for Mobile Communications

GSM, or *Global System for Mobile Communications*, is used by more than 700 million people across 190 countries^[7]. In less than ten years after its introduction, GSM became the most popular and widely used digital mobile wireless communications standard in the world. GSM networks use TDMA technology and are fully digital, employing a unique voice codec known as *GSM codec* to provide relatively good voice quality using narrow bandwidth (usually 9.6 kbps). However, GSM is not as secure as many may think. Although it does use encryption and smartcard technology, this didn't result in strong security. As a result, it is possible to intercept and decrypt GSM communications, fake short text messages (*Short Message Service* [SMS]), and clone *Subscriber Identification Modules* (SIMs), miniature smartcards used to identify subscribers to the GSM network. GSM security is not the subject of this article, but it deserves attention and I hope to cover it in a separate article in this journal.

Wireless Application Environment

Before proceeding further, we should clarify one point. The term "WAP" is usually used to refer to the entire suite of protocols and technologies that are actually called the *Wireless Application Environment* (WAE)^[8]. However, "WAP" is used everywhere to refer to WAE (which includes WAP). Because WAP is the commonly used term, we shall continue to use it as well.

Wireless Application Protocol

WAP protocols were expected to satisfy the following criteria in order to implement the objectives set by the WAP Forum:

- Independent of wireless network standard (bearer technology)
- Open to all
- Will be proposed to the appropriate standards bodies
- Applications scale across transport options
- Applications scale across device types
- Extensible to new networks and transports

The objectives of the WAP as defined by the WAP Forum follow:

- To bring Internet content and advanced data services to digital cellular phones and other wireless terminals
- To create a global wireless protocol specification that will work across differing wireless network technologies
- To enable the creation of content and applications that scale across a very wide range of bearer networks and device types
- To embrace and extend existing standards and technology wherever appropriate

Two major versions of WAP exist—Versions 1 and 2. WAP Version 2 is backward compatible with WAP Version 1 and tends to be more integrated with the newest Internet and Web standards than WAP 1. Although WAP uses many technologies and concepts from the Internet and Web worlds, because of their inherent limitations, WAP devices are unable to directly access Web resources on the Internet^[9]. To do so, they must use a WAP gateway. The following table shows the relationship between the WAP client device, WAP gateway, and Web servers on the Internet, with their protocol layers side by side:

Web Client	WAP Gateway	Web Server
WSP	WSP/HTTP	HTTP
WTP	WTP/HTTP	HTTP
WTLS	WTLS/SSL/TLS	SSL/TLS
WDP	WDP/TCP/UDP	TCP/UDP
Bearer	Bearer/IP	IP

The table shows that the main function of the WAP gateway is to translate between WAP and Web/Internet protocols, conventions, and encodings. In some cases the WAP gateway and the Web server may be the same system, eliminating the need for a separate WAP gateway and possibly improving performance—however, for this setup to work the combined WAP/Web server has to be integrated into the mobile/wireless network provider’s infrastructure. In practice, network operators provide the WAP gateway services and content providers offer WAP content on separate Web servers configured for WAP access (any standards-compliant Web server can do this).

Wireless Session Protocol

The *Wireless Session Protocol* (WSP) is the WAP session-layer protocol for remote operations between a wireless (WAP) client and proxies, gateways, and servers^[10]. It functions above the *Wireless Transaction Protocol* (WTP) and the *Wireless Datagram Protocol* (WDP), and optionally, the *Wireless Transport Layer Security* (WTLS). The WSP provides a way for an organized exchange of data between client/server applications in a wireless environment. It provides such features as establishment and release of sessions between client and server; agreement on common functionality by way of negotiation; and exchange of data between client and server using compact encoding. WSP defines two subprotocols—a connection-oriented session service protocol over WTP and a connectionless service protocol over the WDP.

Wireless Transaction Protocol

WTP runs on top of the WDP and optionally, the WTLS protocol, and provides the request/response protocol used by WAP browsers to request and receive content^[11]. WTP is a reliable transaction-oriented protocol specially designed for wireless networks—in WTP there are no connection setup or release phases.

Reliability in WTP is achieved using transaction IDs, retransmissions, acknowledgments, and removal of duplicates.

Wireless Datagram Protocol

WDP is the transport protocol of WAP^[12]. It operates directly above the bearer technology (such as GSM CSD or GPRS) and directly below WTP described previously. WDP provides a consistent, bearer-independent interface for the upper-level protocols to the transport service provided by WDP. In addition to the GSM *Circuit Switched Data* (CSD) and the *General Packet Radio Service* (GPRS), WDP supports the following wireless bearer technologies:

GSM SMS	IDEN Packet Data
GSM USSD	FLEX
GSM Cell Broadcast	REFLEX
ANSI-I36	PHS CSD
CDPD	DataTAC
CDMA CSD	TETRA Short Data Service
CDMA Packet Data	TETRA Packed Data
CDMA SMS	DECT SMS
PDC Circuit Switched Data	DECT Connection-oriented Service
PDC CSD	DECT Packed Switched Service
PDC Packet Data	Mobitex
IDEN CSD	

When used over GSM CSD, WDP actually uses the *User Datagram Protocol* (UDP) in the following way:

Layer 4: UDP
Layer 3: Internet Protocol (IP)
Layer 2: Point-to-Point Protocol (PPP)
Layer 1: GSM CSD

When used over the GPRS, PPP at Layer 2 is not necessary, because GPRS works at Layers 1 and 2:

Layer 4: UDP
Layer 3: IP
Layers 1 and 2: GSM and GPRS

In all cases when IP is supported over a given bearer, UDP is used by WDP—actually, UDP is the WDP in these cases.

Wireless Control Message Protocol

Not surprisingly, *Wireless Control Message Protocol* (WCMP) resembles and corresponds to the *Internet Control Message Protocol* (ICMP) of TCP/IP networks^[13]. WCMP is used by WDP nodes to report errors and provide network information and diagnostics. However, WCMP is not necessary and is not used with bearers that support IP—the function of WCMP in these circumstances is carried out by ICMP. In particular, this is the case with GSM CSD and GPRS bearers.

Wireless Transport Layer Security

WTLS is the transport layer security protocol of the WAE that provides privacy, integrity, and authentication services^[14]. It is heavily influenced by the *Transport Level Security* (TLS) protocol Version 1 and includes additional support for optimized handshake, connectionless transport, and dynamic key refresh. WTLS, like other WAP protocols, is optimized for low-bandwidth, high-latency wireless networks and supports server and client certificates for mutual authentication. WTLS includes the following three subprotocols:

- Cipher protocol
- Alert protocol
- Handshake protocol

The following cryptographic algorithms are used by the Wireless TLS protocol:

- RSA
- SHA-1
- Diffie-Hellman (DH)
- Elliptic Curve Diffie-Hellman (EC-DH)
- DSA
- Elliptic Curve DSA (EC-DSA)
- MD5
- RC5
- DES
- IDEA

WTLS is tightly linked to and works in conjunction with the *Wireless Public Key Infrastructure* (WPKI).

Wireless Public Key Infrastructure

WPKI tries to reuse the existing *Public Key Infrastructure* (PKI) standards as much as practical to provide an adequate PKI framework for the WAE. Both X.509 and WTLS certificates can be used by WTLS^[15].

Wireless Markup Language Version 1

The *Wireless Markup Language* (WML) Version 1^[16] is used in WAP/WAE 1 and supported in WAE 2. Unlike usual HTML, it is a strict application of the *Extensible Markup Language* (XML), specially designed for use on narrowband devices. Also unlike HTML, WML has a metaphor of *decks* and *cards*. A deck contains one or more cards, and cards in turn contain one or more screens of user interaction. This metaphor helps increase efficiency on low-speed, high-latency wireless networks by bundling several screens into a single WML file (deck). WML supports all basic text display options, such as *italic*, **boldface**, and underlined text, as well as inter-card and inter-deck navigation using hyperlinks. The most apparent difference between HTML and WML noted by HTML developers is the fact that WML is a strict markup language and does not tolerate even seemingly little errors—an incorrectly written WML file will not display at all. Some would say this is an overkill but it is not—this feature of WML is important because compiled versions of WML files are sent to WAP clients by the WAP gateway instead of the source WML text files. This compiled bytecode is known as *WMLC*, and it considerably lessens the time it takes to download a WML document.

WML Version 2

WML version 2 is based on XHTML Basic with additional modules for support of features specific to wireless devices—this extended XHTML is called *XHTML Mobile Profile* (XHTML-MP)^[17]. WML Version 2 is backward compatible with WML Version 1, so devices able to display WML 2 will also display WML 1 content. Use of XHTML shows that WAP in Version 2 is moving toward even closer integration with Internet and Web standards.

WMLScript

WMLScript is a lightweight scripting language based on ECMAScript, which is in turn based on JavaScript^[18]. It is well integrated with WML and has a defined set of standard libraries, including support for cryptographic functions. Like WML, WMLScript files are also compiled into bytecode and only then sent to the requesting WAP device. Another difference between JavaScript and WMLScript is that WMLScript content is not embedded in WML pages but instead is requested separately—the necessary WMLScript functions are only referenced in WML pages. The main use of WMLScript is the client-side validation of user input—accepting only valid input is more crucial for WAP than for Web applications because of the low-speed and usually expensive nature of WAP transport.

Wireless bitmaps

The *Wireless Bitmaps* (WBMP) file format (**.wbmp**) is used by WAP devices to transmit and display small and simple monochrome bitmap images^[19].

GSM CSD

CSD is the traditional data service provided by GSM networks. Also known as a *data call service*, it provides either a 9.6- or 14.4-kbps dialup facility and is supported by all GSM networks. Data calls are possible both from and to a GSM network. When used as a bearer for WAP, it serves at the physical layer of the *Open System Interconnection* (OSI) model, with PPP used in the usual way.

High-Speed Circuit Switched Data

The *High-Speed Circuit Switched Data* (HSCSD) service is similar in nature to CSD, but provides 28.8 or 43.2 kbps of bandwidth. It is not as widespread as the regular CSD, nor it is as asked-for as GPRS.

General Packet Radio Service

GPRS is an always-on, higher-speed alternative to the CSD service of GSM networks. It solves two of the most annoying issues of GSM data users—connection delay (the time it takes to set up a data call before data may be sent or received) and the bandwidth limitation, increasing the supported data rates to 48 kbps, with theoretical maximum of 171.2 kbps. Because GPRS is a connectionless packet service, GPRS terminals are always connected and may send and receive IP packets at any time. This makes possible applications such as instant messaging previously impossible or impractical with GSM CSD. Eight time slots are available for GPRS in GSM networks, but only five may be used simultaneously. The GPRS class supported by the GPRS terminal dictates what data rates are possible:

Class 2:	Uplink 8–12 kbps, downlink 16–24 kbps
Class 4:	Uplink 8–12 kbps, downlink 24–36 kbps
Class 6:	Uplink 16–24 kbps, downlink 24–36 kbps, or Uplink 24–36 kbps, downlink 16–24 kbps
Class 8:	Uplink 8–12 kbps, downlink 32–40 kbps
Class 10:	Uplink 8–12 kbps, downlink 32–48 kbps, or Uplink 16–24 kbps, downlink 24–36 kbps
Class 12:	Uplink 8–12 kbps, downlink 32–48 kbps, or Uplink 16–24 kbps, downlink 24–36 kbps, or Uplink 24–36 kbps, downlink 16–24 kbps, or Uplink 32–48 kbps, downlink 8–12 kbps

In addition to the classes of GPRS service, there are three classes of GPRS terminals:

- Class A terminals can be connected to GSM and GPRS services simultaneously.
- Class B terminals can be connected to both GSM and GPRS services, but can use only one service at a time.
- Class C terminals can be connected to either GSM or GPRS services but the user has to switch between two modes of operation.

When used as a bearer for WAP, GPRS works at the physical and data link layers of the OSI reference model. Because GPRS is connectionless and always on, there is no need for PPP—so IP works directly over GPRS.

So Why Aren't We Happy with WAP?

Many surveys of customer opinion show that the end users of WAP are not as happy as WAP developers and content providers wanted them to be. WAP service and content providers discovered that sign-up and usage rates of WAP services have not reached two-thirds of the total customer base once predicted. In short, WAP didn't change the world, and people still use their mobile phones mainly to talk to each other and send a text message or two. If you have used WAP, you probably know the reasons: the data transfer rate is slow, screens are small, charges are high, and it is tiring to type even a short URL or an e-mail message using the ten keys of a phone.

But wait a moment—are these limitations of WAP or the handsets and networks they use? Remember, WAP was required to work on devices with many limitations? So it does. Is WAP to blame that these devices have these limitations? No, that wouldn't be just. But of course it is not only the today's technology restrictions that stood in the way of the widespread usage and popularity of WAP. Scarcity of WAP content and services also contributed to this. Relatively high charges for WAP/data usage by network operators didn't help either, so the combination of these issues resulted in the situation we have today—most networks support WAP but most users don't use it anyway.

Is the technology dead, as some think? Definitely not—there are millions of WAP handsets and most wireless users will not have 3G for the foreseeable future because of both technical and economic issues, so the only available solution for these users is WAP. On the other side, 3G networks and handsets are coming and will be upon us sooner or later (they are already available in some countries), and only time will show whether tomorrow's WAP will be more popular or less relevant when 3G finally arrives. And, of course, fundamental limits of mobile phones—screen sizes, power consumption, and input methods—will still remain relevant. Other issues, such as the time it takes to set up a CSD connection, are solved by newer technologies such as GPRS, and are not really faults of WAP. You may say that if GPRS is available why would you need WAP? Why not run trusted IP? Well, this is true if you are using GPRS with a laptop or a palmtop computer, but a large majority of mobile phones don't have the resources necessary to run IP, UDP, TCP, HTTP/HTTPS, POP, and SMTP—so even if GPRS is available but your equipment cannot run the full TCP/IP suite, your only choice is still WAP.

Although WAP is clearly not as popular as its proponents and developers hoped, it is still used and developed, and handsets that support only WAP are still sold. But the hype and excitement built up by the media and the industry didn't match the reality, and it is these unrealistic expectations that have broken the promise of WAP.

Additional Acronyms

DataTAC:	<i>Motorola wireless data system</i>
DECT:	<i>Digital Enhanced Cordless Technology</i>
DES:	<i>Data Encryption Standard</i>
DSA:	<i>Digital Signature Algorithm</i>
FLEX:	<i>Motorola one-way paging system</i>
IDEA:	<i>International Data Encryption Algorithm</i>
IDEN:	<i>Integrated Dispatch Enhanced Network</i>
MD5:	<i>Message Digest 5</i>
PDC:	<i>Pacific Digital Cellular System</i>
RC5:	<i>Rivest Cipher 5</i>
REFLEX:	<i>Motorola two-way paging system</i>
SHA-1:	<i>Secure Hash Algorithm 1</i>
TETRA:	<i>TERrestrial TRunked RAdio</i> Nokia open digital professional mobile radio standard
USSD:	<i>Unstructured Supplementary Service Data</i>

For Further Reading

- [1] WAP Forum: <http://www.wapforum.org>
- [2] Open Mobile Alliance: <http://www.openmobilealliance.org>
- [3] Location Interoperability Forum:
<http://www.openmobilealliance.org/lif>
- [4] MMS Interoperability Group (MMS-IOP):
<http://www.openmobilealliance.org>
- [5] SyncML: <http://www.openmobilealliance.org/syncml>
- [6] Wireless Village: <http://wireless-village.org>
- [7] Global System for Mobile Communications (GSM):
<http://www.etsi.org>, <http://www.gsmworld.com>
- [8] Wireless Application Environment (WAE) Version 2.0:
<http://www.wapforum.org>
- [9] Wireless Application Protocol Architecture Specification:
<http://www.wapforum.org>
- [10] Wireless Session Protocol Specification: <http://www.wapforum.org>
- [11] Wireless Transaction Protocol Specification:
<http://www.wapforum.org>

- [12] Wireless Datagram Protocol Specification:
<http://www.wapforum.org>
- [13] Wireless Control Message Protocol Specification:
<http://www.wapforum.org>
- [14] Wireless Transport Layer Security Specification:
<http://www.wapforum.org>
- [15] Wireless Public Key Infrastructure Architecture Specification:
<http://www.wapforum.org>
- [16] Wireless Markup Language Version 1 Specification:
<http://www.wapforum.org>
- [17] Wireless Markup Language Version 2 Specification:
<http://www.wapforum.org>
- [18] WMLScript Specification: **<http://www.wapforum.org>**
- [19] Wireless Bitmap Specification: **<http://www.wapforum.org>**

EDGAR DANIELYAN, CISSP, CCNP Security, CCDP®, SCNA, TICSa, CIWCI Security is the principal partner at Danielyan Consulting LLP (**www.danielyan.com**), an information security consultancy in London and Yerevan. He is a published author and editor specialising in UNIX, networking, and information security, having been a cofounder of a national ISP and manager of a country TLD. His book, *Solaris 8 Security*, was published by New Riders Publishing in English and by Pearson Education in Japanese. He is a member of IEEE, IEEE Standards Association, IEEE Computer Society, ACM, ISACA, USENIX, and the SAGE. E-mail: **edd@danielyan.com**

The IETF IPv6 Operations Group and the Development of a Framework for Deployment of IPv6 into IPv4 Networks

by Bob Fink,
Margaret Wasserman, Wind River,
Jun-ichiro Itojun Hagino, IJF

During 2002, the *Internet Engineering Task Force* (IETF) determined that it was best to focus the introduction of IPv6 into the IPv4 Internet by developing deployment scenarios before further development of transition mechanisms without any clearly identified framework for their place in an IPv6 deployment.

Previously the IPv6 Transition working group of the IETF, called *ngtrans* (for IP next-generation transition), was chartered to develop mechanisms and tools to support an IPv6 transition. This work initially focused, in 1995–1996, on the development of the original IPv6 standards, and it led to the basic Transition Mechanism RFC 1933^[1] and later RFC 2893^[2] that defined dual IPv4 and IPv6 protocol stack operation as well as IPv6-over-IPv4 tunnels.

Subsequent attempts to define a framework for transition in 1998–1999 were not successful because there did not appear to be a single vision for a transition to IPv6. Indeed the focus became one of how to have IPv4 and IPv6 coexist for a long period of time, because most felt that a full transition could take well over 10–15 years, with many believing that it would never completely obsolete IPv4. This led to the development of many transition mechanisms and tools, some of which might possibly be more useful than others, that never fit into a coherent framework for operation of a *dual protocol*, that is, IPv4 and IPv6, network.

v6ops

Thus in 2002 the *ngtrans* working group was disbanded, and the IPv6 Operations working group, *v6ops*, created. The *v6ops* working group was chartered to:

- Solicit input from network operators and users to identify operational or security issues with the IPv4/IPv6 Internet, and determine solutions or workarounds to those issues. This includes identifying standards work that is needed in other IETF working groups or areas and working with those groups or areas to begin appropriate work. These issues will be documented in Informational or *Best Current Practice* (BCP) RFCs, or in Internet-Drafts. For example, important pieces of the Internet infrastructure such as the *Domain Name System* (DNS), the *Simple Mail Transfer Protocol* (SMTP), and the *Session Initiation Protocol* (SIP) have specific operational issues when they operate in a shared IPv4/IPv6 network. The *v6ops* working group will cooperate with the relevant areas and working groups to document those issues, and find protocol or operational solutions to those problems.

- Provide feedback to the IPv6 working group regarding portions of the IPv6 specifications that cause, or are likely to cause, operational or security concerns, and work with the IPv6 working group to resolve those concerns. This feedback will be published in Internet-Drafts or RFCs.
- Publish Informational RFCs that help application developers (within and outside the IETF) understand how to develop IP version-independent applications. Work with the Applications area, and other areas, to ensure that these documents answer the real-world concerns of application developers. This includes helping to identify IPv4 dependencies in existing IETF application protocols and working with other areas or groups within the IETF to resolve them.
- Publish informational or BCP RFCs that identify potential security risks in the operation of shared IPv4/IPv6 networks, and document operational practices to eliminate or mitigate those risks. This work will be done in cooperation with the Security area and other relevant areas or working groups.
- Publish Informational or BCP RFCs that identify and analyze solutions for deploying IPv6 within common network environments, such as *Internet Service Provider* (ISP) networks (including core, *Hybrid Fiber-Coaxial* [HFC] or cable, DSL, and dialup networks), enterprise networks, unmanaged networks (home or small office), and cellular networks. These documents should serve as useful guides to network operators and users on how to deploy IPv6 within their existing IPv4 networks, as well as in new network installations.
- Identify open operational or security issues with the deployment scenarios documented in the previous bullet point and fully document those open issues in Internet-Drafts or informational RFCs. Try to find workarounds or solutions to basic, IP-level operational or security issues that can be solved using widely applicable transition mechanisms, such as dual-stack, tunneling, or translation. If the satisfactory resolution of an operational or security issue requires the standardization of a new, widely applicable transition mechanism that does not properly fit into any other IETF working group or area, the v6ops working group will standardize a transition mechanism to meet that need.
- Assume responsibility for advancing the basic IPv6 transition mechanism RFCs along the standards track, if their applicability to common deployment scenarios is demonstrated.

v6ops has started by creating four efforts to define transition scenarios and subsequently to analyze them for potential solutions to the deployment scenarios. These four efforts follow:

- *Third Generation Partnership Project* (3GPP) defined packet networks, that is, *General Packet Radio Service* (GPRS) that would need IP Version 6 deployment into the IPv4 Internet.

- “Unmanaged networks,” which typically correspond to home networks or small office networks.
- ISP networks, including core, HFC or coaxial, DSL, dialup, public wireless, broadband Ethernet, and Internet exchange points.
- Enterprise networks, which are networks that have multiple links and a router connection to an ISP, and are actively managed by a network operations entity.

During 2003 and 2004 it is expected that these deployment scenario efforts will lead to further analysis and identification of deployment solutions and development of appropriate mechanisms to support them.

In addition to this work, serious efforts are under way to engage the entire IETF standards process in the identification and development of appropriate solutions for an IPv6 deployment. One such effort is the *IPv4 Survey* project, which has reviewed the entire IETF RFC catalog of standards to identify what work might need to be done and to disseminate this information to the appropriate area within the IETF.

As progress is made in v6ops, follow-up articles in IPJ will inform you of these efforts.

For Further Reading

- [1] “Transition Mechanisms for IPv6 Hosts and Routers,” R. Gilligan and E. Nordmark, RFC 1933, April 1996.
- [2] “Transition Mechanisms for IPv6 Hosts and Routers,” R. Gilligan and E. Nordmark, RFC 2893, August 2000.
- [3] v6ops IETF information:
<http://www.ietf.org/html.charters/v6ops-charter.html>
- [4] v6ops Web site:
<http://www.6bone.net/v6ops/http://www.6bone.net/v6ops/>

ROBERT FINK is a retired U.S. national laboratory network researcher working with the IPv6 Forum. He is currently a co-chair of the IETF v6ops (IPv6 Operations) working group, and leads the 6bone project. You can reach him at: bob@thefinks.com

MARGARET WASSERMAN is a Principal Technologist at Wind River. She is currently a co-chair of the IETF IPv6 and v6ops working groups. You can reach her at: mrw@windriver.com

JUN-ICHIRO ITOJUN HAGINO is a network researcher with IIJ Research Laboratory. He is currently a co-chair of the IETF v6ops working group and a member of the IETF IAB. You can reach him at itojun@iijlab.net

Opinion: The Mythology of IP Version 6

by Geoff Huston, Telstra

Disclaimer: This is an opinion piece and, therefore, the author takes some liberties in making his points. I hope you as the reader take this in the spirit in which it is intended—a gentle poke at ourselves that sometimes we oversell ourselves and our technology.

In January 1983, the *Advanced Research Projects Agency Network* (ARPANET) experienced a “flag day,” and the Network Control Protocol, NCP, was turned off, and TCP/IP was turned on. Although there are, no doubt, some who would like to see a similar flag day where the world turns off its use of IPv4 and switches over to IPv6, such a scenario is a wild-eyed fantasy. Obviously, the Internet is now way too big for coordinated flag days. The transition of IPv6 into a mainstream deployed technology for the global Internet will take some years, and for many there is still a lingering doubt that will happen at all.

Let’s look more closely at how IPv6 came about, and then look at IPv6 itself in some detail to try to separate the myth from the underlying reality about the timeline for the deployment of IPv6. Maybe then we can suggest some answers to these questions.

IPv6

The effort that has led to the specification of IPv6 is by no means a recently started initiative. A workshop hosted by the then *Internet Activities Board* (IAB) in January 1991 identified the two major scaling issues for the Internet: a sharply increasing rate of consumption of address space and a similar, unconstrained growth of the interdomain routing table. The conclusion reached at the time was that “if we assume that the Internet architecture will continue in use indefinitely, then we need additional [address] flexibility.”

These issues were considered later that year by the *Internet Engineering Task Force* (IETF) with the establishment of the ROAD (*ROuting and ADdressing*) effort. This effort was intended to examine the issues associated with the scaling of IP routing and addressing, looking at the rate of consumption of addresses and the rate of growth of the interdomain routing table. The ultimate objective was to propose some measures to mitigate the worst of the effects of these growth trends. Given the exponential consumption rates then at play, the prospect of exhaustion of the IPv4 Class B space within two or three years was a very real one at the time. The major outcome of the IETF ROAD effort was the recommendation to deprecate the implicit network/host boundaries that were associated with the Class A, B, and C address blocks. In their place the IETF proposed the adoption of an address and routing architecture where the network/host boundary was explicitly configured for each network, and proposed that this boundary could be altered such that two or more network address blocks may be aggregated into a common, single block.

Side Note:

Some would argue that although CIDR was important, it was not the only reason why IPv4 has been able to defy the earlier predictions of its imminent demise. Dynamic *Network Address Translation*, or NAT, allows a network to use a local private address pool to uniquely number its devices, and then translate these private addresses into public addresses to support transactions involving local and external end points. This way, a small pool of public addresses, or even a single address, is used to service a very much larger local private network. It is difficult to estimate the number of devices that are positioned behind NATs, but a highly conservative estimate would see the Internet being at least three times as large as the directly visible part of the Internet.

Side Note:

At an IETF plenary session from that time, the OSI protocol suite was termed the “Road-kill of the Information Superhighway.” It was not completely clear that the presenter made the comment in jest!

This approach was termed *Classless Interdomain Routing*, or CIDR. This was a short-term measure that was intended to buy some time, and it was acknowledged that it did not address the major issue of defining a longer-term, scalable network architecture. But as a short-term measure it has been amazingly successful, given that almost ten years and one Internet boom later, the CIDR address and routing architecture for IPv4 is still holding out.

The IAB, by then renamed the Internet *Architecture* Board, considered the ROAD progress in June 1992, still with its eye on the longer-term strategy for Internet growth. The board’s proposal was that the starting point for the development of the next version of IP would be *Connectionless Network Layer Protocol* (CLNP). This protocol was an element of the *Open System Interconnection* (OSI) protocol suite, with CLNP being defined by the ISO 8473 standard. It used a variable-length address architecture, where network level addresses could be up to 160 bits long. RFC 1347 contained an initial description of how CLNP could be used for this purpose within the IPv4 TCP/IP architecture and with the existing Internet applications. For the IAB this was a bold step, and considering that the IETF community at the time regarded the OSI protocol suite as a very inferior competitor to its own efforts with IP, it could even be termed a highly courageous step. Predictably, one month later in July 1992, at the IETF meeting this IAB proposal was not well received.

The IETF outcome was not just a restatement of architectural direction for IP, but a sweeping redefinition of the respective roles and membership of the various IETF bodies, including that of the IAB.

Of course such a structural change in the composition, roles, and responsibilities of the bodies that collectively make up the IETF could be regarded as upheaval without definite progress. But perhaps this is an unkind view, because the IAB position also pushed the IETF into a strenuous burst of technical activity. The IETF immediately embarked on an effort to undertake a fundamental revision of the Internet Protocol that was intended to result in a protocol that had highly efficient scaling properties in both addressing and routing. There was no shortage of protocols offered to the IETF during 1992 and 1993, including the fancifully named TUBA, as well as PIP, SIPP and NAT.

This effort was part of a process intended to understand the necessary attributes of such a next-generation protocol.

The IETF formed an *Internet Protocol Next Generation (IPng) Directorate* in 1994, and canvassed various industry sectors to understand the broad dimensions of the requirements of such a protocol. This group selected the IPv6 Protocol from a set of proposals, largely basing its selection on the so-called “Simple Internet Protocol,” or SIP proposal. The essential characteristic of the protocol was that of an evolutionary refinement of the Version 4 protocol, rather than a revolutionary departure from Version 4 to an entirely different architectural approach.

Side Note:

IPv6 has had a variety of names—the original IAB documents refer to IP Version 7, working on the assumption that the protocol numbers 5 and 6 were already in use in research networks. It was renamed IPng, for “next generation.”

The final word from the *Internet Assigned Numbers Authority* (IANA) was that protocol number 6 was unused, and the final specification was named Version 6 of the Internet Protocol.

The major strength of IPv6 is the use of fixed-length, 128-bit address fields. Other packet header changes include the dropping of the fragmentation control fields from the IP header, dropping the header checksum and length, and altering the structure of packet options within the header and adding a flow label. But it is the extended address length that is the critical change with IPv6. A 128-bit address field allows an addressable range of 2 to the 128th power, and 2 to the power of 128 is an exceptionally large number. On the other hand, if we are talking about a world that is currently capable of manufacturing more than a billion silicon chips every year, and recognizing that even a one in one thousand address utilization rate would be a real achievement, then maybe it is not all that large a number after all. There is no doubt that such a protocol has the ability to encompass a network that spans billions of devices, which is a network attribute that is looking more and more necessary in the coming years.

Its not just the larger address fields per se, but also the ability for IPv6 to offer an answer to the address scarcity workarounds being used in IPv4 that is of value here. The side effect of these larger address fields is that there is then no forced need to use NAT as a means of increasing the address scaling factor. NAT has always presented operational issues to both the network and the application. NAT distorts the implicit binding of IP address and IP identity and allows only certain types of application interaction to occur across the NAT boundary. Because the “interior” to “exterior” address binding is dynamic, the only forms of applications that can traverse a NAT are those that are initiated on the “inside” of the NAT boundary. The exterior cannot initiate a transaction with an interior end point simply because it has no way of addressing this remote device. IPv6 allows all devices to be uniquely addressed from a single address pool, allowing for coherent end-to-end packet delivery by the network. This in turn allows for the deployment of end-to-end security tools for authentication and encryption and also allows for true peer-to-peer applications.

IPv6, as a protocol architecture, is not a radical departure from the architecture of IPv4. The same datagram delivery model is used, with the same minimal set of assumptions about the underlying network capabilities, and the same decoupling of the routing and forwarding capabilities. The use of an address field in the IP header to contain the semantics of both location and identity was not altered in any fundamental way. The changes made by IPv6 could be seen as conservative set of decisions, based on falling back to the IPv4 protocol model for guidance, on the principle that IPv4 is an operating proof of concept for this architectural approach.

In such a light, IPv6 can be seen as an attempt to regain the advantage of the original IP network architecture: that of a simple and uniform network service that allows maximal flexibility for the operation of the end-to-end application.

It is often the case that complex architectures scale very poorly, and from this perspective the core of IPv6 appears to be a readily scalable architecture.

The Mythology of IPv6

Good as all this is, these attributes alone have not been enough so far to propel IPv6 into broad-scale deployment, and consequently there has been considerable enthusiasm to discover additional reasons to deploy IPv6. Unfortunately, most of these reasons fall into the category of myth, and in looking at IPv6 it is probably a good idea, as well as fair sport, to expose some of these myths as well.

“IPv6 Is More Secure”

A common claim is that IPv6 is more “secure” than IPv4. It is more accurate to indicate that IPv6 is no more or less secure than IPv4. Both IPv4 and IPv6 offer the potential to undertake secure transactions across the network, and both protocols are potentially highly capable in attempting to undertake highly secure transactions. Yes, the IPv6 specification includes as mandatory support for *Authentication and Encapsulating Security Payload* extension headers, but no, there is no “mandatory to use” sticker associated with these extension headers, and, like IPv4 *IP Security* (IPSec), it is left to the application and the user to determine whether to deploy security measures at the network transport level. So, to claim that IPv6 is somehow implicitly superior to IPv4 is an overly enthusiastic claim that falls into the category of “IPv6 myth.”

Now I should qualify this, because there is a distinction between the protocol and its environment of deployment. In the case of IPv4, this protocol capability is compromised in many environments in the face of various forms of deployed active middleware such as NAT. It’s too early to tell with IPv6, but the line of argument is that NAT-based active middleware has been deployed as a means of address extension, and in a IPv6 world such devices are no longer necessary, and will not be deployed. So perhaps one could say that IPv6 enables a path toward widespread peer-to-peer authentication and transport security at the protocol level, but whether the deployment models faithfully follow along such a path remains an open question.

“IPv6 Is Required for Mobility”

It is also claimed that only IPv6 supports mobility. If one is talking about a world of tens of billions of mobile devices, then the larger IPv6 address fields are entirely appropriate for such large-scale deployments. IPv6 includes a developing concept of stateless autoconfiguration and *Neighbor Discovery* mechanisms.

But if the claim is more about the technology to support mobility than the number of mobile devices, then this claim also falls short. The key issue with mobility is that mobility at a network layer requires the network to separate the functions of providing a unique identity for each connected device, and identifying the location within the network for each device.

As a device “moves” within the network its identity remains constant while its location is changing. IPv4 overloaded the semantics of an address to include both identity and locality within an address, and IPv6 did not alter this architectural decision. In this respect, IPv4 and IPv6 offer the same levels of support for mobility. Both protocols require an additional header field to support a decoupled network identity, commonly referred to as the “home address,” and then concentrate on the manner of the way in which the home agent maintains a trustable and accurate copy of the mobile node or current location of the network. This topic remains the subject of activity within the IETF in both IPv4 and IPv6.

“IPv6 Is Better for Wireless Networks”

Mobility is often associated with wireless, and again there has been the claim that somehow IPv6 is better suited for wireless environments than IPv4. Again this is well in the realm of myth.

Wireless environments differ from wireline environments in numerous ways. One of the more critical differences is that a wireless environment may experience bursts of significant levels of bit error corruption, which in turn will lead to periods of non-congestion-based packet loss within the network. A TCP transport session is prone to interpreting such packet loss as being the outcome of network level congestion. The TCP response is not only retransmission of the corrupted packets, but also an unnecessary reduction of the sending rate at the same time. Neither IPv4 nor IPv6 have explicit signaling mechanisms to detect corruption-based packet loss, and in this respect the protocols are similarly equipped, or ill-equipped as in this case, to optimize the carriage efficiency and performance of a wireless communications subnet.

“IPv6 Offers Better QoS”

Another consistent assertion is that IPv6 offers “bundled” support for differentiated *Quality of Service* (QoS), whereas IPv4 does not. The justification for this claim often points to the 20-bit flow label in the IPv6 header as some kind of instant solution to QoS. This claim conveniently omits to note that the flow identification field in the IPv6 header still has no practical application in large-scale network environments. Both IPv4 and IPv6 support an 8-bit traffic class field, which includes the same 6-bit field for differentiated service code points, and both protocols offer the same fields to an *Integrated Services* packet classifier. From this perspective, QoS deployment issues are neither helped nor hindered by the use of IPv4 or IPv6. Here, again, it is a case of nothing has changed.

“Only IPv6 Supports Auto-Configuration”

Another common claim is that only IPv6 offers “plug-and-play” auto-configuration. Again this is an overenthusiastic statement, given the widespread use of the *Dynamic Host Configuration Protocol* (DHCP) in IPv4 networks these days. Both protocol environments support some level of “plug-and-play” auto-configuration capability, and in this respect the situation is pretty much the same for both IPv4 and IPv6.

“IPv6 Solves Routing Scaling”

It would be good if IPv6 included some novel approach that solved, or even mitigated to some extent, the routing scaling issues. Unfortunately, this is simply not the case, and the same techniques of address aggregation using provider hierarchies apply as much to IPv6 as they do to IPv4. The complexity of routing is an expression of the product of the topology of the network, the policies used by routing entities, and the dynamic behavior of the network—not the protocol being routed. The larger address space does little to improve on capability to structure the address space in order to decrease the routing load. In this respect IPv6 does not make IP routing any easier, nor any more scalable.

“IPv6 Provides Better Support for Rapid Prefix Renumbering”

If provider-based addressing is to remain an aspect of the deployed IPv6 network, then one way to undertake provider switching for multihomed end networks is to allow rapid renumbering of a network common prefix. Again, it has been claimed that IPv6 offers the capability to undertake rapid renumbering within a network to switch to a new common address prefix. Again IPv6 performs no differently from IPv4 in this regard. As long as “rapid” refers to a period of hours or days, then yes, IPv4 and IPv6 both support “rapid” local renumbering. For a shorter time frame for “rapid,” such as a few seconds or even a few milliseconds, this is not really the case.

“IPv6 Provides Better Support for Multihomed Sites”

This leads on to the more general claim that IPv6 supports multihoming and dynamic provider selection. Again this is an optimistic claim, and the reality is a little more tempered. Multihoming is relatively easy if you are allowed to globally announce the network address prefix without recourse to any form of provider-based address aggregation. But this is a case of achieving a local objective at a common cost of the scalability of the entire global routing system, and this is not a supportable cost. The objective here is to support some form of multihoming of local networks where any incremental routing load is strictly limited in its radius of propagation. This remains an active area of consideration for the IETF and clear answers, in IPv4 or IPv6, are not available at present. So at best this claim is premature, and more likely the claim will again fall into the category of myth rather than firm reality.

“IPv4 Has Run Out of Addresses”

Again, this is in the category of myth rather than reality. Of the total IPv4 space, some 6 percent is reserved and another 6 percent is used for multicast. Forty-one percent of the space has already been allocated, and the remaining 37 percent (or some 1.5 billion addresses) is yet to be allocated. Prior to 1994, some 36 percent of the address space had been allocated. Since that time, and this includes the entire Internet boom period, a further 15 percent of the available address space was allocated. With a continuation of current policies it would appear that IPv4 address space will be available for many years yet.

So Why IPv6 Anyway ?

The general observation is that IPv6 is not a “feature-based” revision of IPv4—there is no outstanding capability of IPv6 that does not have a fully functional counterpart in IPv4. Nor is there a pressing urgency to deploy IPv6 because we are about to run out of available IPv4 address space in the next few months or even years within what we regard as the “conventional” Internet.

It would appear that the real drivers for network evolution lurk in the device world. We are seeing the various wireless technologies, ranging from Bluetooth for personal networking through the increasingly pervasive IEEE 802.11 “hot-spot” networking to the expectations arising from various forms of *third-generation* (3G) large radius services being combined with consumer devices, control systems, identification systems, and various other forms of embedded dedicated function devices. The silicon industry achieves its greatest advantage through sheer volume of production, and it is in the combination of Internet utility with the production volumes of the silicon industry that we will see demands for networking that encompasses tens, if not hundreds, of billions of devices. This is the world where IPv6 can and will come into its own, and I suspect that it is in this device and utility mode of communications that we will see the fundamental drivers that will lead to widespread deployment of IPv6 support networks.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. Huston is currently the Chief Scientist in the Internet area for Telstra. He is also the Executive Director of the Internet Architecture Board, and is a member of the APNIC Executive Committee. He is author of *The ISP Survival Guide*, ISBN 0-471-31499-4, *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, and coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. All three books are published by John Wiley & Sons. E-mail: gih@telstra.net

Letters to the Editor

SIP Typos Dear Mr. Stallings, and Mr. Jacobsen,

The *Session Initiation Protocol* article by Mr. Stallings in the *Internet Protocol Journal*, Volume 6, Number 1, March 2003, provides an excellent tutorial on the subject, IMHO.

The article does an extraordinary job at presenting what is quite a complicated protocol (SIP) in simple terms. However, there seem to be some typographical errors in the article, which I wanted to bring to your attention:

- In Figure 2, message number 10 should be “180 Ringing” as opposed to “100 Ringing.”
- In Figure 2, the line under message number 14 should be pointing in the opposite direction (that is *from* Bob’s proxy *to* Alice’s proxy).
- In Figure 2, message number 16 should read only “ACK” not “180 ACK.”
- In Figure 2, message number 15 should perhaps read as “200 OK” as opposed to just “OK”
- In Figure 3, message number 5 should read “200 OK” as opposed to “200 Trying”
- Figure 4 message number 5 and 7 should perhaps read as “NOTIFY <Signed In>” as opposed to “<Not Signed-In>”
- Figure 4 “User Agent Bob” should be labelled as “(signed in)” as opposed to “(not signed in)”
- There are missing closing angular brackets in the SIP INVITE message listing on page 27:

To: Bob <sip:bob@biloxi.com>

From: Alice <sip:alice@atlanta.com>;tag=...

- There are missing closing angular brackets in the SIP 200 OK message listing on page 28:

To: Bob <sip:bob@biloxi.com>;tag=....

From: Alice <sip:alice@atlanta.com>;tag=...

Sincerely,

—Rajnish Jain, Excel Switching Corp.
rajnishjain@xl.com

The author responds:

Rajnish,

Thanks for the comments. I am embarrassed that so many errors slipped through, even though I and several reviewers for Ole checked the paper.

—Bill Stallings
ws@shore.net

After reading your article, I couldn't help but notice the U.S. Department of Defense's announcement concerning their intentions to adopt IPv6 in the coming years (see "Fragments," page 38). Given that you've made some strong statements about the value of IPv6 in your article, would you care to offer some views about this announcement?

—Ole

Dear Editor,

As I said in the article, the true value of IP v6 lies in the massive amount of coherent address space that allows literally billions of devices to be uniquely addressed. Address uniqueness is a strong value proposition when you want an identifier space to cover a very large deployment space. As an example of this, one of the two properties of the original Digital-Intel-Xerox Ethernet II specification that remains in today's 10 Gigabit Ethernet specification is unique MAC addresses. All of that highly innovative CSMA/CD thinking that at the time we thought was the fundamental property of Ethernet has been dispensed with.

The general observation is that any communications systems requires any party to be able to uniquely identify any other party in order to initiate a private communication session. If you cannot perform that most basic of communications functions, then you simply do not have a functional peer-to-peer communications network.

But doesn't that mean that the stories of IPv4 address exhaustion have some substance? With the large amount of addressable devices hidden behind NATs, and the associated move to using domain names as the underlying identifier space for many communications applications, the pressure on consumption of IPv4 address space has been reduced considerably. This has implied that in a world of human-driven screens and keyboards we see some considerable lifetime left in the admittedly comfortable world of IPv4 as we know it. To support this model we've actually moved away from the IP address as the unique identifier token for many applications, and substituted an application model that is driven from domain names. As an example, consider the virtual hosting mechanism as implemented in Apache Web servers to see this shift in communications identifiers from address to domain name. And both as consumers of the technology and as an industry we can live with this for some time yet, because we appear to concentrate our use IP addresses as a routing and forwarding framework and increasingly use the DNS as the identifier realm of an application.

But our world is a world where the device is subservient to the user, and the applications we associate with the Internet of today are applications that are essentially human pastimes, such as e-mail, Web browsing, or high-value automated transactions, such as those commonly bracketed into the e-commerce area. And we've now established a highly valuable global industry upon these foundations.

But in so doing we should recognize the emergence of a second set of communications realms populated by uniquely identified devices that number in their billions, where the inter-device traffic is not human mediated, and the value of the device transactions are, on an individual transactions value level, far lower than the value of the human-driven realm of IPv4. In other words, in a device rich communications realm, it's likely that the human value we'd ascribe on average to each packet is far lower than our current Internet IPv4 world of human-mediated communications. And it's this extravagantly device-equipped world that we see the U.S. Department of Defense heading. If your stock in trade is one of quite astounding feats of logistical deployment of large numbers of people and large numbers of items of equipment, then the communications requirement is of a different order of scale to that of the retail Internet markets, and, yes, I'm sure that there are entirely effective arguments behind that decision to look forward to a communications realm with a uniform base protocol identifier domain in a scale that is 2 to the power 96 times larger than the entire IP address identifier domain of IPv4.

But I would be cautious about high levels of expectation that this immediately translates into an impetus in the market where you and I converse. My host here where I'm typing this message is already IPv6 capable, and if you are running a recent version of host software, then it's a reasonable assumption that yours is too. But I'll send this message over IPv4 and you'll receive it over IPv4, and between my mail sender and your mail receiver the transport channel will also be IPv4. Should we use IPv6 instead? Would I pay my provider additional money to compensate it for part of its additional expenditure to support a simultaneous IPv6 capable network between you and me? To send precisely the same message? In precisely the same time? Along the same path? Using the same transport TCP session? Obviously, to me, as a (hopefully) economically rational consumer of such services, and no doubt to you, in a similar role, there is no value in spending more money to achieve outcomes in IPv6 that are identical to what we can already do today in IPv4. And in the retail Internet world that remains the basic IPv6 conundrum. Why should any provider spend additional resources to service the same market with identical services, and in so doing be unable to raise additional revenue to offset their additional service costs? One interpretation is that there is no natural motivation for such activities in today's market, otherwise it would already be very widespread indeed.

What we've seen in the mainstream Internet world is an emerging mythology about IPv6 that somehow this additional expenditure, ultimately on the part of the consumer, provides some additional benefit for the consumer, motivating them to switch from IPv4-only services to some hybrid of mixed v4 and v6 and ultimately to a v6 world, and thereby funding the additional provider expenditure associated with such a massive transition.

The reality is more sobering in that in the retail Internet world there is so far nothing obvious in the "additional benefit" category. I'm using *Network Address Translation* (NAT) right now, using an *ssh* session back to my mail server that drives through NAT boxes to make a secure SMTP session, across a first step of 802.11 wireless in order to send this message to you.

I've auto-configured in the wireless world, and for me I'm living in a plug-and-play world that supports my level of roaming access. Would IPv6 make this session any more secure? Any different in terms of *Quality of Service* (QoS) ? In plug-and-play models of roaming? Would there be any visible difference in terms of my ability to communicate with you? To all of these questions the basic answer is still "no."

So, for you and I, we look inside the IPv6 technology box, and find nothing new there to motivate us to spend more money for our existing Internet-based communications services, and for some time to come it would appear that this will still hold.

On the other hand there are circumstances where there is a need to operate in a much larger base protocol address space. These include situations where one wants to take advantage of Internet applications that operate across a world of literally billions of devices, large and small. The application space may want to gather constant reports on the characteristics of the "thing" it is attached to, from a ration pack to a component of a large naval vessel. You may want to use supply channels for such devices such that the deployment is a plug-and-play world without a massive variety of detailed configuration processes. You may be looking to an architecture that would be stable for many years. In such circumstances you really want take advantage of a uniform set of Internet application technologies that potentially span massive numbers of addressable devices. Here a large base address space is a definite asset. And for such industry sectors in voicing such requirements where there is also a somewhat different ultimate value proposition for the supported communications activity, then it's quite understandable that there can be an attractive proposition offered by immediate adoption of IPv6.

But back in the communications realm where you and I currently exchange our messages, such requirements remain in a future framework that is still waiting for relevant value propositions that allow it to gain traction with you and me. And as I attempted to point out in the article, adding some elements of mythology and over-stating the IPv6 value case won't help here.

Maybe we just need to be patient. Steam ships did not halt operation the first day a diesel powered vessel appeared. It was a much slower process that lead to an outcome of the change of the maritime fleet—the next generation of mechanization offered cheaper services, and, as often happens, market price won in that commodity market.

Market price often wins in competitive commodity markets. And the Internet retail market is, in many parts of the world and in many sectors, a strongly competitive space with all the characteristics of a commodity offering. In addressing such initial specialized dedicated communications requirements with IPv6 technology as represented by the U.S. DoD, there is a distinct possibility that there may be some effective use of initial investment that translates into the retail world in some form of efficiency gain for IPv6-capable providers.

And there no doubt that if you and I could communicate in precisely the same fashion as we do today, with precisely the same applications and service environment, using precisely the same host devices and operating systems as we do today, but at some attractive fraction of today's price, then I'm sure that neither of us would care in the slightest that our data was encapsulated using a packet framing format and address tokens that used the IPv6 protocol specifications.

Kind regards,

—*Geoff Huston, Telstra*
gih@telstra.net

Book Review

Google Hacks *Google Hacks: 100 Industrial-Strength Tips & Tools*, by Tara Calishain and Rael Dornfest, ISBN 0-596-00447-8, O'Reilly & Associates, 2003, 329 pages.

Hmm, this is a hard one. This is the second go at writing a review—the first one made me sound like a grumpy luddite and I don't want my secret identity to be revealed yet. So, put on some suitable music ("So What" from "Kind of Blue" by Miles Davis) and this time, to start with, "just the facts, ma'am" and we'll get back to the grumpiness later.

What we have here are "100 Industrial-Strength Tips & Tools" for using the Google search engine (or g**gling as we are not allowed to say). All the usual O'Reilly positives about layout and presentation apply so we can take those as read (and the usual negative about murky grey scale illustrations). The tips/tools are gathered into separate sections dealing with searching (surprise!), services, scraping, using the API, games and Web mastering. All the tips have some description, some have code and others have URLs that take you to the code or the service described. And indeed some of these are quite interesting and useful, but, and the grumpiness is starting to creep in again, many of them are really not. Tip #1 for instance—"Setting Preferences." Since when has a brief description of how what you can find on the Google preferences page been "Industrial-strength"? Too many of the tips are like this—simple stuff that you can get from many places on the Web (including Google itself) with little added value. Someone starting out using Google is not going to buy a book called *Google Hacks* because its title is off-putting, and someone who is a regular user of the service is going to know (or not be interested in) most of the content. Why do we need a 300 page paper copy of this information? Much of what is in here could be boiled down into a small, cheap guide just like those O'Reilly have for programming languages, and the rest of the stuff is irrelevant anyway (for instance the TouchGraph browser is fun and interesting, but it isn't really that useful—everyone I know has played with it for 5 minutes and then never returned).

I had better hopes of the API programming material, but it was not to be. I know I am in a tiny minority here, so don't complain, but most of the program examples provided in the book use *Perl*. "Hurrah" say you, "Boo" say I—I don't like Perl, never have and never will. Just like celery. I can put up with it, but I won't pick it when I have a choice.

Note, I am not knocking the Google APIs (though they are a bit baroque, and it would be nice to be able to get more than 10 results at time, and...). Being able to call up a search engine from within a program is a good thing, even if you do have to use Web Services (I'm not that keen on them either—are you surprised?). This book certainly tells you how to do that (at least from within Perl) but again you can pick that info up from the Web for free and it doesn't run to more than twenty pages tops. Most of the programming examples may have been fun to write and think up but are about as useful as a flowchart stencil.

Oh, and “Googlehacking”^[1] is not new—people were doing that on AltaVista long (in Internet terms) before Google appeared.

All things considered, I don’t see this book being worth \$25. If you know how to use Google even a little bit you ought to be able to use it to find all this information without it. And what of the stablemate book *Amazon Hacks* which is due to appear soon? I fear a miracle of padding there.

—*Lindsay Marshall, University of Newcastle upon Tyne*
Lindsay.Marshall@newcastle.ac.uk

- [1] Googlehacking is the art of finding a two-word query that has only one result. The two words may not be enclosed in quotes, and the words must be found in Google’s own dictionary (no proper names, made-up words, etc).

Would You Like to Review a Book for IPJ?

We receive numerous books on computer networking from all the major publishers. If you’ve got a specific book you are interested in reviewing, please contact us and we will make sure a copy is mailed to you. The book is yours to keep if you send us a review. We accept reviews of new titles, as well as some of the “networking classics.” Contact us at **ipj@cisco.com** for more information.

Several Landmarks Define Push toward IPv6 Deployment in Japan

In April 1998, the KAME Project, <http://www.kame.net/>, an extension of the WIDE Project (<http://www.wide.ad.jp/>; representative Professor Jun Murai, Keio University), was established with eight core members from seven Japanese vendors. Work began under a two-year timeframe to provide free IPv6/IP Security (IPSec) reference code for UNIX BSD variants. The KAME Project remains active today.

The Japanese government's commitment to taking a leadership role in worldwide IPv6 research and deployment was outlined in a speech to open the September 2000 Diet session by then Prime Minister Mori. Mori identified IPv6 as a key discussion area for the national IT Strategy Council—a strategic pillar toward the “rebirth of the nation.”

The *IPv6 Promotion Council of Japan* was established shortly thereafter, in Oct. 2000. Its founding members numbered only 18. As of March 2003 the Council's membership body consisted of 320 organizations from a variety of business fields; carriers, *Internet Service Providers* (ISPs), hardware vendors, software vendors, finance companies, general trading companies, automobile manufacturers, etc.

The Council is the most active and influential IPv6 organization in Japan, and is the formal contact point appointed by the Japanese government to handle requests from overseas private IPv6 promotion bodies, such as the various regional IPv6 Task Force bodies, for technical and deployment cooperation.

The Promotion Council is currently running the “IPv6 Appli-Contest 2003.” The contest awards developers of applications and software who help to create new possibilities in the IPv6 Internet world, see: <http://www.v6pc.jp/apc/en/concept.html>

Supported by the Ministry of Public Management, Home Affairs, Posts and Telecommunications, and the WIDE Project, the contest is drawing on the cooperation of IPv6 bodies in the EU, North America, India, Korea, Taiwan, and China with the goal of creating a library of freely available IPv6 software.

Details on rules and regulations for entry can be found at the following URL: <http://www.v6pc.jp/apc/en/regulations.html>.

The deadline for entries is August 31, 2003.

Six entries will be selected as “Award of Excellence” winners and will share 1,500,000 JPY in prize money. Award of Excellence winners will also be eligible for the “Grand Prize” of 1,000,000 JPY to be presented at a ceremony during WPC EXPO 2003 to be held September 17–20, 2003, in Tokyo.

An excellent, up-to-date overview of the current status of IPv6 research and commercial service offerings in Japan, including IPv6 case studies and technology tutorials, can be found at IPv6style: <http://www.ipv6style.jp/en/index.shtml>

US Department of Defense adopts IPv6

Implementation of the next-generation Internet protocol that will bring the Department of Defense closer to its goal of net-centric warfare and operations was announced on June 13, 2003 by John P. Stenbit, assistant secretary of defense for networks and information integration and DoD chief information officer.

The new Internet protocol, known as IPv6, will facilitate integration of the essential elements of DoD's Global Information Grid—its sensors, weapons, platforms, information and people. Secretary Stenbit is directing the DoD-wide transition.

The current version of the Internet's operating system, IPv4, has been in use by DoD for almost 30 years. Its fundamental limitations, along with the world-wide explosion of Internet use, inhibit net-centric operations. IPv6 is designed to overcome those limitations by expanding available IP address space, improving end-to-end security, facilitating mobile communications, enhancing quality of service and easing system management burdens.

"Enterprise-wide deployment of IPv6 will keep the warfighter secure and connected in a fast-moving battlespace," Secretary Stenbit said. "Achievement of net-centric operations and warfare depends on effectively implementing the transition."

Secretary Stenbit signed a policy memorandum on June 9 that outlines a strategy to ensure an integrated, timely and effective transition. A key element of the transition minimizes future transition costs by requiring that, starting in October 2003, all network capabilities purchased by DoD be both IPv6-capable and interoperable with the department's extensive IPv4 installed base.

For more information, see:

<http://www.dod.gov/news/Jun2003/d20030609nii.pdf>

<http://www.dod.gov/releases/2003/nr20030613-0097.html>

http://www.dod.gov/news/Jun2003/n06132003_200306134.html

<http://www.dod.gov/transcripts/2003/tr20030613-0274.html>

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Architecture and Technology
MCI, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
The Alfred Fitler Moore Professor of Telecommunication Systems
University of Pennsylvania, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.
Copyright © 2003 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID Cisco Systems, Inc.

The Internet Protocol Journal

September 2003

Volume 6, Number 3

A Quarterly Technical Publication for
Internet and Intranet Professionals

In This Issue

From the Editor	1
Securing BGP: S-BGP	2
Securing BGP: soBGP	15
Virus Trends	23
IPv6 Behind the Wall	34
Call for Papers	40
Fragments	41

FROM THE EDITOR

The task of adding security to Internet protocols and applications is a large and complex one. From a user's point of view, the security-enhanced version of any given component should behave just like the old version, just be "better and more secure." In some cases this is simple. Many of us now use a *Secure Shell Protocol* (SSH) client in place of *Telnet*, and shop online using the secure version of HTTP. But there is still work to be done to ensure that *all* of our protocols and associated applications provide security. In this issue we will look at *routing*, specifically the *Border Gateway Protocol* (BGP) and efforts that are underway to provide security for this critical component of the Internet infrastructure. As is often the case with emerging Internet technologies, there exists more than one proposed solution for securing BGP. Two solutions, S-BGP and soBGP, are described by Steve Kent and Russ White, respectively.

The Internet gets attacked by various forms of viruses and worms with some regularity. Some of these attacks have been quite sophisticated and have caused a great deal of nuisance in recent months. The effects following the *Sobig.F* virus are still very much being felt as I write this. Tom Chen gives us an overview of the trends surrounding viruses and worms.

Closely related to the virus attacks is *spam*. Unfortunately, I know of no complete technical, or even legal, solutions to this growing problem, but I would love to hear your views and solutions. Send your comments to: ipj@cisco.com, but don't use the string "spam" in the subject field or it may get filtered out!

Following Geoff Huston's opinion piece "The Myth of IPv6" in our previous issue, we received a response from *The IPv6 Forum*. The article is entitled "IPv6 Behind the Wall" and is by Jim Bound.

I was very pleased to hear that professor Peter T. Kirstein of University College London had been awarded the Internet Society's *Jonathan B. Postel Service Award* for 2003. I have known Peter since about 1977, when we collaborated on SATNET packet voice conferences between Oslo, London, Boston, and Marina del Rey. Peter is truly an Internet pioneer. (See "Fragments," page 41).

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Securing the Border Gateway Protocol

by Stephen T. Kent, BBN Technologies

Routing in the public Internet is based on a distributed system composed of many routers, grouped into management domains called *Autonomous Systems* (ASes). ASes are operated by *Internet Service Providers* (ISPs) and by multihomed subscribers. (Throughout the remainder of this article, for brevity, we will talk in terms of ISPs, usually omitting references to multihomed subscribers.) Routing information is exchanged between ASes using the *Border Gateway Protocol* (BGP)^[1], via UPDATE messages.

BGP is used in two different contexts. *External BGP* (eBGP) propagates routes between ISPs. BGP also is used within an AS to propagate routes acquired from other ASes. This latter use is referred to as *internal BGP* (iBGP). eBGP is the primary focus of this article, because failures of eBGP can adversely affect large portions of the Internet, well beyond the administrative boundary of the source of the failure. Nonetheless, some ISPs have expressed interest in protecting the distribution of routes within an ISP. The security technology discussed in this article can be used to secure iBGP, but eBGP is the focus of this article. We use the term “BGP” to refer to eBGP throughout the article.

BGP is highly vulnerable to a variety of attacks^[2]. In some cases, this vulnerability arises because of a lack of integrity and authentication for BGP messages. However, the more substantive and harder problem is the lack of a secure means of verifying that BGP traffic is authorized, a concept explored in more detail in this article. In April 1997, BBN began work on the security architecture described here, a system we refer to as *S-BGP*, to address the vulnerabilities of BGP. This article begins by reviewing the problem, discusses a model for correct operation of BGP, presents a threat model, and states the goals and assumptions that underlie our proposed security architecture.

Before we begin the discussion of BGP in more detail, a few definitions are in order. A *route* is defined as an *address prefix* and a set of *path attributes*. One of the path attributes is an AS path, and that is the primary focus of BGP security considerations. The AS path specifies the sequence of ASes that subscriber traffic should traverse if forwarded via this route. When propagating an UPDATE to a neighboring AS, the BGP router prepends its AS number to the sequence, and may update certain other path attributes. The first AS included in the path is referred to as the *origin AS*.

Each BGP router (other than at the edges of the Internet) maintains a complete routing table, capable of routing traffic to any reachable destination, and sends its best route for each prefix to each neighbor. In BGP, “best” is very locally defined. The BGP route selection algorithm has few criteria that are universal, thus limiting the extent to which any security mechanism can detect and reject “bad” routes emitted by a neighbor.

Each ISP makes use of local policies that it need not disclose, and this gives BGP route selection a “black box” flavor, which has significant adverse implications for security.

Correct Operation of BGP

Security for BGP should be defined as the correct operation of BGP routers. This definition is based on the observation that any successful attack against BGP will result in other than correct operation, presumably yielding degraded routing. Correct operation of BGP depends upon the integrity, authenticity, and timeliness of the routing information it distributes, as well as each BGP router processing, storing, and distributing this information in accordance with both the BGP specification and local routing policies. Many statements could be made in an effort to characterize correct operation, but they rest on two simple assumptions.

First, control (vs. subscriber traffic) communication between neighbor BGP routers must be authenticity and integrity secure. This is easily achieved through the use of a point-to-point security protocol capable of protecting BGP traffic; for example, *IP Security* (IPSec). Second, BGP routers must execute the route selection algorithm correctly and communicate the results. There are two parts to this assumption: processing received UPDATES, and generation and transmission of UPDATES. In terms of an AS trying to protect itself against external attacks, correct operation of its own BGP routers is mostly a local security issue, but not an Internet-wide security issue. However, an AS should not rely on other ASes to operate properly; such reliance permits a failure in one AS to propagate to others, a domino failure effect. Thus it is important for a BGP router to be able to verify that each UPDATE it receives from a peer is valid (authorized) and timely.

The validity of an UPDATE message is based on four primary criteria:

- The router that sent the UPDATE was authorized to act on behalf of the AS it claims to represent; that is, the AS at the front of the AS path.
- The AS from which the UPDATE emanates was authorized by the preceding AS in the AS path (in the UPDATE message) to advertise the prefixes in the UPDATE.
- The first AS in the AS path was authorized, by the owner of the set of prefixes that are represented in the UPDATE, to advertise those prefixes.
- If the UPDATE withdraws one or more routes (specified by the prefixes for the routes), then the sender must have advertised each route prior to withdrawing it.

There are some limitations to the ability of any practical security mechanism to detect all BGP security failures. The local policy feature of BGP allows each ISP considerable latitude in how UPDATES are processed, making it difficult for an external observer—for example, a router in a neighboring AS—to determine if a router is operating properly.

This is because such behavior might be attributed to local policies not visible outside an AS. To address such attacks, the semantics of BGP itself would have to change. Moreover, because UPDATEs do not carry sequence numbers, a BGP router can emit an UPDATE based on authentic, but old, information; for example, withdrawing or reasserting a route based on outdated information. Thus the temporal accuracy of UPDATEs, in the face of Byzantine failures, is hard to enforce, except in a very coarse fashion. (Simply speaking, a *Byzantine failure* is one in which a nominally trusted or authorized entity misbehaves.)

Threat Model and BGP Vulnerabilities

Routers exhibit both architectural and implementation vulnerabilities. Implementation vulnerabilities are the result of errors that arise in developing design details or coding; for example, translating the BGP specs into software. Architectural vulnerabilities permit various forms of attack, independent of implementation details, and thus are potentially more damaging, because they persist across all implementations. To make Internet routing robust, both forms of vulnerabilities must be addressed. BGP vulnerabilities can be exploited to cause improper routing or nondelivery of subscriber traffic, network congestion, and traffic delays. Misrouting attacks can be used to facilitate both passive and active wiretapping of subscriber traffic. Often an attack against BGP may be part of a larger attack against subscriber computers. For example, there have been BGP attacks that seek to misroute queries to *Domain Name System* (DNS) root servers, as part of an attack against subscriber systems.

BGP can be attacked in many ways. Communication between BGP peers can be subjected to active or passive wiretapping. The BGP software, configuration information, or routing databases of a router may be modified or replaced via unauthorized access to a router, or to a server or management workstation from which router software is downloaded. These latter attacks transform routers into hostile insiders, so security measures must address such Byzantine failures.

Improved physical and procedural security for network management facilities, and routers, and cryptographic security for BGP traffic between routers would help reduce some of these vulnerabilities. However, physical and procedural security is expensive and imperfect, and these countermeasures would not protect the Internet against accidental or malicious misconfiguration by operators, nor against attacks that mimic such errors. Misconfiguration of this sort has been a source of Internet outages in the past and seems likely to persist. Any security approach that relies on ISPs to act properly violates the “principle of least privilege” and leaves the Internet routing system vulnerable at its weakest link. In contrast, the security approach described in this article satisfies this principle, so that any attack on any component of the routing system is limited in its impact on the Internet as a whole.

Routers also are susceptible to resource exhaustion attacks based on delivery of large quantities of management traffic, BGP or otherwise. This vulnerability arises because these devices are designed with the not unreasonable model that management traffic is a very tiny percentage of all the traffic that arrives at a router. Router interfaces can deliver traffic to the management processor at very high rates, because they are designed to accommodate subscriber traffic flows. Solutions to this problem need to be generic, to accommodate all types of router management traffic, and thus are outside the scope of the BGP security measures discussed in this article.

Goals, Constraints, and Assumptions

Any proposed security architecture must exhibit dynamics consistent with the existing BGP system; for example, responding automatically to topology changes, including the addition of new networks, routers, and ASes. These actions take place on different time scales and have different scopes. For example, in the current BGP system, if an ISP replaces a failed router, the action can take place fairly quickly and has only local impact, because ISPs are not aware of the identity of routers in other, non-neighboring, ISPs. The issuance of new AS numbers, representing new nets, is not a fast process, nor is the allocation of new blocks of address space (new prefixes). But both of these actions are globally visible. Changes in routes also may have global impact, and they may occur very quickly.

Solutions also must scale in a manner consistent with the growth of the Internet. The countermeasures must be consistent with the BGP protocol standards and with the likely evolution of these standards. This includes packet size limits and features such as path aggregation, communities, and multiprotocol support (for example, *Multiprotocol Label Switching* [MPLS]). The security measures must be incrementally deployable; there cannot be a “flag day” when all BGP routers suddenly begin executing a new security protocol. It is desirable to not create new organizational entities that must be accepted as authorities by ISPs and subscribers, in order to make routing secure.

S-BGP Architecture

S-BGP consists of four major elements:

- A *Public Key Infrastructure* (PKI) that represents the ownership and delegation of address prefixes and AS numbers
- *Address Attestations* that the owner of a prefix uses to authorize an AS to originate routes to the prefix
- *Route Attestations* that an AS creates to authorize a neighbor to advertise prefixes
- *IPSec* for point-to-point security of BGP traffic transmitted between routers

These elements are used by an S-BGP router to secure communication with neighbors, and to generate and validate UPDATE messages relative to the authorization model represented by the PKI and address attestations. Together, the combination of these security mechanisms prevents a compromised AS from propagating erroneous routing data to other, secured ASes. Each element is described in more detail in the following section.

S-BGP Public Key Infrastructure

S-BGP uses a PKI based on X.509 (v3) certificates to enable routers to validate the authorization of other routers to represent ASes (ISPs). The PKI also allows routers to verify the authorization of each ISP as the owner of one or more prefixes (contiguous blocks of address space). This PKI was described in^[14], and the reader is referred to that paper for additional details. The PKI parallels the existing IP address and AS number assignment delegation system and takes advantage of this infrastructure. Because the PKI mirrors existing infrastructure, it avoids most of the “trust” issues that often complicate the creation of a PKI. This PKI is unusual in that it emphasizes authorization, not authentication. The names used in the certificates in this PKI are not employed to determine whether a given ISP or router is authorized to do anything, and the names are not even meaningful outside of S-BGP.

S-BGP calls for a certificate to be issued to each ISP (or subscriber) that owns (more properly, has a right to use) a portion of the IP address space. This certificate is issued through the same procedures employed for address allocation, starting with the *Internet Assigned Numbers Authority* (IANA) and continuing through a *Regional Internet Registry* (RIR), and, if applicable, an ISP. If an ISP owns multiple prefixes, we issue a single certificate containing a list of prefixes, to minimize the number of certificates in the system. The PKI represents address-space ownership by binding prefixes to a public key belonging to the ISP to which the prefixes have been assigned. Each certificate contains a private extension that specifies the set of prefixes that has been allocated to the ISP. Certificates issued under this PKI also represent the binding between an ISP and the AS numbers allocated to it. The PKI allows each ISP to issue certificates to its routers, certifying that these routers represent the ISP and hence, the ASes owned by the ISP. Here too, the PKI parallels the existing AS allocation system; that is, the IANA allocates AS numbers to RIRs, which in turn assign AS numbers to ISPs that run S-BGP.

Attestations

An *attestation* is a digitally signed datum asserting that its target (an AS) is authorized by the signer (an ISP) to advertise a path to one or more specified prefixes. There are two types of attestations, address and route, which share a common format. For an *Address Attestation* (AA), the signer is the ISP or subscriber that controls the prefixes in the AA, and the target is a set of ASes that the ISP/subscriber authorizes to originate a route to the prefixes. AAs are relatively static data items, because relationships between address-space owners and ISPs change relatively slowly.

For a *Route Attestation* (RA), the signer is an S-BGP router (operating on behalf of an ISP), and the target is an AS or set of ASes, representing the neighbors to which the UPDATE containing the RA will be sent. RAs, unlike AAs, are very dynamic, possibly changing for each transmitted UPDATE.

UPDATE Validation

Attestations and certificates are used by S-BGP routers to validate routes asserted in UPDATE messages; that is, to verify that the first AS in the route has been authorized to advertise the prefixes by the prefix owner(s), and that each subsequent AS has been authorized to advertise the route for the prefixes by the preceding AS in the route. To validate a route received from AS_n , AS_{n+1} requires:

- An AA for each organization owning a prefix represented in the UPDATE (not for prefixes in the UPDATE that represent routes being withdrawn)
- A certified public key for each organization owning a prefix in the UPDATE
- An RA corresponding to each AS along the path (AS_n to AS_1), where the RA generated and signed by the router in AS_n encompasses the *Network Layer Reachability Information* (NLRI) and the path from AS_{n+1} through AS_1
- A certified public key for each S-BGP router that signed an RA along the path (AS_n to AS_1), to check the signatures on the corresponding RAs

An S-BGP router verifies that the advertised prefixes and the origin AS are consistent with AA information. The router verifies the signature on each RA and verifies the correspondence between the signer of the RA and the authorization to represent the AS in question. There also must be a correspondence between each AS in the path and an appropriate RA. If all of these checks pass, the UPDATE is valid.

AAs are not used to check withdrawn routes in an UPDATE. Use of IP-Sec to secure communication between each pair of S-BGP routers, plus the fact that BGP uses a separate *Adjacency Routing Information Base* (Adj-RIB-In) for each neighbor, ensures that only the advertiser of a route can withdraw it.

Distribution of S-BGP Data

Each S-BGP router must have the public keys required to validate the RAs in UPDATES, a scenario that translates into securely distributed keys for every router that implements S-BGP (and that is reachable via an S-BGP path). Each router also needs access to all AA information, to verify that the origin AS is authorized to originate a route to the prefixes in the UPDATE. S-BGP does not distribute certificates, *Certificate Revocation Lists* (CRLs), or AAs via UPDATE messages; transmission of these items via UPDATES would be very wasteful of bandwidth, because each BGP router would receive many redundant copies from its neighbors.

Also, an UPDATE is limited to 4096 bytes and thus generally could not carry all of this data for the route represented by the UPDATE. Instead, S-BGP distributes this data to routers via out-of-band means. The data is relatively static and thus is a good candidate for caching and incremental update. Moreover, the certificates and AAs can be validated and reduced to a more compact format by ISP operation centers prior to distribution to routers. This avoids the need for each router to perform this processing, saving both bandwidth and storage space. It also means that routers do not need to be able to parse X.509 certificates and validate certificate paths for S-BGP purposes, although some capability in this area may be required for IPsec key management.

S-BGP uses *repositories* for distribution of this data. We initially described a model in which a few replicated, loosely synchronized repositories were operated by the RIRs. Discussions with ISPs suggest a model in which major ISPs and Internet exchanges operate repositories, and smaller ISPs and subscribers make use of these repositories. In either model, each ISP periodically, for example daily, uploads new/changed certificates, its current CRL, and AAs. Each ISP also downloads all of this data for all other ISPs that are running S-BGP. The repositories periodically transfer new data to one another to maintain loose synchronization. ISPs process the repository information to create more compact files that contain the AA data and the public keys and prefix and AS data from the certificates, but none of the certificate management information or CRLs. These resulting “extracted” files are transferred to the routers executing S-BGP under the control of the ISP.

Because certificates, AAs, and CRLs are signed and carry validity interval information, they require minimal additional security while in transit to or from a repository or while stored on a repository. Nonetheless, S-BGP employs the *Secure Sockets Layer* (SSL) protocol, with both client and server certificates, to protect access to the repositories, as a countermeasure to denial-of-service attacks. The simple, hierarchic structure of the PKI allows repositories to automatically effect access control checks on the uploaded data, for example, to prevent one ISP from accidentally or maliciously overwriting the certificates, CRLs, and AAs from another ISP.

Distribution of Route Attestations

S-BGP distributes RAs with BGP UPDATES in a newly defined, optional, *transitive path attribute*. Because routes may change quickly, it is important that RAs accompany the UPDATES that are validated using them. If any other means of distribution is employed for this data, there is a likelihood that the UPDATES and the data will be out of synch, creating a conundrum for a router; that is, what should the router do when the UPDATE and the security data differ? RAs employ a compact encoding scheme to help ensure that they fit within the BGP packet size limits, even when route or address aggregation occurs. (S-BGP accommodates aggregation by explicitly including signed attribute data that otherwise would be lost when aggregation occurs.) An S-BGP router receiving an UPDATE from a peer caches the RAs with the route in the Adj-RIB for the peer, and in the *Local Routing Information Base* [Loc-RIB] (if the route is selected).

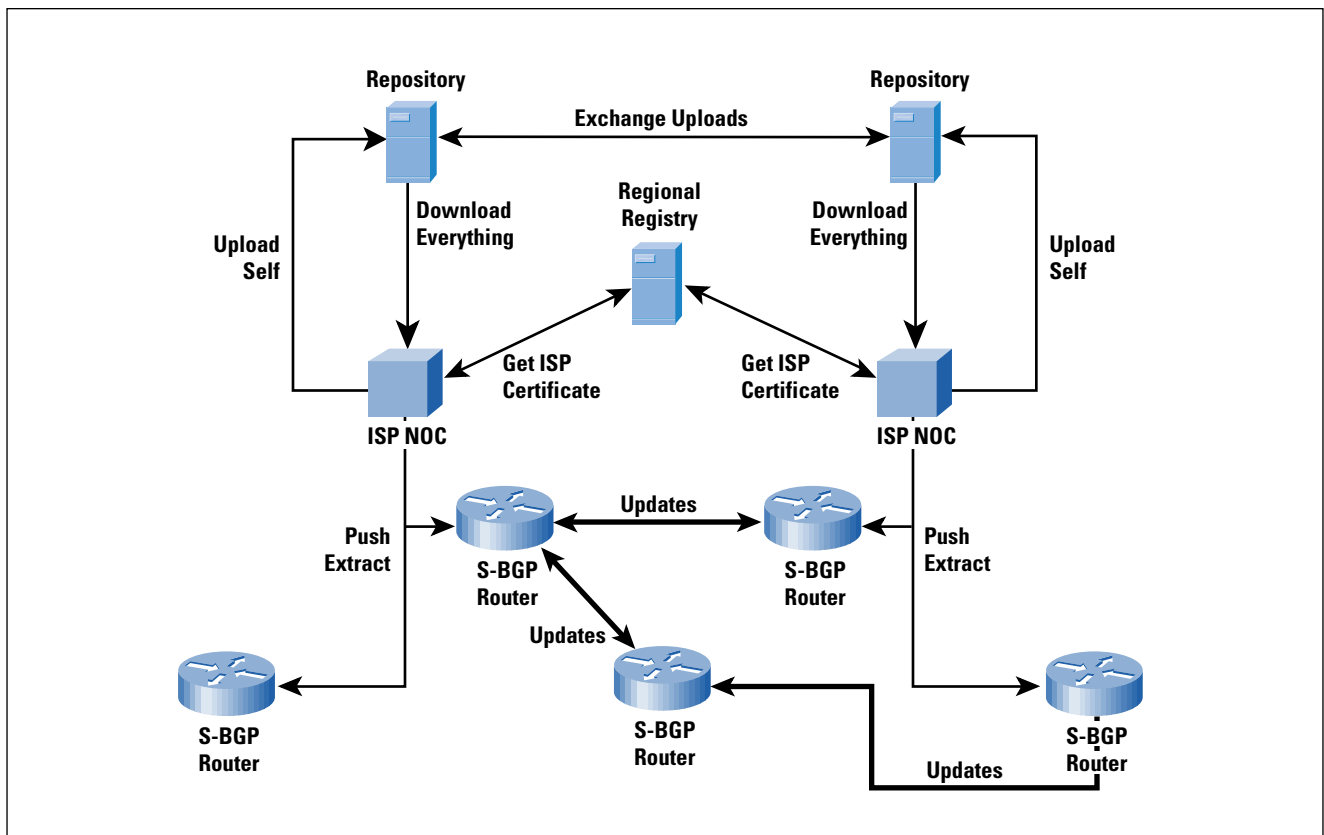
As noted in the following discussion, the bandwidth required to support in-band distribution of route attestations is negligible (compared to subscriber traffic).

Although the RA mechanism was designed to protect AS path data, it can also accommodate other new path attributes; for example, communities^[11] and confederations^[12]. Specifically, there is a provision to indicate what data, in addition to the AS path, is covered by the digital signature that is part of the RA.

Putting It All Together

Figure 1 illustrates how the major elements of S-BGP interact, using a simplified example. The figure shows two ISPs, each with a *Network Operations Center* (NOC), a repository, and three routers. A third ISP is represented by a single (S-BGP-enabled) router. Each ISP interacts with an RIR to acquire a certificate representing the prefixes and AS numbers assigned to the ISP. Each NOC interacts with a repository to upload data (certificates, CRLs, and AAs) from that ISP, and to download the same data acquired from all other ISPs. The repositories interact with one another to exchange uploaded ISP data, to make that data available to all other ISPs. Within an ISP, the NOC pushes a copy of the extracted certificate and AA data, produced from the downloads acquired from a repository, to each router. Routers exchange UPDATE messages, containing RAs, that enable validation of each received UPDATE.

Figure 1: S-BGP Element Interactions



IPSec and Router Authentication

S-BGP uses IPSec^[6,7,8], specifically the *Encapsulating Security Payload* (ESP) protocol, to provide authentication, data integrity, and antireplay for all BGP traffic between neighboring routers. The *Internet Key Exchange* (IKE) protocol^[9,10] is used for key management services in support of ESP. The S-BGP PKI includes certificates for IKE, separate from those used for RA processing.

The use of IPSec is preferable to the current option of the *Message Digest Algorithm 5* (MD5) TCP checksum option^[15], in several respects. IPSec uses keyed hash functions in a way that is cryptographically more secure than the MD5 checksum option, and IKE provides automated key management, a feature sorely lacking in the option. Protecting BGP traffic at the IP layer, vs. the TCP layer, counters more vulnerabilities, because the TCP implementation is protected as well, for example, including SYN flooding and spoofed RSTs (resets), are rejected.

Residual Vulnerabilities in S-BGP

Despite the extensive security offered by S-BGP, architectural vulnerabilities exist that are not eliminated by its use. For example, an S-BGP router may reassert a route that was withdrawn earlier, even if the route has not been readvertised. The router also may suppress UPDATES, including ones that withdraw routes. These vulnerabilities exist because BGP UPDATES do not carry sequence numbers or time stamps that could be used to determine their timeliness. However, RAs do carry an expiration date and time, so there is a limit on how long an attestation can be misused this way. S-BGP restricts malicious behavior to the set of actions for which a router or AS is authorized, based on externally verifiable, authoritative constraints.

Performance and Operational Issues

In developing the S-BGP architecture, we paid close attention to the performance and operational impact of the proposed countermeasures, and reported our analysis in earlier papers. In preparing this article, we updated our data, utilizing a variety of sources; for example, the *Route Views* project. Although much data about BGP and associated infrastructure is available, other data is difficult to acquire in a fashion that is representative of a “typical” BGP router. This is because each AS in the Internet embodies a slightly different view of connectivity, as a result of local policy filters applied by other ASes.

It is important that the transmission, storage, and processing requirements imposed by S-BGP not be so great as to overwhelm routers. Each of these requirements must be analyzed separately.

The transmission of RAs in UPDATES does significantly increase the size of these messages, by about 800 percent. However, because the volume of this traffic is minuscule relative to subscriber traffic, the increase is negligent. The set of files containing certificates, AAs, and CRLs would be about 75–85 MB. Daily transmission of these files between ISPs and repositories would not represent a significant increase in traffic volume for the Internet.

Although the transmission overhead is not a concern, storage of the RAs in each Adj-RIB and the Loc-RIB is a problem. The additional space required to hold these RAs is estimated at about 30–35 MB per peer, if S-BGP were fully deployed today. This is a modest amount of memory for a typical router with a few peers, but a significant amount of storage for routers at Internet exchanges, where a router may have tens or even hundreds of peers.

Thus the management CPU in a router might need a gigabyte or more of RAM under these conditions. (When a large ISP peers with many other ISPs at an exchange, the peering is not symmetric; that is, the large ISP accepts only a few routes from each of the smaller ISPs, filtering out the rest. Thus the amount of additional memory required for RAs in Adj-RIBs for each of these small ISP peers may be considerably less than for symmetric peer relationships.) This requisite memory seems modest by current workstation standards, but most deployed routers cannot be configured with this much memory.

The computational burden of router processing of RAs in UPDATES is a function of the path length in each UPDATE and the rate at which UPDATES arrive. The arrival rate is a function of the number of S-BGP peers the router sees, and the rate at which each peer sends UPDATES. Our analysis suggests that the long-term (24-hour) UPDATE rate for a router with 30 peers is about 0.5 UPDATES per second. On average, each UPDATE would contain about 3.7 RAs. We originally estimated the busy minute rate as about 10 times the average rate. At this rate, a router could probably perform the requisite signature verification in software (about 18 signature verifications per second). Recent evidence suggests a factor of 100–200 might be a better estimate, in light of experience with major worm attacks, and at that rate it would be hard for software to keep pace.

Heuristics are available to reduce this burden. Analysis shows that about 50 percent of all UPDATES are sent as a result of route “flaps”; that is, transient communication failures that, when remedied, result in a return to the former route. Thus if a router maintained a depth-two cache for each Adj-RIB-In, it could avoid signature validation about 50 percent of the time. However, this would double the storage requirements for these RIBs, and that would exacerbate the storage problem cited previously.

Our previous analysis also assumed that receipt of each UPDATE would result in transmission of an UPDATE with one new signature. This was an oversimplification; a router generates and transmits an UPDATE only if the newly received route is “better” than the current best route (for the prefix), or if the best route for the prefix is withdrawn by the UPDATE. When a router has many peers, most of the UPDATES it receives may not yield a better route, and thus will not trigger transmission of a new UPDATE.

On the other hand, when a router does select a new route, an UPDATE may be constructed and sent to each neighbor, requiring one signature per neighbor. This is because an RA specifies the AS number of the neighbor to which it is directed. It is possible to construct an RA that identifies the next hop as a set of AS numbers, corresponding to all the neighbors to which an UPDATE is authorized to be sent. The downside of this strategy is that it makes the RAs larger, contributing to the storage problem noted previously.

The observation made previously suggests a heuristic for UPDATE processing to mitigate signature validation costs. A router can defer validation of the RAs in any UPDATE that it receives, if the UPDATE would not represent a new best route. This optimization could be especially helpful for routers that receive the greatest number of UPDATES; that is, routers with many neighbors. One might worry that this strategy allows an attacker to force processing, by sending what would be considered “very good” routes, but an S-BGP router could detect such fraudulent UPDATES and could choose to drop its connection to a peer that behaved this way, in order to counter such an attack.

Initialization/reboot of a BGP router also results in a surge in UPDATE processing, and the deferred processing heuristic is applicable here too, even though reboots are relatively infrequent. Saving RIBs in nonvolatile storage addresses this problem. Most deployed routers do not have sufficient nonvolatile storage to adopt this strategy, but some do have hard drives that would easily accommodate the RIBs.

It is reasonable to assume that next-generation routers could be configured with enough RAM for the RIBs, but this analysis shows that full deployment is not feasible with the currently deployed router base. To add RAM, and possibly to add nonvolatile storage, router vendors will have to upgrade the processor boards where net management processing takes place. That suggests that addition of a crypto accelerator chip would be prudent as part of the board redesign process, for example, to deal with surge conditions noted previously.

Deployment and Transition Issues

Adoption of S-BGP requires cooperation among several groups. ISPs and subscribers running BGP must cooperate to generate and distribute AAs. Major ISPs must implement the S-BGP security mechanisms in order to offer significant benefit to the Internet community. The IANA and RIRs must enhance operational procedures to support generation of prefix and AS number allocation certificates. Router vendors need to offer additional storage in next-generation products, or offer ancillary devices for use with existing router products, and revise BGP software to support S-BGP.

There is some good news; S-BGP can be deployed incrementally. Only neighboring ASes receive full benefit from such deployment. Although we chose a transitive path attribute syntax to carry RAs, and thus it might be possible for non-neighbor ASes to exchange RAs, it seems likely that intervening ASes would not have sufficient storage for the RAs in their RIBs.

Also, the controls needed in routers to take advantage of noncontiguous deployment of S-BGP are quite complex, hence our suggestion that only contiguous deployment of S-BGP be attempted.

External routes received from S-BGP peers need to be redistributed within the AS, both to interior routers and to other border routers, in order to maintain a consistent and stable view of the exterior routes across the AS. Thus an AS must switch to using S-BGP for all its border routers at once, to avoid route loops within the AS.

Status

As of early 2003, an implementation of S-BGP has been developed and demonstrated on small numbers of workstations representing small numbers of ASes. We also developed software for a simple repository, and for NOC tools that support secure upload and download of certificates, CRLs, and AAs to and from repositories, and for certificate management for NOC personnel and routers. This suite of software, plus CA software from another *Defense Advanced Research Projects Agency* (DARPA) program, provide all of the elements needed to represent a full S-BGP system. All of this software is available in open source form.

Summary

S-BGP represents a comprehensive approach to addressing a wide range of security concerns associated with BGP. It detects and rejects unauthorized UPDATE messages, irrespective of the means by which they arise; for example, misconfiguration, active wiretapping, compromise of routers or management systems, etc. S-BGP is not perfect; it has a few residual vulnerabilities, but these pale in comparison to the security features S-BGP provides, and removal of these vulnerabilities would require more fundamental changes to BGP semantics.

The S-BGP design is based on a top-down security analysis, starting with the semantics of BGP and factoring in the wide range of attacks that have or could be launched against the existing infrastructure.

Acknowledgements

Many individuals contributed to the design and development of S-BGP, including Christine Jones, Charlie Lynn, Joanne Mikkelsen, and Karen Seo.

References

- [1] Y. Rekhter and T. Li, "A Border Gateway Protocol 4 (BGP-4)," RFC 1771, March 1995.
- [2] S. Kent, C. Lynn, and K. Seo, "Secure Border Gateway Protocol (S-BGP)," *IEEE Journal on Selected Areas in Communications*, Vol. 18, No. 4, April 2000.
- [3] C. Villamizar, R. Chandra, and R. Govindan, "BGP Route Flap Damping," RFC 2439, November 1998.

- [4] B.R. Smith, and J.J. Garcia-Luna-Aceves, "Securing the Border Gateway Routing Protocol," Proceedings of Global Internet '96, November 1996.
- [5] S. Murphy, panel presentation on "Security Architecture for the Internet Infrastructure," Symposium on Network and Distributed System Security, April 1995.
- [6] S. Kent and R. Atkinson, "Security Architecture for the Internet Protocol," RFC 2401, November 1998.
- [7] R. Glenn and S. Kent, "The NULL Encryption Algorithm and Its Use with IPsec," RFC 2410, November 1998.
- [8] S. Kent and R. Atkinson, "IP Encapsulating Security Payload (ESP)," RFC 2406, November 1998.
- [9] D. Maughan, M. Schertler, M. Schneider, and J. Turner, "Internet Security Association and Key Management Protocol (ISAKMP)," RFC 2408, November 1998.
- [10] D. Harkins and D. Carrel, "The Internet Key Exchange (IKE)," RFC 2409, November 1998.
- [11] R. Chandra, P. Traina, and T. Li, "BGP Communities Attribute," RFC 1997, August 1996.
- [12] P. Traina, "Autonomous System Confederations for BGP," RFC 1965, June 1996.
- [13] T. Bates, R. Chandra, D. Katz, and Y. Rekhter, "Multiprotocol Extensions for BGP-4," RFC 2283, February 1998.
- [14] K. Seo, C. Lynn, and S. Kent, "Public-Key Infrastructure for the Secure Border Gateway Protocol (S-BGP)," DARPA Information Survivability Conference and Exposition, June 2001.
- [15] A. Heffernan, "Protection of BGP Sessions via the TCP MD5 Signature Option," RFC 2385, August 1998.

STEPHEN KENT received the S.M., E.E., and Ph.D. degrees in computer science from MIT, and a B.S. in mathematics from Loyola University of New Orleans. He has worked at BBN for over 25 years, where he serves today as Chief Scientist-Information Security. He served on the IAB for over a decade, and chaired the Privacy & Security Research Group of the IRTF and the PEM WG in the IETF, where he currently co-chairs the PKIX WG. He has served on several committees for the National Research Council, and chairs a committee on authentication and privacy for the NRC. His current work focuses on PKI issues, BGP security, and very high speed IP encryption. He is a Fellow of the ACM, and a member of the Internet Society and Sigma Xi. His e-mail address is: kent@bbn.com

Securing BGP Through Secure Origin BGP

by Russ White, Cisco Systems

Networks have come under increasing scrutiny in the area of security. Routing, the part of the network that provides information on how to reach destinations within the network, has been gaining attention from a security perspective as well. *The Internet Engineering Task Force* (IETF) has, in fact, formed a new working group, the *Routing Protocols Security Requirements Working Group* (<http://www.rpsec.org>), to analyze security in routing systems.

Of course, the biggest network in existence is the Internet, and the routing protocol that provides reachability and path information for the Internet is the *Border Gateway Protocol* (BGP), specified in RFC 1771. Several methods of securing the information carried within BGP have been proposed:

- *Internet Route Verification* (IRV), described in “Working Around BGP: An Incremental Approach to Improving Security and Accuracy of Interdomain Routing,” Symposium on Network and Distributed Systems Security, February 2003, by Geoffrey Goodell, William Aiello, Timothy Griffin, John Ioannidis, Patrick McDaniel, and Aviel Rubin. IRV relies on out-of-band communication with a route originator to verify the correctness of a route.
- S-BGP, described in the companion article and at: www.net-tech.bbn.com/projects/s-bgp
- *Domain Name System* (DNS)-based *Network Layer Reachability Information* (NLRI) origin *Autonomous System* (AS) verification in BGP, which is the oldest attempt at validating the information carried within BGP, is described in [draft-bates-bgp4-nlri-origin-verif-00.html](#),

This article discusses *Secure Origin BGP* (soBGP), a solution recently proposed by a group (including me) mostly within Cisco Systems. We believe soBGP to be a deployable mechanism for validating the correctness and authorization of the data carried within BGP, and also for preventing the sorts of attacks resulting from misconfiguration or intentional insertion of bad data into the Internet routing system.

We address four goals when we consider security in terms of BGP:

- Is the AS originating the destination (prefix) authorized to advertise it? In other words, if a router receives an advertisement for the 10.1.1.0/24 network originating in AS65500, is there any way to verify that AS65500 is supposed to be advertising 10.1.1.0/24?
- Does the AS advertising the destination actually have a path to the destination? In other words, if a router is receiving an advertisement from a BGP peer in AS65501 that it can reach 10.1.1.0/24, is there any way to verify that AS65501 actually has a path to the AS originator 10.1.1.0/24?

- Is the peer advertising the route authorized by the originator, or owner, of the destination, to advertise a path to the destination?
- Does the path advertised by a peer AS fall within the policies the local network administrators have set forward? The most obvious issue is whether or not the AS Path advertised by the peer is an acceptable path to send the traffic along.

We argue elsewhere that the second two goals cannot be fully met within an operational internetwork, for many reasons; see **draft-white-pathconsiderations-00.txt** for further discussion on this point. In this article, then, we discuss how soBGP can meet the first two goals in operational networks.

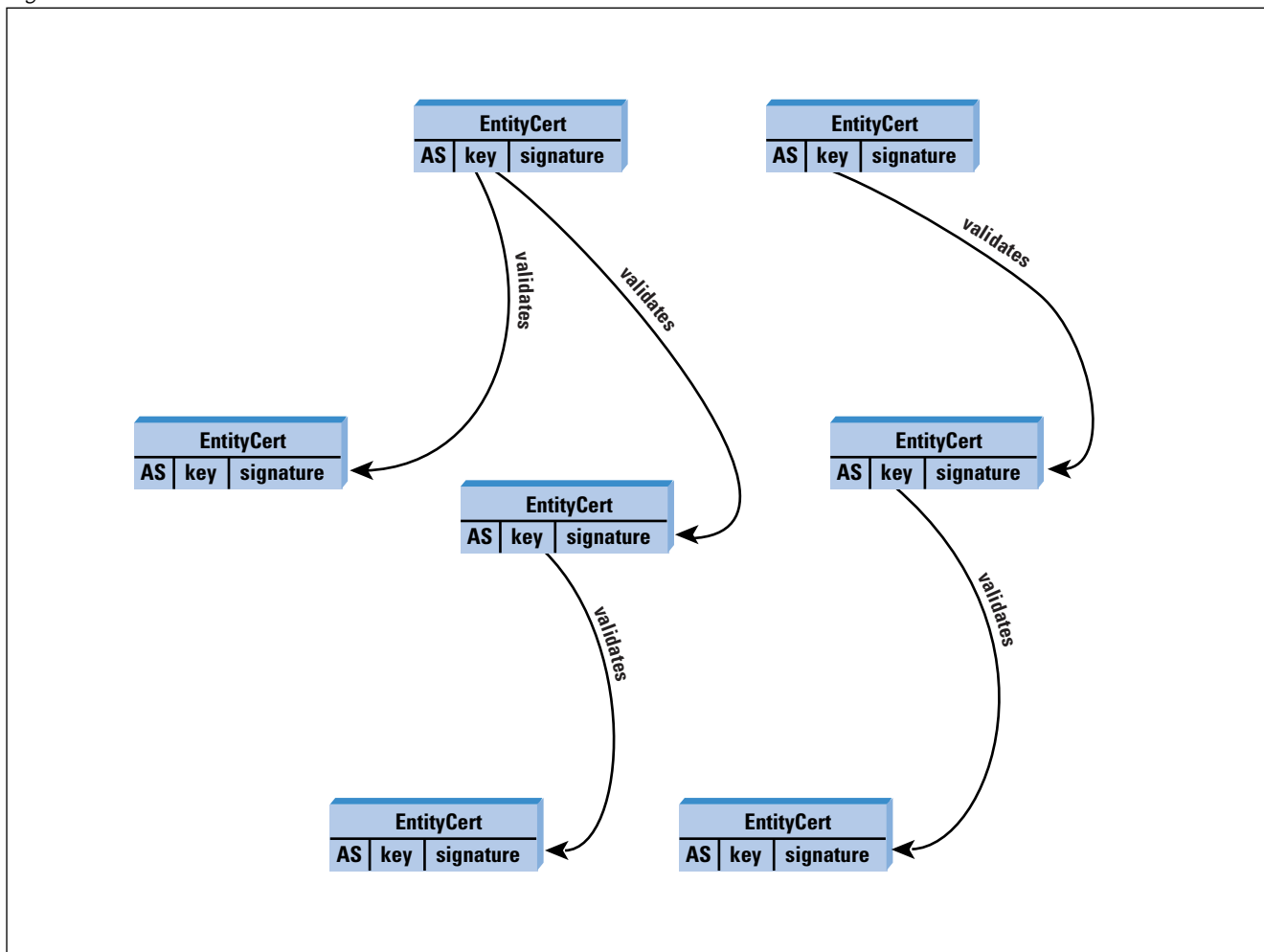
Begin at the Beginning: Who Are You?

The first step in securing anything is authentication; each participant must have some way of knowing who the other participants are, and what information they will be using to sign or encrypt their data. This is a classic problem in cryptography, called *key distribution*. There must be some way to receive keys used to sign or encrypt data, and then to validate that the keys received actually belong to the participant we believe they belong to.

This problem is addressed in soBGP using an *EntityCert*, which ties an AS number to a public key (or a set of public keys) corresponding to a private key the AS will be using to sign various other certificates. An EntityCert is defined in soBGP to be an X.509v3 certificate, similar to those used by *Transport Layer Security* (TLS) and *IP Security* (IPSec). The main problem we face when accepting an EntityCert is knowing whether or not the key carried within the certificate is actually the key of the advertising AS.

soBGP resolves this by requiring the EntityCert to be signed by a third party, validating that this AS actually belongs with this key. A small number of “root keys” distributed out of band could then be used to validate a set of advertised EntityCerts. These are used in turn to build up the database of known good ASm/key pairs in the system, allowing even more EntityCerts to be validated. Thus, EntityCerts can form a web of trust, built on the public keys of a small number of well-known entities, such as top-level backbone service providers, key authentication service providers (such as Verisign), and others.

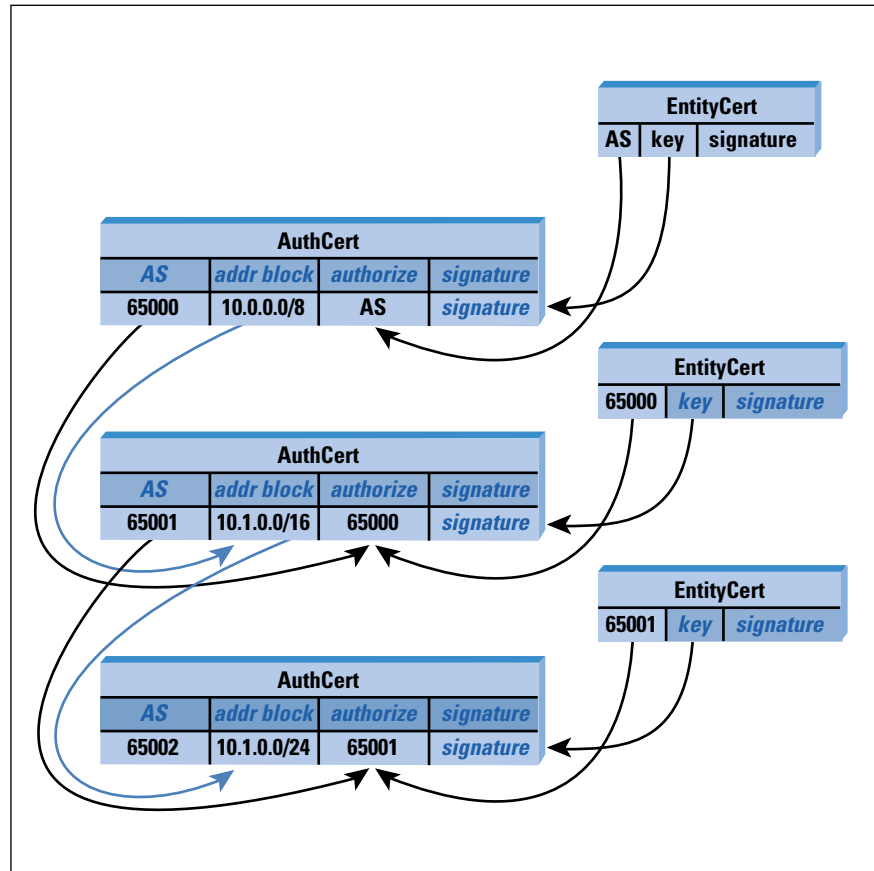
Figure 1: Web of Trust



The key each AS distributes in its EntityCert is actually the public half of a private/public key pair. An AS would keep its private key entirely private, holding it on one highly secure device in its network (which is not even required to be online), and generating signatures for other certificates as needed. Only an AS public key is ever exposed in this way, so no special protection mechanisms (for example, tamper-resistant hardware) are required at any border to prevent private keys from being compromised.

The First Goal: Are You Authorized?

Now that we have distributed a public key per AS, we can build a certificate that will provide authorization for an AS to advertise a specific block of addresses. This authorization is provided through an *Authorization Certificate*, or *AuthCert*. An AuthCert ties an AS to a block of addresses that the AS may advertise, as Figure 2 illustrates.

Figure 2: Authorization
Example

Starting at the top of the illustration, we find that some AS has authorized AS65000 to advertise prefixes within the block 10.0.0.0/8. The AuthCert is signed using the authorizing AS key. To delegate some part of this block of address space to another AS, AS65001, AS65000 builds an AuthCert tying 10.1.0.0/16 to AS65001. AS65001, in turn, suballocates a smaller part of this address space to AS65002, by building an AuthCert tying AS65002 to 10.1.1.0/24.

Any device receiving these three AuthCerts can check them by:

- Looking up the public key of the authorizer, and verifying the signature on the AuthCert
- Making certain the authorizer is permitted to advertise the address space it has suballocated this block of address space from

The device then builds a local table of address blocks and corresponding ASs authorized to advertise prefixes within those address blocks. Received updates can be checked against this database to verify authorization of the originating AS to advertise a prefix.

Blocks of address space are used here, rather than individual prefixes; an AuthCert can authorize an AS to advertise any number of prefixes within a block of addresses. This reduces the number of certificates within the system, thereby reducing overall cryptographic processing requirements. If a specific AS desires per-prefix authorization, it can build individual AuthCerts for each allocated prefix, rather than for blocks of address space.

Per-Prefix Policy

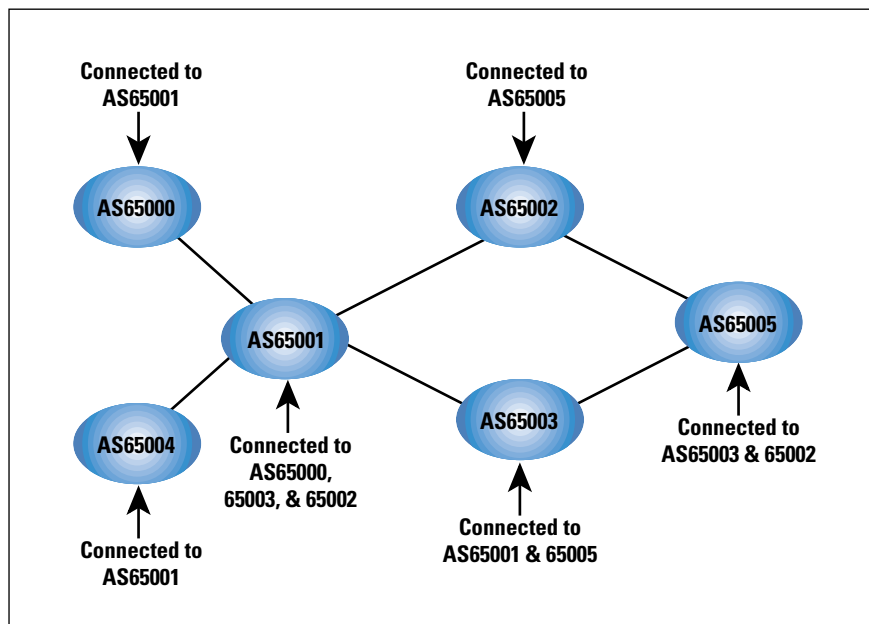
AuthCerts are not advertised as independent certificates within soBGP; instead, they are wrapped in a *PrefixPolicyCert*. *PrefixPolicyCerts* contain an AuthCert, a set of policies the originator would like to apply to prefixes advertised within this block of addresses, and a signature generated using the private key of the authorized AS. Policies that may be included in the *PrefixPolicyCert* include the longest prefix length allowed within the address block, and possibly other policies, such as a list of ASs that may not be or must be in the AS Path of routes to destinations within the address block.

In reality, the per-prefix policies available to the originator are limitless; the main problem is enforcing those policies when they are received by other ASs.

The Second Goal: Do You Really Have a Path?

Our second goal is to be able to verify that the advertiser of a given route actually has a path to the destination. This goal is met in soBGP by building a topology map of the paths of the entire internetwork. Each AS attached to the internetwork builds an *ASPolicyCert*, which contains, primarily, a list of its peers, and signed using the originator's private key. Using this list of transit peers, a map of the internetwork topology may be built, as Figure 3 illustrates.

Figure 3: Connectivity
Graph Example



If AS65005 receives an update from AS65002, claiming it can reach a destination in AS65000 through the path {65002, 65001, 65000}, it can:

- Check to make certain AS65002 claims to be connected to AS65001 in its *ASPolicyCert*, and that AS65001 claims to be connected to AS65002 in its *ASPolicyCert*
- Check to make certain AS65001 claims to be connected to AS65000 in its *ASPolicyCert*, and that AS65000 claims to be connected to AS65001 in its *ASPolicyCert*

If, for instance, AS65002 claims a path to a destination inside AS65000 through the path (65002, 65000), AS65002 would be able to discover that the path is invalid, because AS65000 does not claim to be connected to AS65002. This simple two-way connectivity check along a graph can be mixed with various policy statements—stating a specific peer is not a transit, not advertising certain peers, etc.—to provide a much wider range of policies than AS Path-based methods.

Transporting Certificates

One of the primary problems any security system such as soBGP is going to face is transporting security information through the internet-work. We would like to make certain we do not rely on the routing system to provide information about the security of the routing system. In other words, we would not like to rely on unsecured routing information in order to reach a server providing the information required to secure the path to the server itself.

soBGP resolves this by proposing to advertise certificates in much the same way as routing information is propagated today—through an interdomain protocol. Currently the soBGP drafts specify a new type of BGP message, the SECURITY message, which can be used to transport the required certificates, the EntityCert, the PrefixPolicyCert, and the ASPolicyCert, throughout an internetwork. Other methods of transporting data such as these certificates throughout an internetwork are currently being pursued by the IETF; if other methods are offered, soBGP could transport certificates across any such distribution mechanism.

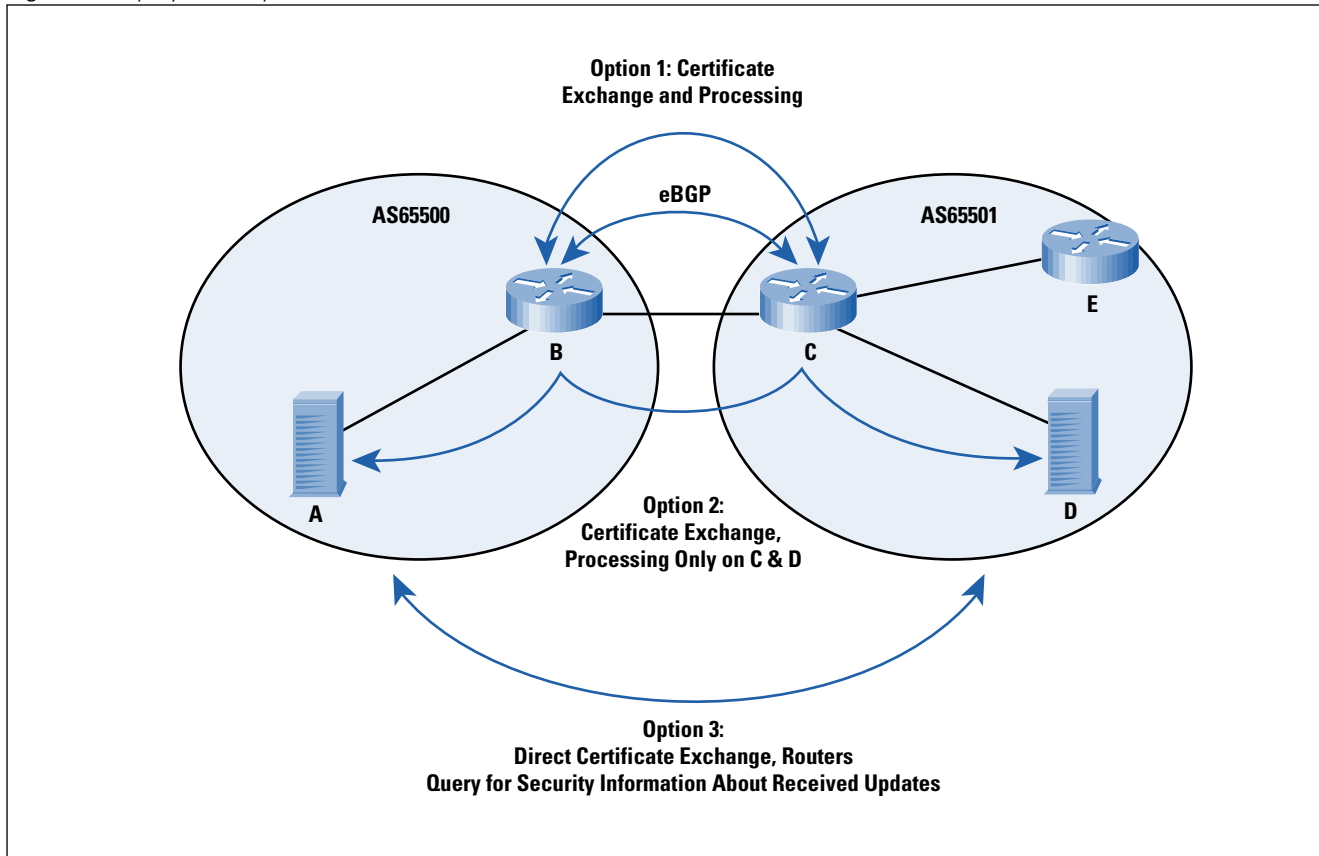
Deployment

Finally, we come to the hardest problem any routing security system is going to face: actually getting it operating in the field, with useful results, with a minimum of equipment changes, and a minimum number of participants. Here, soBGP provides a wide variety of options, primarily because it is not transport-dependent, nor dependent on a yet-to-be constructed centralized set of servers.

Although deployment options abound, here we discuss three, just to show the range of options available. Figure 4 illustrates these options.

The first option shown in this network is direct certificate exchange and processing between border routers. With this option, routers that are capable of the cryptographic processing required to validate received certificates exchange certificates with their peers in other ASs (just as they exchange routing information today), process those certificates, and build local databases from which they perform security checks on received updates.

Figure 4: Deployment Options



Although it may appear that processing, in this situation, would be extensive, it is actually possible to spread the processing required out among the border routers in a large AS. For instance, each certificate that router C receives and processes can be subsequently sent over an encrypted link to Router E. Router E could treat these certificates as though they had been validated locally, because they are received across an encrypted link from a trusted peer within the same administrative domain. Thus, only the edge router that has learned a certificate would actually process the certificate. This spreads the processing along all the edges in the AS.

A second option is for the edge routers, B and C, to exchange the certificates, but not process them. Instead, each edge router would relay the not-yet-validated certificates to internal servers A and D, respectively, thereby validating the certificates by performing the necessary cryptographic operations. As the border routers receive updates, they can query the server about the validity of each update, and take action based on the reply received.

Finally, it is possible for the servers to exchange certificates directly, over a multihop session. Servers A and D would then process the certificates, and the border routers, B and C, would query these servers to determine if received updates are valid or invalid.

Summary

Through this short survey of soBGP, we have shown it to be a flexible, moderately lightweight, yet strong system for validating the information carried through BGP in a large internetwork. It has low overhead processing requirements and very flexible deployment options, but no reliance on centralized servers. We are currently working to develop prototypes of soBGP on several platforms, to show how the technology will work on a wide range of devices.

For more information on soBGP, refer to:

<ftp://ftp-eng.cisco.com/sobgp/index.html>

You will find the most recent versions of the drafts, several slide shows, and other information about soBGP at this site.

RUSS WHITE works for Cisco Systems in the Routing Protocols Deployment and Architecture (DNA) team in Research Triangle Park, North Carolina. He has worked in the Cisco Technical Assistance Center (TAC) and Escalation Team in the past, has coauthored several books on routing protocols, including *Advanced IP Network Design*, *IS-IS for IP Networks*, and *Inside Cisco IOS® Software Architecture*. He is currently in the process of publishing a book on BGP deployment, and is the cochair of the Routing Protocols Security Working Group within the IETF. E-mail: riw@cisco.com

Trends in Viruses and Worms

by Thomas M. Chen, Southern Methodist University

The modern computer *virus* was conceived and demonstrated by Fred Cohen in 1983. Like biological viruses, computer viruses reproduce by attaching to a normal program or document and taking over control of the execution of that program to infect other programs. Early viruses could spread slowly mostly by floppies (such as the 1986 *Brain* virus), but the Internet has made it much easier for viruses to move among computers and spread rapidly. Networks have created a fertile environment for worms, which are related to viruses in their ability to self-replicate but are not attached to other programs. Worms are particularly worrisome as standalone automated programs designed to exploit the network to seek out vulnerable computers. The term *worm* was originated by John Shoch and Jon Hupp during their experiments on mobile software at Xerox PARC in 1979, inspired by the network-based *tapeworm* monster in John Brunner's novel, *The Shockwave Rider*^[1]. Shoch and Hupp thought of worms as multi-segmented programs distributed across networked computers.

The Internet increases the vulnerability of all interconnected machines by making it easier for malicious programs to travel between computers by themselves. Recent virus and worm outbreaks, such as the *Blaster* worm in August 2003 and the SQL *Sapphire/Slammer* worm in January 2003, have demonstrated that networked computers continue to be vulnerable to new attacks despite the widespread deployment of antivirus software and firewalls. Indeed, a review of the history of viruses and worms shows that they have continually grown in sophistication over the years. This article highlights a series of significant past innovations in virus and worm technology. The purpose is to show that viruses and worms continue to pose a major risk today and most likely into the future as their creators persist in seeking ways to exploit security weaknesses in networked systems.

Stealth

The earliest viruses attempted to hide evidence of their presence, a trend that continues to today. The 1986 DOS-based *Brain* virus hid itself in memory by simulating all of the DOS system calls that normally detect viruses, causing them to return information that gave the appearance that the virus was not there.

The 2001 *Lion* worm installed a rootkit called *t0rn*, which is designed to make the actions of the worm harder to detect through numerous system modifications to deceive *syslogd* from properly capturing system events (*syslogd* is often used to detect worm activity)^[2]. More recently, viruses and worms have attempted to hide by actively attacking antivirus software on the infected computer (refer to the section "Armoring").

Social Engineering

The 1987 *Christma Exec* virus was an early example of social engineering, spreading by e-mail among IBM mainframes. An arriving message tricks the user into executing the virus by promising to draw a Christmas tree graphic. The virus does produce a Christmas card graphic on the computer display (drawn using a scripting language called *Rexx*) but sends a copy of itself in the user's name to that user's list of outgoing mail recipients. The recipients believe the e-mail is from the user, so they are more likely to open the e-mail.

Social engineering continues to be common practice in today's viruses and worms, particularly those spread by e-mail. In January 1999, the *Happy99/Ska* worm/Trojan horse hybrid spread by e-mail with an attachment called **Happy99.exe**^[3]. When the attachment was executed, it displayed fireworks on the screen to commemorate New Year's Day, but secretly modified the **WSOCK32.DLL** file (the main Windows file for Internet communications) with a Trojan horse program that allowed the worm to insert itself into the Internet communications process. Every e-mail sent by the user generated a second copy without any text but carried the worm to the same recipients.

The 1999 *PrettyPark* worm propagated as an e-mail attachment called **Pretty Park.exe**. The attachment is not explained, but it bears the icon of a character from the television show, *South Park*. If executed, it installs itself into the Windows System folder and modifies the Registry to ensure that it runs whenever any **.EXE** program is executed. In addition, the worm e-mails itself to addresses found in the Windows Address Book. It also mails some private system data and passwords to certain *Internet Relay Chat* (IRC) servers. Reportedly, the worm also installs a backdoor to allow a remote machine to create and remove directories, and send, receive, and execute files.

In February 2001, the *Anna Kournikova* virus demonstrated social engineering again, pretending to carry a JPG picture of the tennis player. If executed, the virus e-mails a copy of itself to all addresses in the Outlook address book.

In March 2002, the *Gibe* worm spread as an attachment in an e-mail disguised as a Microsoft security bulletin and patch. The text claimed that the attachment was a Microsoft security patch for Outlook and Internet Explorer. If the attachment is executed, it displays dialog boxes that appear to be patching the system, but a backdoor is secretly installed on the system.

Macro Viruses

The *Concept* virus was the first macro virus, written for Word for Windows 95. The vast majority of macro viruses are targeted to Microsoft Office documents that save macro code within the body of documents. Macro viruses have the advantages of being easy to write and independent of computing platform. However, macro viruses are no longer widespread after people have become more cautious about using the Office macro feature.

Mass E-Mailers

In March 1999, the *Melissa* macro virus spread quickly to 100,000 hosts around the world in three days, setting a new record and shutting down e-mail for many organizations using Microsoft Exchange Server^[4]. It began as a newsgroup posting promising account names and passwords for erotic Web sites. However, the downloaded Word document actually contained a macro that used the functions of Microsoft Word and the Microsoft Outlook e-mail program to propagate. Up to that time, it was widely believed that a computer could not become infected with a virus just by opening e-mail. When the macro is executed in Word, it first checks whether the installed version of Word is infectable. If it is, it reduces the security setting on Word to prevent it from displaying any warnings about macro content. Next, the virus looks for a certain Registry key containing the word “Kwyjibo” (apparently from an episode of the television show, *The Simpsons*). In the absence of this key, the virus launches Outlook and sends itself to 50 recipients found in the address book. Additionally, it infects the Word **NORMAL.DOT** template using the Microsoft *Visual Basic for Applications* (VBA) macro auto-execute feature. Any Word document saved from the template would carry the virus.

In June 1999, the *ExploreZip* worm appeared to be a WinZip file attached to e-mail but was not really a zipped file^[5]. If executed, it appears to display an error message, but the worm secretly copies itself into the Windows Systems directory or loads itself into the Registry. It sends itself via e-mail using Outlook or Exchange to recipients found in unread messages in the inbox. It monitors all incoming messages and replies to the sender with a copy of itself.

In May 2000, the fast-spreading *Love Letter* worm demonstrated a social engineering attack^[6]. It propagated as an e-mail message with the subject “I love you” and text that encourages the recipient to read the attachment. The attachment is a Visual Basic script that could be executed with Windows Script Host (present if the computer has Windows 98, Windows 2000, Internet Explorer 5, or Outlook 5). Upon execution, the worm installs copies of itself into the Windows System directory and modifies the Registry to ensure that the files are run when the computer starts up. The worm also infects various types of files (for example, **.VBS**, **.JPG**, **.MP3**, etc.) on local drives and networked shared directories. If Outlook is installed, the worm e-mails copies of itself to addresses found in the address book. In addition, the worm makes a connection to IRC and sends a copy of itself to anyone who joins the IRC channel. The worm has a password-stealing feature that changes the startup URL in Internet Explorer to a Website in Asia. The Website downloads a Trojan horse designed to collect various passwords from the computer.

In 2002, 90 percent of the known viruses were mass e-mailers. Two of the most prevalent ones, *Bugbear* and *Klez*, began a trend of carrying their own *Simple Mail Transfer Protocol* (SMTP) engines. Although e-mail continues to be the most common infection vector, recent worms have been exploring new vectors (see the section “New Infection Vectors”).

In addition, mail servers are becoming more powerful in their capabilities to detect and filter malicious code. For these reasons, mass e-mailing may decline as an infection vector for future viruses.

Polymorphism

Polymorphism is based on the simpler idea of encryption, which makes a virus harder to detect by antivirus software scanning for a unique virus signature (byte pattern). Encryption attempts to hide a recognizable signature by scrambling the virus body. To be executable, the encrypted virus is prepended with a decryption routine and encryption key. However, encryption is not effective because the decryption routine remains the same from generation to generation, although the key can change, scrambling the virus body differently. Antivirus scanners can detect a sequence of bytes identifying a specific decryption scheme.

Polymorphic viruses permute continuously to avoid detection by antivirus scanning^[7]. The earliest polymorphic virus might have been a virus found in Europe in 1989. This virus replicated by inserting a pseudorandom number of extra bytes into the decryption algorithm, preventing any common sequence of more than a few bytes between two successive infections. Polymorphism became practical when a well-known hacker, *Dark Avenger*, developed a user-friendly *Mutation Engine* program to provide any virus with variable encryption. With a static signature so small, the risk of false positives by antivirus scanners became very high. Other hackers soon followed with their own versions of so-called mutation engines. The 1995 *Pathogen* and *Queeg* viruses were polymorphic DOS file-infecting viruses produced by Black Baron's *Simulated Metamorphic Encryption enGine* (SMEG)^[7].

Blended Attacks

The famous 1988 *Morris* worm was the first to use a combination of attacks (or blended attacks) to spread quickly to 6000 UNIX computers in a few hours (10 percent of the Internet at that time)^[8].

- It captured the password file and ran a password-guessing program on it using a dictionary of common words.
- It exploited the debug option in the UNIX *sendmail* program, allowing it to transfer a copy of itself.
- It carried out a buffer overflow attack through a vulnerability in the UNIX *fingerd* program.

In May 2001, the *Sadmind/IIS* worm spread by targeting two separate vulnerabilities on two different operating systems. It first exploited a buffer overflow vulnerability in Sun Solaris systems and installed software to carry out an attack to compromise Microsoft *Internet Information Services* (IIS) Web servers.

The July 2001 *Sircam* worm uses two ways to propagate. First, it e-mails itself as an attachment using its own SMTP engine, and if the attachment is executed, e-mails a copy of itself to addresses found in the Windows address book. Second, it spreads by infection of unprotected network shares.

In September 2001, *Nimda* raised new alarms by using five different ways to spread to 450,000 hosts within the first 12 hours^[9]. *Nimda* seemed to signal a new level of worm sophistication.

- It found e-mail addresses from the computer Web cache and default *Messaging Application Programming Interface* (MAPI) mailbox. It sent itself by e-mail with random subjects and an attachment named **readme.exe**. If the target system supported the automatic execution of embedded MIME types, the attached worm would be automatically executed and infect the target.
- It infected Microsoft IIS Web servers, selected at random, through a buffer overflow attack called a *unicode* Web traversal exploit.
- It copied itself across open network shares. On an infected server, the worm wrote *Multipurpose Internet Mail Extensions* (MIME)-encoded copies of itself to every directory, including network shares.
- It added JavaScript to Web pages to infect any Web browsers going to that Website.
- It looked for backdoors left by previous *Code Red II* and *Sadmind* worms.

Armoring

In November 2002, the *Winevar* worm was an example of an “armored” worm that contained special code designed to disable antivirus software using a list of keywords to scan memory to recognize and stop antivirus processes and scan hard drives to delete associated files^[10].

Klez and *Bugbear* are recent examples of worms that attack antivirus software by stopping active processes and deleting registry keys and database files used by popular antivirus programs. The 2003 *Fizzer* and *Lirva* worms also attempt to disable antivirus software.

Dynamic Software Updates

In October 2000, the *Hybris* worm propagated as an e-mail attachment^[11]. It connected to the **alt.comp.virus** newsgroup to receive encrypted plug-ins (code updates). The method is sophisticated and potentially very dangerous, because the worm payload (destructive capability) can be modified dynamically.

The 2003 *Lirva* worm attempted to connect to a Website on **web.host.kz** to download BackOrifice, a notorious remote-access software package that gives complete control to a remote attacker. It also attempted to download another unknown file that was not found on the Website.

This technique was given an interesting twist by the *Welchia* or *Nachi* worm, which began spreading on August 18, 2003, soon after the *Blaster* worm. Apparently, its creator intended *Welchia* as a “good” worm to remove *Blaster*. It attempted to download and install a fix for *Blaster* from a Microsoft Website.

New Infection Vectors

The Linux *Slapper* worm, appearing in September 2002, was among the first to exploit *peer-to-peer* (P2P) technology^[12]. It spread to Linux computers by exploiting the long *Secure Sockets Layer 2* (SSL2) key argument buffer overflow in the *libssl* library, used by the *mod_ssl* module of the Apache 1.3 Web server. When the worm infects a new machine, it binds to *User Datagram Protocol* (UDP) port 2002 and becomes part of a P2P network. The parent of the worm on the attacking machine sends to its offspring the list of all hosts on the P2P network and broadcasts the address of the new worm on the network. Then periodic updates to the host list are exchanged between machines on the network. The new worm also scans the network for other vulnerable machines, sweeping randomly chosen class B networks.

In March 2003, the *AimVen* worm spread by the *America OnLine Instant Messenger* (AIM) by modifying the AIM program. Whenever an **.EXE** file is sent through AIM, the worm overwrites the file with a copy of itself.

The *Fizzer* worm discovered in May 2003 is a mass e-mailer that includes its own SMTP engine like *Klez* and *Bugbear*. It also tries to spread via *KaZaa*, a popular P2P file-sharing application, and shared directories.

The 2003 *Lirva* worm, named after the singer, Avril Lavigne, is a mass e-mailer taking advantage of the same MIME header exploit as *Badtrans* and *Klez*, but also tries to spread by IRC, “I seek You” (ICQ), *KaZaa*, and open network shares^[13].

Data-Stealing Payloads

Most fast-spreading worms in the past have not carried destructive payloads. Instead, they have tended to appear to be proof-of-concepts to demonstrate a particular security weakness. Some worms, though, such as *Code Red*, have installed *Denial-of-Service* (DoS) agents or backdoors on infected machines. Recently worms have begun to carry keyloggers and password-stealing Trojans in their payloads.

The 2003 *Fizzer* worm includes a keystroke logging Trojan horse that stores the data in an encrypted file. It establishes its own accounts on IRC and AIM to wait for instructions from the virus writer, who could conceivably fetch the keystrokes data.

The 2003 *Lirva* worm e-mails cached Windows dialup networking passwords to the virus writer, and e-mail random **.TXT** and **.DOC** files to various addresses.

Bugbear installs a keystroke logging tool into the Windows System folder that e-mails the keystrokes data to preprogrammed addresses^[14]. It listens on port 36794 for commands from a remote hacker.

Fast and Furious Worms

A particularly worrisome new trend is extremely fast worms targeted to specific (usually Windows-related) vulnerabilities that might saturate their target population within a few hours or even less than an hour. These worms tend to be simpler and targeted to single rather than multiple vulnerabilities, in order to be highly efficient in their probing for other vulnerable machines.

The first example might be the *Code Red* worm, which actually appeared in three different versions^[15]. The first version of *Code Red I* appeared on July 12, 2001, targeted to a buffer overflow vulnerability in Microsoft IIS Web servers. However, a programming error in its pseudorandom address generator caused each worm copy to probe the same set of IP addresses and prevented the worm from spreading quickly. A week later on July 19, a second version of *Code Red I* with the programming error apparently fixed was able to infect more than 359,000 servers within 14 hours. At its peak, the worm was infecting 2000 hosts every minute. A more complex and dangerous *Code Red II* targeted to the same IIS vulnerability appeared on August 4.

More recently, the *Structured Query Language (SQL) Sapphire/Slammer* worm appeared on January 25, 2003, targeted to Microsoft SQL Server machines not running *Service Pack 3 (SP3)*, such as SQL Server 2000 and *Microsoft Desktop Engine (MSDE) 2000*^[16]. It reportedly infected 90 percent of vulnerable hosts within 10 minutes (about 120,000 servers)^[17]. The spreading rate was surprisingly fast and resulted in DoS effects (network outages and high packet loss) due to traffic overloading servers and routers. In the first minute, the infection doubled every 8.5 seconds, and hit a peak scanning rate of 55,000,000 scans per second after only 3 minutes. In comparison, *Code Red* infection doubled in 37 minutes (slower but infected more machines). *Slammer* was able to spread so quickly because it appeared to be designed simply for efficient replication. The worm carried no payload and consisted of a single 404-byte UDP packet (including 376 bytes for the worm) that could be sent without having to wait for responses from targeted machines. In contrast, *Code Red* was about 4000 bytes and *Nimda* was 60,000 bytes, and their scanning depended on the time to establish TCP connections to targeted machines. The *Slammer* worm was much more efficient, simply generating copies of itself at the full rate of the infected machine.

Latest Developments

The week of August 12–19, 2003, has been called the worst week for worms in history, seeing *MS Blaster*, *Welchia* (or *Nachi*), and *Sobig.F* in quick succession. *MS Blaster* or *LovSan* was another fast worm, which appeared on August 12, 2003, targeted to a *Windows Distributed Component Object Model (DCOM) Remote Procedure Call (RPC)* vulnerability announced on July 16, 2003^[18]. The worm probes for a DCOM interface with RPC listening on TCP port 135 on Windows XP and Windows 2000 PCs. Through a buffer overflow attack, the worm causes the target machine to start a remote shell on port 4444 and send a notification to the attacking machine on UDP port 69.

A *Trivial File Transfer Protocol* (TFTP) “get” command is then sent to port 4444, causing the target machine to fetch a copy of the worm as the file **MSBLAST.EXE**. In addition to a message against Microsoft, the worm payload carries a DoS agent (using TCP SYN flood) targeted to the Microsoft Website **windowsupdate.com** on August 16, 2003. Although *Blaster* has reportedly infected about 400,000 systems, experts reported that the worm did not achieve near its potential spreading rate because of novice programming.

Six days later on August 18, 2003, the apparently well-intended *Welchia* or *Nachi* worm spread by exploiting the same RPC DCOM vulnerability as *Blaster*. It attempted to remove *Blaster* from infected computers and download a security patch from a Microsoft Website to repair the RPC DCOM vulnerability. Unfortunately, its scanning resulted in a DoS effect on some networks, such as Air Canada’s check-in system and the U.S. Navy and Marine Corps computers.

The very fast *Sobig.F* worm appeared on the next day, August 19, 2003, only seven days after *Blaster*^[19]. The original *Sobig.A* version was discovered in January 2003, and apparently underwent a series of revisions until the most successful *Sobig.F* variant. Similar to earlier variants, *Sobig.F* spreads among Windows machines by e-mail with various subject lines and attachment names, using its own SMTP engine. The worm size is about 73 kilobytes with a few bytes of garbage attached to the end to evade antivirus scanners. It works well because it grabs e-mail addresses from a variety of different types of files on the infected computer and secretly e-mails itself to all of them, pretending to be sent from one of the addresses. At its peak, *Sobig.F* accounted for 1 in every 17 messages, and reportedly produced over 1 million copies of itself within the first 24 hours. Interestingly, the worm was programmed to stop spreading on September 10, 2003, suggesting that the worm was intended as a proof-of-concept. This is supported by the absence of a destructive payload, although the worm is programmed with the capability to download and execute arbitrary files to infected computers. The downloading is triggered on specific times and weekdays, which are obtained via one of several *Network Time Protocol* (NTP) servers. The worm sends a UDP probe to port 8998 on one of several preprogrammed servers, which responds with a URL for the worm to download. The worm also starts to listen on UDP ports 995–999 for incoming messages, presumably instructions from the creator.

Conclusions

Why does the Internet remain vulnerable to large-scale worm outbreaks? Since at least 1983, the Internet community has understood the risks and mechanics of viruses. The 1988 Morris worm taught the community to be watchful for potentially dangerous worms. Over the years, a variety of antivirus software, firewalls, intrusion detection systems, and other security equipment have been installed. Moreover, the *Computer Emergency Response Team* (CERT) at CMU was established as the first computer security incident response team, which later joined an expansive global coalition of security incident response teams called the *Forum of Incident Response and Security Teams* (FIRST)^[20].

Despite our knowledge and infrastructure defenses, many viruses and worms have broken out regularly in the Internet over the years. By some reports, 5 to 15 new viruses and worms are released every day, although a fraction of that number are not released in the wild and most do not spread well. Still, fast-spreading viruses and worms continue to appear with regularity. Outbreaks have become so commonplace that most organizations have come to view them as a routine cost of operation.

The problem is sometimes portrayed as a perpetual struggle between virus writers who keep innovating (as described here) and the antivirus industry, which tries to keep up. However, the problem is actually larger, involving the entire computer industry. Viruses and worms are successful because computers have security vulnerabilities that can be exploited. Clearly, the Internet itself is simply serving its purpose of interconnecting computer systems. The security vulnerabilities exist in the host end systems. Security vulnerabilities continue to exist for many reasons. First, software is often written in an unsecure manner, for example, vulnerable to buffer overflow attacks that are commonly used by worms. Buffer overflow attacks have been widely known since 1995, but this type of vulnerability continues to be found very often (on every operating system.) Second, when vulnerabilities are announced with corresponding software patches, many people are slow to apply patches to their computer for various practical reasons. Weakly protected computers can be compromised, putting the entire community at risk, including secured computers that can still be impacted by the traffic effects of a worm outbreak.

However, there is reason to be hopeful for a solution. Fortunately, worms typically have a weakness of exploiting vulnerabilities that have been known for some time. Worm writers do not invent new exploits for the simple reason that they want to ensure that their worm will spread after it is released. For example, the *Code Red I* worm took advantage of a buffer overflow vulnerability in Microsoft IIS servers that had been known for a month. The *Nimda* worm exploited a unicode Web traversal vulnerability in Microsoft IIS servers that was published a year earlier. The SQL *Slammer/Sapphire* worm exploited a buffer overflow vulnerability in Microsoft SQL servers that had been known for six months. The recent *Blaster* worm exploited a Windows DCOM RPC vulnerability announced two months earlier. Watching for probing activity attempting to exploit known vulnerabilities could help detect and block worm outbreaks at an early stage. Ideas for automatic detection and quarantine of new epidemics is attracting research^[21].

Aside from technological considerations, an important issue is accountability. The most obvious parties to hold liable are the virus creators, but it has been observed many times that few virus writers have been prosecuted, and sentences have tended to be light. The author of the 1988 Internet worm, Robert Morris, was sentenced to three years of probation, 400 hours of community service, and a \$10,000 fine.

Chen Ing-hau was arrested in Taiwan for the 1998 *Chernobyl* virus, but he was released when no official complaint was filed. Onel de Guzman was arrested for writing the 2000 *LoveLetter* virus, which resulted in \$7 billion of damages, but he was released because of the lack of relevant laws in the Philippines. Jan De Wit was sentenced for the 2001 *Anna Kournikova* virus to 150 hours of community service. David L. Smith, creator of the 1999 *Melissa* that caused at least \$80 million of damages, was sentenced to 20 months of custodial service and a \$7500 fine.

It is notoriously difficult to trace a virus or worm to its creator from analysis of the code, unless inadvertent clues are left in the code. In addition, cases are difficult to prosecute, and malicious intention (as opposed to just recklessness) is difficult to prove. Moreover, long prison sentences have been perceived as overly harsh for arrested virus creators, who have tended to be teenagers and university students. In addition, in the absence of a serious legal deterrent, the general perception persists that virus creators can easily avoid the legal consequences of their actions. Perhaps to address this problem, authorities have been diligently investigating the creators of *Blaster* and *Sobig*. So far, a teenager, Jeffrey Lee Parson, has been arrested for writing the *Blaster.B* variant, a slight modification of the original *Blaster*. Soon afterward, Dan Dumitru Ciobanu was arrested in Romania for writing the *Blaster.F* variant.

Some have argued wishfully that software vendors should be held financially liable for damages resulting from the security vulnerabilities in their products. The assumption is that accountability would increase motivation to write and sell more secure software, a solution that would result in a less inviting environment for viruses and worms. So far, software vendors have managed to acknowledge their role but avoid accountability.

References

- [1] J. Shoch and J. Hupp, "The 'worm' programs—early experience with a distributed computation," *Communications of ACM*, Volume 25, pp. 172–180, March 1982.
- [2] A. Kasarda, "The Lion worm: king of the jungle?" SANS reading room, <http://www.sans.org/rr>
- [3] CERT incident note CA-1999-02, "Happy99.exe trojan horse," http://www.cert.org/incident_notes/IN-99-02.html
- [4] CERT advisory CA-1999-04, "Melissa macro virus," <http://www.cert.org/advisories/CA-1999-04.html>
- [5] CERT advisory CA-1999-06, "ExploreZip trojan horse program," <http://www.cert.org/advisories/CA-1999-06.html>
- [6] CERT advisory CA-2000-04, "Love letter worm," <http://www.cert.org/advisories/CA-2000-04.html>
- [7] D. Harley, R. Slade, and R. Gattiker, *Viruses Revealed*, Osborne/McGraw-Hill, 2001.

- [8] E. Spafford, "The Internet worm program: an analysis," *ACM Computer Communications Review*, Volume 19, pp. 17–57, January 1989.
- [9] CERT advisory CA-2001-26, "Nimda worm,"
<http://www.cert.org/advisories/CA-2001-26.html>
- [10] Virus Bulletin, "W32/WineVar,"
<http://www.virusbtn.com/resources/viruses/winevar.xml>
- [11] CERT incident note IN-2001-02, "Open mail relays used to deliver Hybris worm,"
http://www.cert.org/incident_notes/IN-2001-02.html
- [12] F-Secure, "F-Secure virus descriptions: Slapper,"
<http://www.f-secure.com/v-descs/slapper.shtml>
- [13] Symantec Security Response, "W32.Lirva.C@mm,"
<http://securityresponse.symantec.com/avcenter/venc/data/w32.lirva.c@mm.html>
- [14] Sophos, "W32/Bugbear-A,"
<http://www.sophos.com/virusinfo/analyses/w32bugbeara.html>
- [15] H. Berghel, "The Code Red worm," *Communications of ACM*, Volume 44, pp. 15–19, December 2001.
- [16] CERT advisory CA-2003-04, "MS-SQL server worm,"
<http://www.cert.org/advisories/CA-2003-04.html>
- [17] D. Moore, et al., "The spread of the Sapphire/Slammer worm,"
<http://www.caida.org/outreach/papers/2003/sapphire/sapphire.html>
- [18] CERT advisory CA-2003-20, "W32/Blaster worm," Aug. 11, 2003,
<http://www.cert.org/advisories/CA-2003-20.html>
- [19] Symantec Security Response, "W32.Sobig.F@mm,"
<http://securityresponse.symantec.com/avcenter/venc/data/w32.sobig.f@mm.html>
- [20] Forum of Incident Response and Security Teams (FIRST),
<http://www.first.org>
- [21] D. Moore, C. Shannon, G. Voelker, and S. Savage, "Internet quarantine: requirements for containing self-propagating code," IEEE Infocom 2003, San Francisco, April 2003.

THOMAS M. CHEN holds BS and MS degrees in electrical engineering from MIT, and a PhD in electrical engineering from the University of California, Berkeley. From 1989 to 1997, he worked on ATM networking research at GTE Laboratories (now Verizon). He is currently an Associate Professor in the Department of Electrical Engineering at SMU in Dallas, Texas. He is the associate editor-in-chief of *IEEE Communications Magazine*, a senior editor of *IEEE Network*, an associate editor of *ACM Transactions on Internet Technology*, and founding editor of *IEEE Communications Surveys*. He is the coauthor of *ATM Switching Systems* (Artech House, 1995). E-mail: tchen@engr.smu.edu

IPv6 Behind the Wall

by Jim Bound

IPv6 has technology advantages over IPv4, and most of them will not be seen by the end user any more than users see features added to other extensions to the Internet Protocol suite, sensors on their automobiles, or from any core technology evolution. This article focuses on three of those IPv6 technology advantages “Behind the Wall.”

An essential catalyst for the Next-Generation Internet is the *Internet Protocol Version 6* (IPv6), which will provide an evolution to a more pervasive use of the Internet and networking in general. The current Internet, using IPv4, is insufficient to support the business and operational preconditions for peer-to-peer applications and security, billions of mobile devices, sensor networks, and the requisite distributed computing infrastructure to support a mobile society. The “band aids” applied to permit the current Internet to keep it operating has created additional operational costs and reduced operational capabilities for users and networks.

This article is an IPv6 Forum (www.ipv6forum.com) statement of the technology advantages of IPv6.

IPv6 Supports End-to-End Applications and Security

There are several schools of thought and opinions on the issue of address space and all project different results, depending on one’s mathematical view and philosophy regarding use models. There is also the effect of disruptive technology, which can make moot any projections of IPv4 address space. In that sense, rationing is justified and intelligent. The IPv6 Forum believes we already are experiencing the initial quake of disruptive technology, and that there is a need for users and markets to evolve further with a basic tenet that end-to-end applications and security are a priori for that evolution to begin. The IPv6 Forum believes that *Network Address Translation* (NAT) is about control, but that control comes at a cost of the freedom to use peer-to-peer computing over client to server-only computing.

Two users on the Internet today generally cannot each initiate peer-to-peer communications with each other because their location and identity are not available to each other from two disparate networks. In addition, security between them must trust a third party, and absolute private communications is impossible. The reason is that the Internet has evolved so that users are generally behind NATs that preclude peer-to-peer communications, or the exchange of private security credentials. Some will say this affords users security on the Internet. Although NAT does provide a denial-of-service perimeter, it also provides a denial of service to a direct trust relationship between peers. IPv6 is the only way to have peer-to-peer security for the Next-Generation Internet at a reasonable cost and a true privacy trust model on the Internet.

In the field of network computer science when engineers and architects implement translation functions in a solution, a cost is incurred that would not exist without translation. This is due to the need to keep *state* before, during, and after the translation. In software engineering terminology, these *state machines* add time and space costs to the entire operation. In addition, a NAT box is a single point of failure, because it is the only point on the network where a user can exit or enter when translation exists. Translation also does not permit the use of all functions possible without translation because too many participants need to know the mappings, and each function requires a separate state to be maintained, and the time + space costs increase exponentially. The time + space costs of NAT to keep the Internet operational have been passed on to every part of the current Internet business, consumer, and government market sectors, and cannot even support the original functions of the Internet before NAT. The current Internet has no hope of supporting the functions of the Next-Generation Internet required or of offering a solution to the great digital divide that exists currently and is increasing daily.

The good news is that IPv6 is evolving, early adopter deployment has begun, and vendors have delivered initial IPv6 products to the market. IPv6 will not require NAT, and the infrastructure supports a stateless architecture for the Internet, using statefull properties only where they can be used without a translation attribute or policy. IPv6 inherently supports mobile communications, billions of devices, and sensor networks that will be pervasive at a reasonable cost and provide the option to eliminate the digital divide within the current Internet.

IPv6 Supports a Stateless Node Discovery Architecture

A Next-Generation Internet base technology advantage for mobile user devices, ad hoc networks, mobile network providers, and generally for all users is the *Stateless Node Discovery Architecture* inherent within IPv6.

IPv6 nodes can discover each other and form IPv6 addresses to communicate on a network using what is called *Neighbor Discovery* and *Stateless Autoconfiguration*. IPv6 supports an extensible stateless node discovery paradigm, which provides the following features:

- Discover presence of nodes on the network
- Discover Datalink Layer nodes on the network
- Discover routers on the network
- Discover link configuration parameters on the network

These features permit an IPv6 node to obtain and maintain information about the accessibility of another node on the network for communications. Node Discovery is the predecessor to the node obtaining an address from IPv6 autoconfiguration. This core IPv6 technology framework also permits nodes to communicate on networks where there are no routers within an ad hoc network.

A host, when booted on an IPv6 link, first creates a *link-local* address by taking the architecturally defined prefix in Neighbor Discovery **FE80**, and appending an *End User Identifier* (EUI), determined by the host, to that prefix. This link-local address is then verified on the link that it is not duplicated with other link-local addresses on that host's link. This host communication is performed using link IPv6 multicast packets, to avoid duplicate link-local addresses, which are not permitted on an IPv6 Link.

The host then uses the link-local address to send on the IPv6 link *Neighbor Solicitations*, and all other hosts on that link see those multicast solicitations, and then return *Neighbor Advertisements* to the host. After this communications process, all nodes on the IPv6 link can now communicate, and communication was accomplished without the use of servers or routers in a stateless manner.

The host also listens for *Router Advertisements* on the IPv6 link (or sends *Router Solicitations*), which provide address prefixes, link configuration parameters, and information as to whether or not to use a stateless or stateful method for address assignment, and additional network configuration parameters using the *Dynamic Host Configuration Protocol for IPv6* (DHCPv6)^[1].

If the host is instructed to use the stateless method for address configuration, then it can use the router prefixes announced to form IPv6 addresses from those prefixes by appending the EUI determined from the link-local address to that prefix to create an IPv6 Address. IPv6 supports multiple address types within the address architecture^[2,3]. If the host is instructed to use the stateful method for address configuration, then DHCPv6 can be used to configure additional hosts' addresses.

Users will not see these IPv6 stateless advantages for network communications, but they will exist behind the wall of the user to provide a new and improved set of mechanisms for Node Discovery and Address Autoconfiguration far more robust and efficient than using the current IP Version 4 (IPv4) protocol. The IPv6 Stateless Architecture for Node Discovery permits a new model for node communications on links.

The Mobile IPv6 Technology Value Proposition

Mobile IPv6 offers many improvements over Mobile IPv4. Mobile IP as a technology permits users to remain connected across wireline (for example, Ethernet, xDSL) and wireless (for example, 802.11, cellular, satellite) networks, while roaming between networks. This permits users to stay connected while on the way to the airport from home, rather than shutting down their personal digital assistant (PDA)/laptop at home, and reconnecting at the WiFi location at the airport.

Figure 1: Route Optimization with Built-In Security

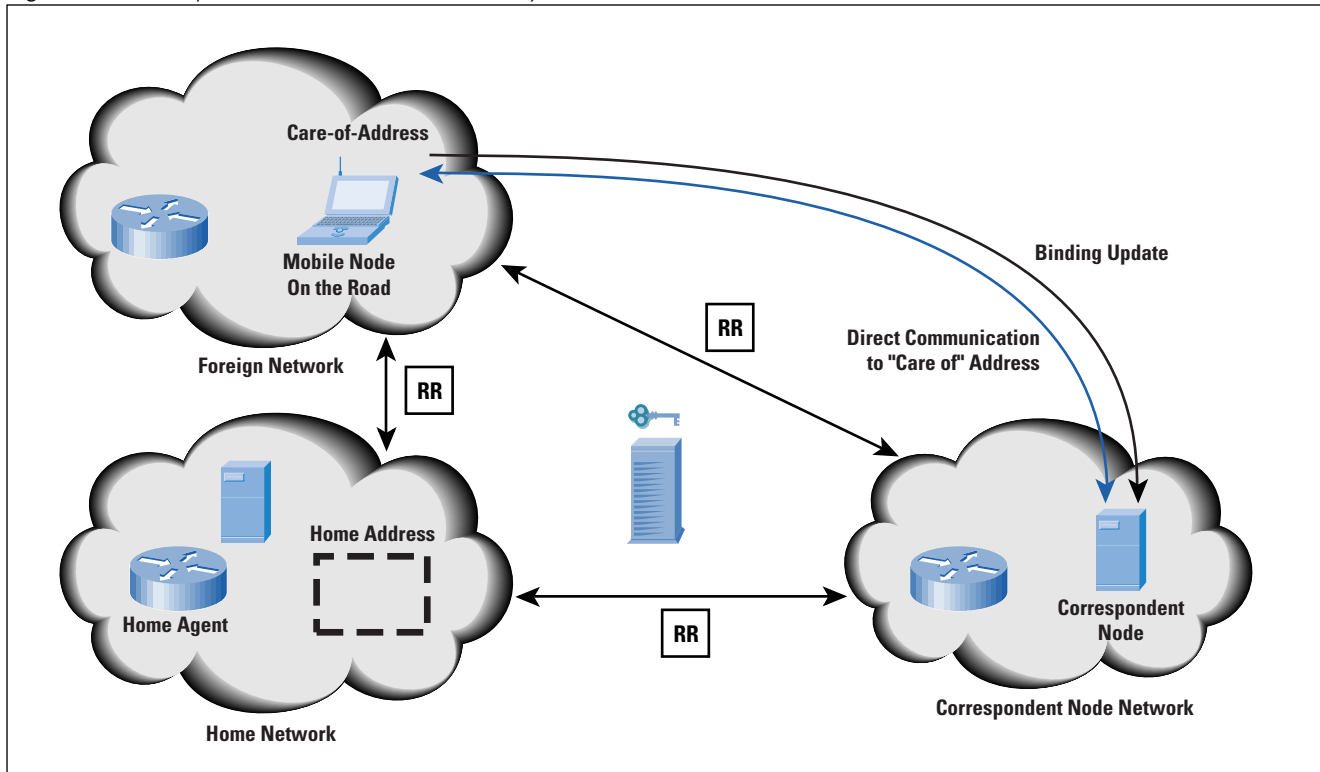


Figure 1 depicts the multiple phases of a mobile IPv6 connection. On the home network, a mobile node receives its home address as any IPv6 node. The mobile node registers that address with the *Home Agent*, which is a router that keeps the location information for the mobile node when it moves to a foreign network, stores the mobile-node *care-of address* when the mobile node is away from home, and performs other functions on behalf of the mobile node when it is away from home. A peer node that the mobile node communicates with is defined as the *Correspondent Node* (which may be stationary or mobile).

Security between the mobile node and home agent can be accomplished using the *IP Security Protocol* (IPSec) architecture. This permits secure communications between the mobile node and the home agent. When a correspondent node receives a packet from a mobile node, it first checks its binding caches to see if it has a cache of the mobile-node care-of address, and if it does not, the correspondent node sends the packet to the mobile-node home address. The home agent receives all packets sent to the mobile node when it is away from home and then tunnels the packets to the mobile-node care-of address.

To permit a mobile node and correspondent node to communicate directly, without going through a home agent, requires the use of *Mobile IPv6 Route Optimization*. First the connection to the correspondent node needs to be secure from the home agent and directly from the mobile node. In the figure, that is done using a procedure defined as *Return Routability* (RR) within the Mobile IPv6 protocol. The network path between the mobile node and correspondent node is secured through the RR procedure.

Mobile IPv6 uses the extensibility of the IPv6 protocol defining new Neighbor Discovery messages and types, *Routing Header*, and the use of the *Destination Option* in an IPv6 packet, which does not exist in IPv4. Discussion of those extensions is beyond the scope of this article, and is left as an exercise for readers to read the actual Mobile IPv6 specification.

Mobile IPv6 has core technical operational advantages over Mobile IPv4, as follows:

- There is no need to deploy special routers as “foreign agents,” as in Mobile IPv4. Mobile IPv6 operates in any location without any special support required from the local router.
- Support for route optimization is a fundamental part of the protocol, rather than a set of nonstandard extensions.
- Mobile IPv6 route optimizations can operate securely even without prearranged security associations. It is expected that the route optimizations can be deployed on a global scale among all mobile-node correspondent nodes.
- Support is also integrated into Mobile IPv6 for allowing route optimizations to coexist with routers that perform ingress filtering.
- The IPv6 *Neighbor Unreachability Detection* assures symmetric reachability between the mobile node and its default router in the current location.
- Most packets sent to a mobile node away from home in Mobile IPv6 are sent using an IPv6 routing header rather than IP encapsulation, reducing the amount of resulting overhead compared to Mobile IPv4.
- Mobile IPv6 is decoupled from any particular link layer because it uses IPv6 Neighbor Discovery instead of IPv4 *Address Resolution Protocol* (ARP). This also improves the robustness of the protocol.
- The use of IPv6 encapsulation (and the routing header) removes the need in Mobile IPv6 to manage tunnel soft state.
- The dynamic home-agent address discovery mechanism in Mobile IPv6 returns a single reply to the mobile node. The directed broadcast used in IPv4 returns separate replies from each home agent.

Summary

This article has presented three of the key technology advantages of IPv6 behind the wall. There are others, but they are technically too complex to define in a short article, but rather the subject of IPv6 implementation white papers. The IPv6 architecture extends the potential for the Next-Generation Internet to support rapid renumbering of networks, Quality of Service, extensions for ad hoc networks, and the hope of extending the Internet beyond the capabilities and functions today with IPv4. Most important is that IPv6 enhancements will be developed without using “band aids,” as is currently being done with today’s IPv4 architecture. The author of this article would like to thank Tony Hain and Patrick Grossetete from Cisco Systems for their review.

For Further Reading

- [1] R. Droms, Ed., J. Bound, B. Volz, T. Lemon, C. Perkins, M. Carney, "Dynamic Host Configuration Protocol for IPv6 (DHCPv6)," RFC 3315, July 2003.
- [2] R. Hinden, S. Deering, "Internet Protocol Version 6 (IPv6) Addressing Architecture," RFC 3513, April 2003.
- [3] R. Hinden, S. Deering, E. Nordmark, "IPv6 Global Unicast Address Format," RFC 3587, August 2003.

Additional information regarding IPv6 can be found at the International IPv6 Forum Web site www.ipv6forum.com and the North American IPv6 Task Force Web site www.nav6tf.org. Specifically, readers can view the IPv6 Forum basic value proposition at:

http://www.nav6tf.org/summit_slides/IPv6_Value_Proposition_June_2003final.ppt

JIM BOUND works at Hewlett Packard Corporation as an HP Fellow and is a Network Technical Director within the Enterprise UNIX (HP-UX) Division's Network and Security Lab Engineering Group. Jim was a member of the Internet Protocol Next Generation (IPng) Directorate within the IETF, which selected IPv6, among several proposals, to become the basis of the IETF's work on an IPng in 1994. Jim has been a key designer and implementor of IPv6, and contributor and coauthor of IPv6 specifications. Jim founded an ad-hoc IPv6 deployment group working with implementors across the Internet in 1998, which became the IPv6 Forum, where Jim is now Chair of the IPv6 Forum Technical Directorate and Member of the Board of Directors. Jim is also Chair of the North American IPv6 Task Force. Jim is a pioneer member of the Internet Society, and member of the Institute of Electrical and Electronics Engineers (IEEE). In July 2001, Jim received the IPv6 Forum Internet IPv6 Pioneer Award as the IPv6 Forum's "Lead Plumber." Jim has been working in the field of networking as engineer and architect since 1978, and is a subject matter expert to government and industry, for IPv6 and network-centric technology. E-mail: jim.bound@hp.com

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Peter T. Kirstein Receives Postel Award

Peter Kirstein is this year's recipient of the prestigious *Jonathan B. Postel Service Award*. A founding member of the Internet Society, Professor Kirstein is one of the pioneers of the Internet and was directly involved with its development and evolution. He was awarded the Postel Service Award in recognition of his foresight, persistence and innovation in navigating international technical and political complexities, and thus enabling the global propagation of the Internet. The Postel Award was presented on July 16, during the 57th meeting of the *Internet Engineering Task Force* (IETF) in Vienna, Austria.

"The Internet Society is pleased to recognize Peter's significant contribution to the development of the Internet by awarding him this year's Postel Award," said Internet Society President/CEO Lynn St. Amour. "His commitment to the evolution and growth of the Internet, particularly during the 1970s, made possible the global infrastructure we have today. And, his efforts continue, most recently working in the Southern Caucasus and Central Asia regions." Steve Crocker, noted Internet authority and chair of this year's Postel award committee, commented on Kirstein's foresight in laying the groundwork for the Internet's global scope. "Peter Kirstein saw that the future of networking lay in international cooperation and interconnection, and deftly organized the steps to make it happen. He used both technical and personal skills and enabled many others to do magnificent work."

In 1973, Kirstein established one of the first two international nodes of the ARPANET, playing a very active part in the ensuing SATNET activity, which covered five countries. His group continued to provide the principal Internet link between the UK and the US throughout the 1980s, during which time he was responsible for both the **.UK** and **.INT** domains. He continues to collaborate in US *Defense Advanced Research Agency* (DARPA) programs. He has led six European projects in computers and communications funded by the European Commission, and participated in twelve more. Currently, he is leading the *Silk Project*, which is providing satellite-based Internet access to the Newly Independent States in the Southern Caucasus and Central Asia. In June, he was awarded a *Commander, Order of the British Empire*, for his services to Internetworking research.

He has chaired the International Collaboration Board, which currently involves six NATO countries, since 1983, and served on the Networking Panel of the *NATO Science Committee* (serving as chair in 2001). He has been on Advisory Committees for the *Australian Research Council*, the *Canadian Department of Communications*, the German GMD, and the Indian *Education and Research Network* (ERNET) Project. Kirstein obtained his undergraduate degree in Mathematics and Engineering from Gonville and Caius College, Cambridge University, his PhD in Electrical Engineering from Stanford University, and was awarded a DSc in Engineering from the University of London.

Kirstein expressed his appreciation for the award and respect for Jon Postel's work, explaining, "Postel's efforts to ensure the successful development and deployment of the Internet was an inspiration to us all. His stewardship of the RFC series was essential to the successful development of the Internet. His conscientious and painstaking operation of the Domain Name System and the Internet Assigned Numbers Authority were indispensable to the international growth of the system. I am particularly pleased to be recipient of an award in his name, and feel greatly honored to be considered worthy of having my activities linked with his memorial."

The Jonathan B. Postel Service Award was established by the Internet Society to honor those who have made outstanding contributions in service to the data communications community. The award is focused on sustained and substantial technical contributions, service to the community, and leadership. With respect to leadership, the nominating committee places particular emphasis on candidates who have supported and enabled others in addition to their own specific actions.

The award is named after Dr. Jonathan B. Postel, who embodied all of these qualities during his extraordinary stewardship over the course of a thirty-year career in networking. He served as the editor of the RFC series of notes from its inception in 1969, until 1998. He also served as the ARPANET "numbers Czar" and the Internet Assigned Numbers Authority over the same period of time. He was a founding member of the *Internet Architecture* (nee *Activities*) *Board* (IAB) and the first individual member of the Internet Society, where he also served as a trustee.

Previous recipients of the Postel Award include Jon himself (posthumously and accepted by his mother), Scott Bradner, Daniel Karrenberg and Stephen Wolff. The award consists of an engraved crystal globe and \$20,000.

The *Internet Society* (ISOC) (www.isoc.org) is a not-for-profit membership organization founded in 1991 to provide leadership in Internet related standards, education, and policy. With offices in Washington, DC, and Geneva, Switzerland, it is dedicated to ensuring the open development, evolution and use of the Internet for the benefit of people throughout the world. ISOC is the organizational home of the IETF, the IAB, the *Internet Engineering Steering Group* (IESG) and other Internet-related bodies who together play a critical role in ensuring that the Internet develops in a stable and open manner. For over 12 years ISOC has run international network training programs for developing countries and these have played a vital role in setting up the Internet connections and networks in virtually every country connecting to the Internet during this time.

Deployment of Internationalized Domain Names

The *Internet Corporation for Assigned Names and Numbers* (ICANN) recently announced the commencement of global deployment of *Internationalized Domain Names* (IDNs)^[2,3,4], which will allow use on the Internet of domain names in languages used in all parts of the world.

In October 2002, the IESG approved the publication of a standardized way of integrating IDNs into the Internet's *Domain Name System* (DNS). After the proposed technical standard was published in March 2003, the ICANN Board endorsed an approach for implementation of the technical standard that had been developed cooperatively by ICANN and leading IDN registries.

Following up on the Board's endorsement, ICANN and the leading IDN registries finalized an agreed text of the principles to be followed in IDN registration activities. Those "Guidelines for the Implementation of Internationalized Domain Names"^[1] were published. IDN registries adhering to the Guidelines will employ language-specific registration and administration rules that are documented and publicly available. These IDN registries will work collaboratively with each other and with interested stakeholders to develop the language-specific policies, with the objective of achieving consistent approaches to IDN implementation to maintain Internet interoperability for the benefit of DNS users worldwide.

The registries for the **.cn** (China), **.jp** (Japan), and **.tw** (Taiwan) country codes, as well as for the **.info** and **.org** generic top-level domains, have committed to adhere to the Guidelines. As authorized by the ICANN Board in March, registries seeking to deploy IDNs under their agreements with ICANN will be authorized to do so on the basis of the Guidelines. In addition, the ICANN Board has recommended the Guidelines to other registries, and encourages broad participation by registries, language experts, and others in consultative, collaborative, community-based processes to study and develop appropriate language-specific IDN registration rules and policies.

As the deployment of IDNs proceeds, ICANN and the participating IDN registries have agreed to work together to review Guidelines at regular intervals based on their deployment experience, and to make any necessary adjustments.

[1] <http://www.icann.org/general/idn-guidelines-20jun03.htm>

[2] P. Faltstrom, P. Hoffman, A. Costello, "Internationalizing Domain Names in Applications (IDNA)," RFC 3490, March 2003.

[3] P. Hoffman, M. Blanchet, "Nameprep: A Stringprep Profile for Internationalized Domain Names (IDN)," RFC 3491, March 2003.

[4] A. Costello "Punycode: A Bootstring encoding of Unicode for Internationalized Domain Names in Applications (IDNA)," RFC 3492, March 2003.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Technology Strategy
MCI, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.
Copyright © 2003 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID Cisco Systems, Inc.

The Internet Protocol Journal

December 2003

Volume 6, Number 4

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
IPv4: How long do we have? ...	2
Low-tech Network Maintenance	16
Letters to the Editor	23
Book Review	25
Fragments	28

FROM THE EDITOR

I will remember 2003 as the year when high-speed Internet access became widely available in public locations such as airports, hotels, and coffee shops. As a frequent traveler, I really appreciate not having to find a suitable telephone jack and corresponding country-specific telephone adapter plug in order to get my e-mail. The IEEE 802.11 “WiFi” standard has truly arrived. I even stayed in a new hotel in Norway that provided WiFi access in every room by placing base stations in the hallways. When I first stepped into my hotel room and noticed that it had only a *digital* telephone and no sign of any Ethernet jacks I worried, but a quick check revealed that I could purchase a scratch-off card at reception that provided me with a username and password valid for 24 hours. A clear example of a “technology generation leap.”

The year 2003 was also the year in which unsolicited e-mail, or “spam,” became a major problem for all Internet users. Various filtering systems have thankfully been devised and deployed, but this problem has no easy solution. It will be interesting to see what impact new antispam legislation will have over the coming months and years.

The first article presents an in-depth look at the IP Version 4 address space and its measured and projected consumption rate. When work first started on the design of IP Version 6, projections indicated that we’d run out of IPv4 addresses within a few years. Geoff Huston takes a fresh look at this in an article entitled “IPv4—How long do we have?”

The job of System Administrator, or “sysadmin,” is a challenging one, and if your job includes keeping the network running 24 hours a day, you will probably appreciate some of the tips in our second article, entitled “Low-Tech Network Maintenance.”

For the second time recently, Queen Elizabeth II has honored an Internet pioneer. Tim Berners-Lee, the inventor of the World Wide Web and director of the *World Wide Web Consortium* (W3C), was made a *Knight Commander, Order of the British Empire* in the 2004 New Years Honours list. (See “Fragments,” page 28).

Which brings us to the IPJ publication schedule. If you are a regular subscriber to the IPJ, you probably have noticed a somewhat irregular publishing schedule in 2003. This December 2003 issue is indeed being published in January 2004. This results from our effort to produce timely quality articles in a world where the experts are not staff writers. Of course, you should still expect to receive four issues per year, and your feedback to ipj@cisco.com will help make IPJ even better.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

IPv4—How long do we have?

by Geoff Huston, Telstra

One of those stories that keeps on appearing from time to time is the claim that somewhere in the world, or even all over the world, we are “running out of IP addresses,” referring to the consumption of unallocated IPv4 addresses^[1]. In one sense this is a pretty safe claim, in that the IPv4 address pool is indeed finite, and, as the IPv4 Internet grows it makes continual demands on previously unallocated address space. So the claim that the space will be exhausted at some time in the future is a relatively safe prediction. But the critical question is not “if” but “when,” because this is a question upon which many of our current technology choices are based.

Given this revived interest in the anticipated longevity of the IPv4 address space, it is timely to revisit a particular piece of analysis that has been a topic of some interest at various times over the past decade or more. The basic question is: “How long can the IPv4 address pool last in the face of a continually growing network?” This article looks at one approach to attempt to provide some indication of “when.” Like all predictive exercises, many assumptions have to be made, and the approach described here uses just one of numerous possible predictive models—and, of course, the future is always uncertain.

The IPv4 Address Space

The initial design of IPv4 was extremely radical for its time in the late 1970s. Other contemporary vendor-based computer networking protocols were designed within the constraints of minimizing the packet header overhead in order to improve the data payload efficiency of each packet. At the time address spans were defined within the overall assumption that the networks were deployed as a means of clustering equipment around a central mainframe. In many protocol designs 16 bits of address space in the packet headers was considered to be extravagant. To use a globally unique address framework of 32 bits to address network hosts was, at the time, a major shift in thinking about computer networks from a collection of disparate private facilities into a truly public utility.

To further add to the radical nature of the exercise, the Internet Network Information Center was prepared to hand out unique blocks of this address space to anyone who submitted an application. Address deployment architectures in other contemporary protocols did not have the address space to support such address distribution functions, nor did they even see a need for global uniqueness of computer network addresses. Network administrators numbered their isolated corporate or campus networks starting at the equivalent of “1,” and progressed onward from there. Obviously network splits and mergers caused considerable realignment of these private addressing schemes, with consequent disruption to the network service.

By comparison, it seemed, the address architecture of the Internet was explicitly designed for interconnection. But even with 32 bits to use in an address field, getting the right internal structure for addresses is not as straightforward as it may initially seem.

The Evolution of the IPv4 Address Architecture

IP uses the address to express two aspects of a connected device: the identity of this device (endpoint identity) and the location within the network where this device can be reached (location or forwarding identity). The original IP address architecture used the endpoint identity to allow devices to refer to each other in end-to-end application transactions, whereas within the network the address is used to direct packet-forwarding decisions. The address was further structured into two fields: a *network* identifier and a *host* identifier within that network. The first incarnation of this address architecture used a division at the first octet: the first 8 bits were the network number and the following 24 bits were the host identifier. The underlying assumption was one of deployment across a small number of very large local networks. This view was subsequently refined, and the concept of a class-based address architecture was devised for the Internet. Half of the address space was left as a 8/24-bit structure, called the *Class A* space (allowing for up to 127 networks each with 16,777,216 host identities). A quarter of the remaining space used a 16/16-bit split (allowing for up to 16,128 networks, each with up to 65,536 hosts), defining the *Class B* space. A further eighth of the remaining space was divided using a 24/8-bit structure (allowing for 2,031,616 networks, each with up to 256 hosts), termed the *Class C* space. The remaining eighth of the space was held in reserve.

This address scheme was devised in the early 1980s, and within a decade it was pretty clear that there was a problem with impending exhaustion. The reason was an evident run on Class B addresses. Although very few entities could see their IP network spanning millions of computers, the personal desktop computer was now a well-established part of the landscape, and networks of just 256 hosts were just too small. So if the Class A space was too big, and the Class C too small, then Class B was the only remaining option. In fact, the Class B blocks were also too large, and most networks that used a Class B address consumed only a few hundred of the 65,535 possible host identities within each network. The addressing efficiency of this arrangement was very low, and a large amount of address space was being consumed in order to number a small set of devices. Achieving even a 1 percent host density (expressed as a ratio of number of addressed hosts to the total number of host addresses available) was better than normal at the time, and 10 percent was considered pretty exceptional.

Consequently, Class B networks were being assigned to networks at an exponentially increasing rate. Projections from the early 1990s forecast exhaustion of the Class B space by the mid-1990s. Obviously there was a problem, and the *Internet Engineering Task Force* (IETF) took on the task of finding some solutions. Numerous responses were devised by the IETF.

As a means of mitigation of the immediate problem, the IETF altered the structure of an IP address. Rather than having a fixed-length network identifier of 8, 16, or 24 bits, the network part of the address could be any length at all, and a network identifier was now the couplet of an IP address field containing a network part and the bit length of the network part. The boundary between the network and host part could change across the network, so rather than having “networks” and “subnetworks” as in the class-based address architecture, there was the concept of a variable length network mask. This was termed the “classless” address architecture (or “CIDR”), and the step was considered to be a short-term expediency to buy some additional time before address exhaustion. The longer-term plan was to develop a new IP architecture that could encompass a much larger connectivity domain than was possible with IPv4.

We now have IPv6 as the longer-term outcome. But what has happened to the short-term expediency of the classless address architecture in IPv4? It appears to have worked very well indeed so far, and now the question is: how long can this supposedly short-term solution last?

Predictions of Address Consumption

Predicting the point of IPv4 address exhaustion has happened from time to time since the early 1990s within the IETF^[2]. The initial outcomes of these predictive exercises were clearly visible by the mid-1990s: the classless address architecture was very effective in improving the address utilization efficiency, and the pressures of ever-increasing consumption of a visibly finite address resource were alleviated. But a decade after the introduction of CIDR addressing, it is time to understand where we are heading with the consumption of the underlying network address pool.

Dividing up the Address Space

There are three stages in address allocation. The pool of IP addresses is managed by the *Internet Assigned Numbers Authority* (IANA). Blocks of addresses are allocated to *Regional Internet Registries* (RIRs), who in turn allocate smaller blocks to *Local Internet Registries* (LIRs) or *Internet Service Providers* (ISPs).

Currently 3,707,764,736 addresses are managed in this way. It is probably easier to look at this in terms of the number of “/8 blocks,” where each block is the same size as the old Class A network, namely 16,777,216 addresses. The total address pool is 221 /8s, with a further 16 /8s reserved for multicast use, 16 /8s held in reserve, and 3 /8s designated as not for use in the public Internet.

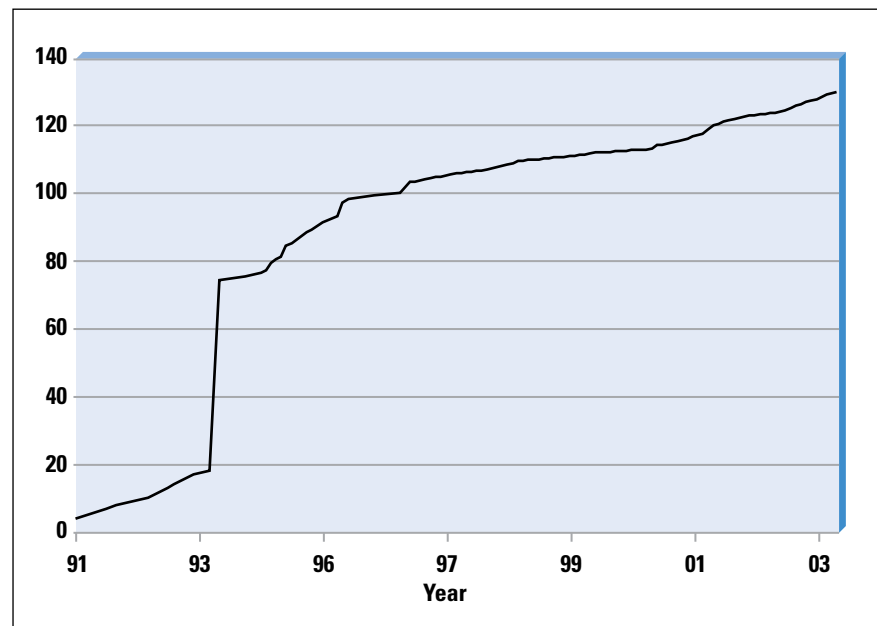
In looking at futures, there are three sources of data concerning address consumption:

- How quickly is the IANA passing address blocks to the RIRs, and when will IANA run out?
- How quickly are the RIRs passing address blocks to LIRs, and when will this run out?
- How much address space is actually used in the global Internet, and how quickly is this growing? When will this run out?

The IANA Registry

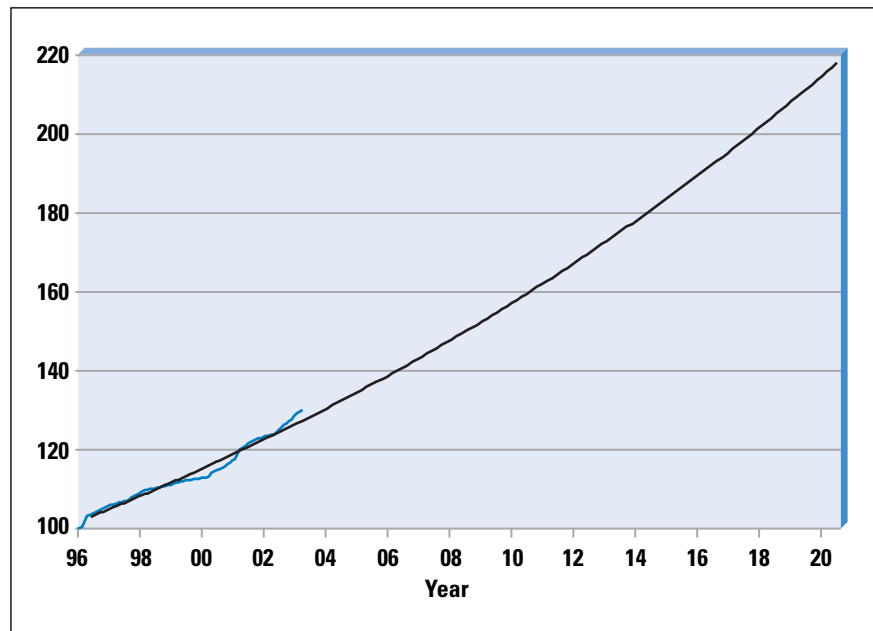
So the first place to look is the IANA registry file^[3]. This registry reveals that of these 221 /8 blocks, 89 /8 blocks are still held as unallocated by the IANA, 129.9 /8 blocks have been allocated, and the remaining 2.1 /8 blocks are reserved for other uses. The IANA registry also includes the date of allocation of the address block, so it is possible to construct a time series of IANA allocations, as shown in Figure 1.

Figure 1: IANA Allocated IPv4 /8 Address Blocks



Interestingly, there is nothing older than 1991 in this registry. This exposes one of the problems with analyzing registry data, in that there is a difference between the current status of a registry and a time-stamped log of the transactions that were made to the registry over time. The data published by the IANA is somewhere between the two, and the log data is incomplete; in addition, the current status of some address blocks is unclear. It appears that the usable allocation data starts in 1995. So if we take the data starting from 1995 and perform a linear regression to find a best fit of an exponential projection, it is possible to make some predictions as to the time it will take to exhaust the remaining unallocated 89 /8s. (Figure 2).

Figure 2: IANA Allocated IPv4 /8 Address Blocks



It is worth a slight digression into the method of projection being used here. The technique is one of using a best fit of an exponential growth curve to the data. The underlying assumption behind such a projection is that the growth rate of the data is proportional to the size of the data, rather than being a constant rate. In network terms, this assumes that the rate of consumption of unallocated addresses is a fixed proportion of the number of allocated addresses, or, in other words, the expansion rate of the network is a proportion of its size, rather than being a constant value. Such exponential growth models may not necessarily be the best fit to a network growth model, although the data since 1995 does indicate an underlying exponential growth pattern. Whether this growth model will continue into the future is an open issue.

The projection of 2019 as the date for consumption of the unallocated address space using this technique is perhaps surprising, because it seems that the network is bigger now than ever, yet the amount of additional address space required to fuel further accelerating growth for a further decade is comparatively small. This is true for many reasons, and the turning point when these aspects gained traction in the Internet appeared to be about 1995. They include:

- The first 1.6 billion addresses (equivalent to some 100 /8 blocks) were allocated using the class-based address architecture. Since this date address allocation has used a classless architecture, and this has enabled achievement of significantly improved efficiencies in using the address space.
- The RIRs came into the picture, and started using conservation-based policies in address allocations. The RIR process requires all address applicants to demonstrate that they can make efficient and effective use of the address space, and this has dampened some of the wilder sets of expectations about the address requirements of an enterprise.

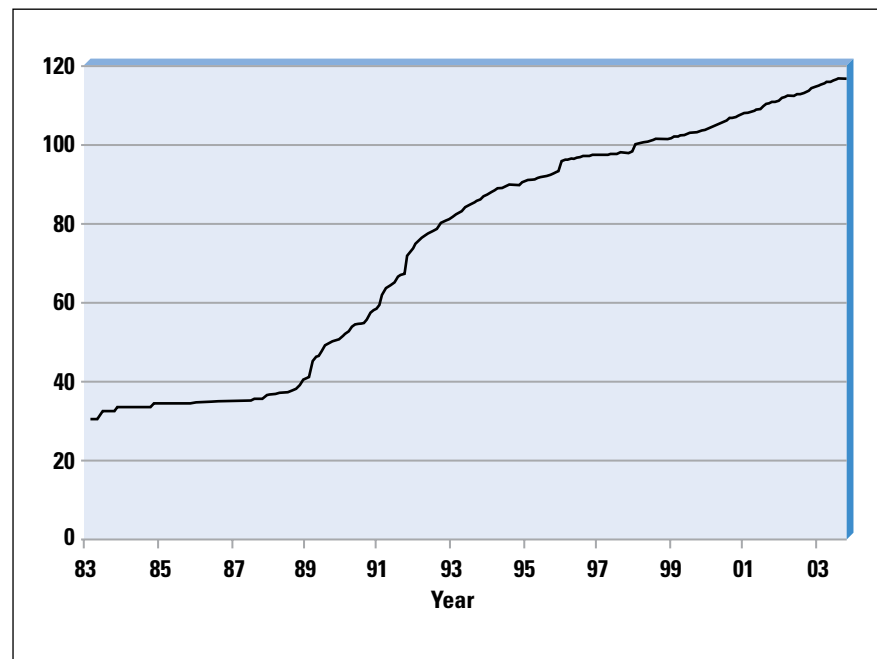
- Address compression technologies became widely deployed. Dynamic *Network Address Translation* (NAT) devices have, for better or worse, become a common part of the network landscape. NAT devices allow large “semi-private” networks to use a very small pool of public addresses as the external view of the network, while using private address space within the network. *Dynamic Host Configuration Protocol* (DHCP) has allowed networks to recycle a smaller pool of addresses across a larger set of intermittently connected devices.

Whether these factors will continue to operate in the same fashion in the future is an open question. Whether future growth in the use of public address space operates from a basis of a steadily accelerated growth is also an open question. The assumption made in this exercise is that the projections depend on continuity of effectiveness of the RIR policies and their application, continuity of technology approaches, and absence of disruptive triggers. Although the RIRs have a very well-regarded track record and there are strong grounds for confidence that this will continue, obviously the latter two assumptions about technology and disruptive events are not all that comfortable. With that in mind, the next step is to look at the RIR assignment data.

The RIR Registries

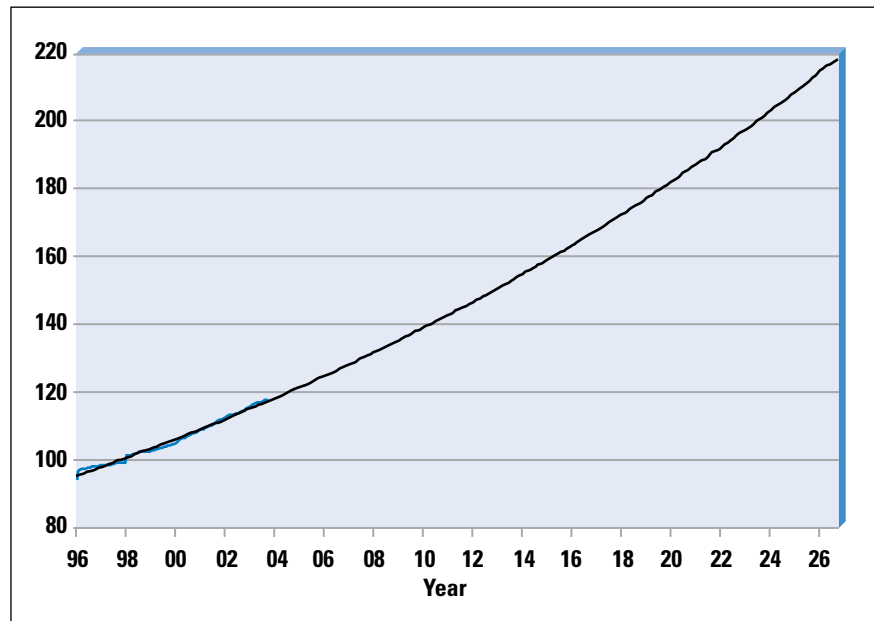
The RIRs also publish a registry of their transactions in “stats” files. For each currently allocated or assigned address block the RIRs have recorded, among other items, the date of the RIR assignment transaction that assigned an address block to a LIR or ISP. Using this data we can break up the 129.9 /8 blocks further, and it is evident that the equivalent of 116.7 /8 blocks have been allocated or assigned by the RIRs, and the remaining space, where there is no RIR allocation or assignment record, is the equivalent of 13.2 /8 blocks. These transactions can again be placed in a time series, as shown in Figure 3.

Figure 3: RIR Assigned IPv4 /8 Address Blocks



The post-1995 data used to extrapolate forward using the same linear regression technique described previously to find a curve of best fit using the same underlying growth model assumptions yields:

Figure 4: RIR Assigned IPv4 /8 Address Blocks—Projection



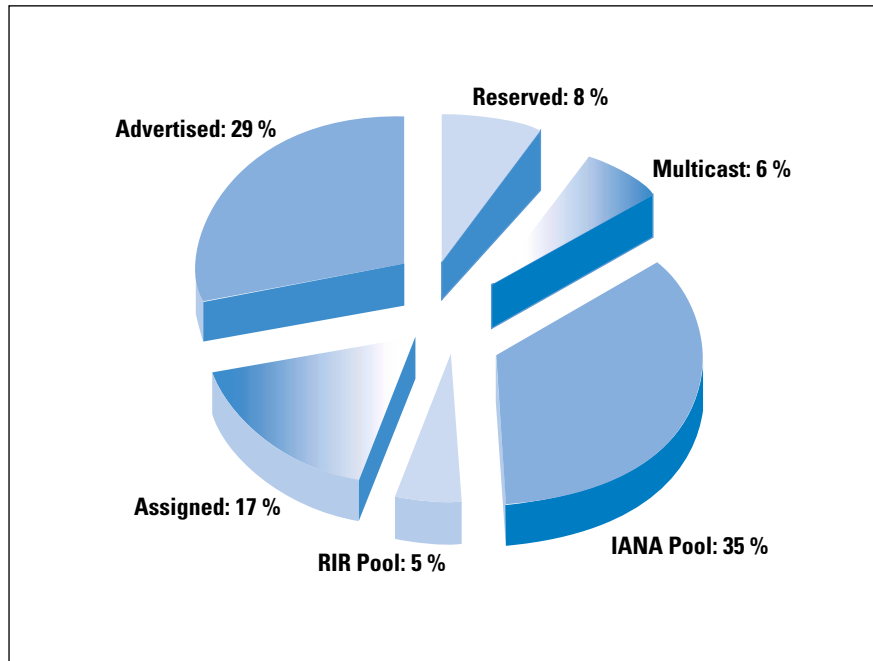
This form of extrapolation gives a date of 2026 for the time at which the RIRs will exhaust the number pool. Again the same caveats about the use of this approach as a reliable predictor apply here, and the view forward is based on the absence of large-scale disruptions, or some externally induced change in the underlying growth models for address demand.

The BGP Routing Table

When addresses are assigned to end networks, the expectation is that these addresses will be announced to the network in the form of routing advertisements. So some proportion of these addresses is announced in the Internet routing table. The next task is to establish the trends of the amount of address space covered by the routing table. The approach used has been to take a single view of the address span of the Internet. This is the view from one point, inside the AS1221 network operated by Telstra.

The data as of October 2003 shows that some 29 percent of the total IPv4 address space is announced in the *Border Gateway Protocol* (BGP) routing table, whereas 17 percent has been allocated to an end user or LIR but is not announced on the public Internet as being connected and reachable. A total of 5 percent of the address space is held by the RIR's pending assignment or allocation (or at least there is no RIR recorded assignment of the space), while 35 percent of the total space remains in the IANA unallocated pool. A further 8 percent of the space is held in reserve (Figure 5).

Figure 5: IPv4 /8 Address Space



This BGP data is based on an hourly inspection of the amount of address space advertised within the Internet routing table. The data collection commenced in late 1999, and the data gathered so far is shown in Figure 6. The problem with this data is that there is some considerable amount of fluctuation in the amount of address space advertised over time. The major step changes are due to a small number of /8 advertisements that periodically are announced and withdrawn in BGP. In order to obtain reasonable data for generating projections, some noise reduction on this data needs to be undertaken. The approach used has been to first filter the data using a constant value of 18 /8 prefix announcements, and then use a sliding average function to create a smoothed time series. This is indicated in Figure 7.

The critical issue when using this data for projection is to determine what form of function can provide a best fit to the data. A good indication of the underlying trends in the data can be found by analyzing the first-order differential of the data. An underlying increasing growth model would have an increasing first-order differential, whereas a decreasing growth model would have a negatively inclined differential. A least-squares best-fit analysis of the data shows that the growth rates have not been consistent over the past three years. A reasonable fit for this data appears to be a constant growth model, or a linear growth projection, with a consumption rate of some 3 /8 blocks per year.

Figure 6: Advertised IPv4 /8 Address Space (/8 Blocks)

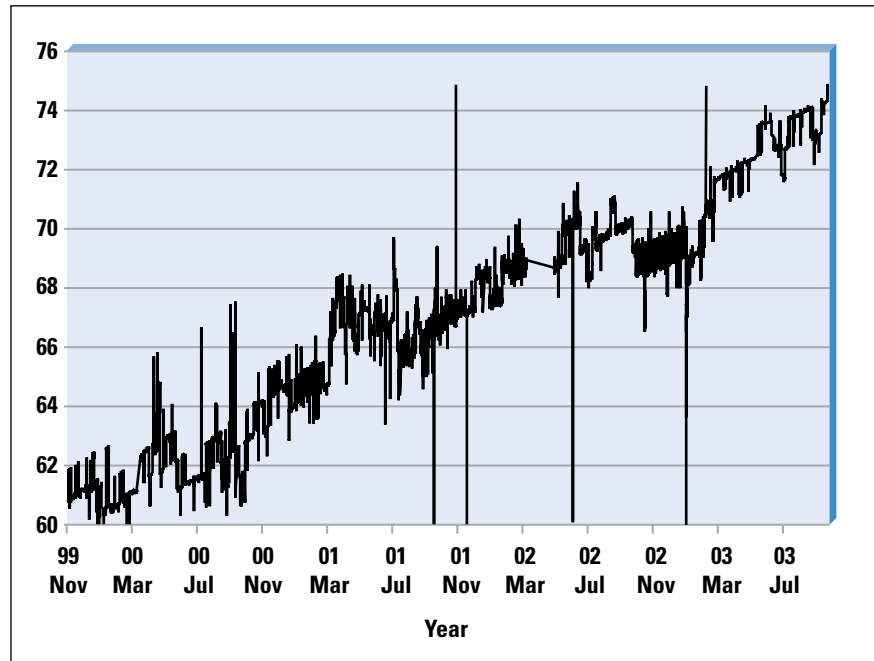
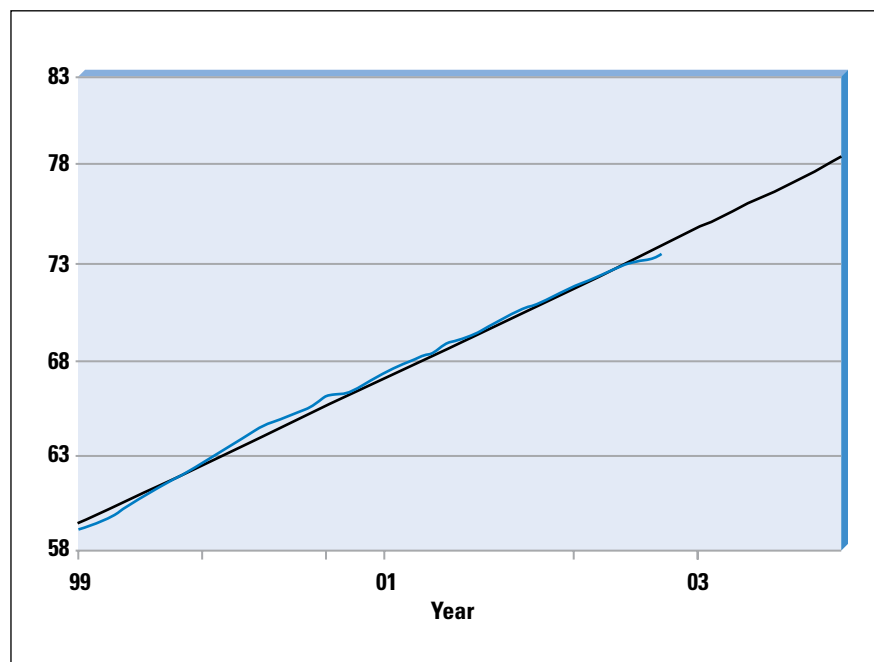


Figure 7: Smoothed IPv4 /8 Advertised Address Blocks



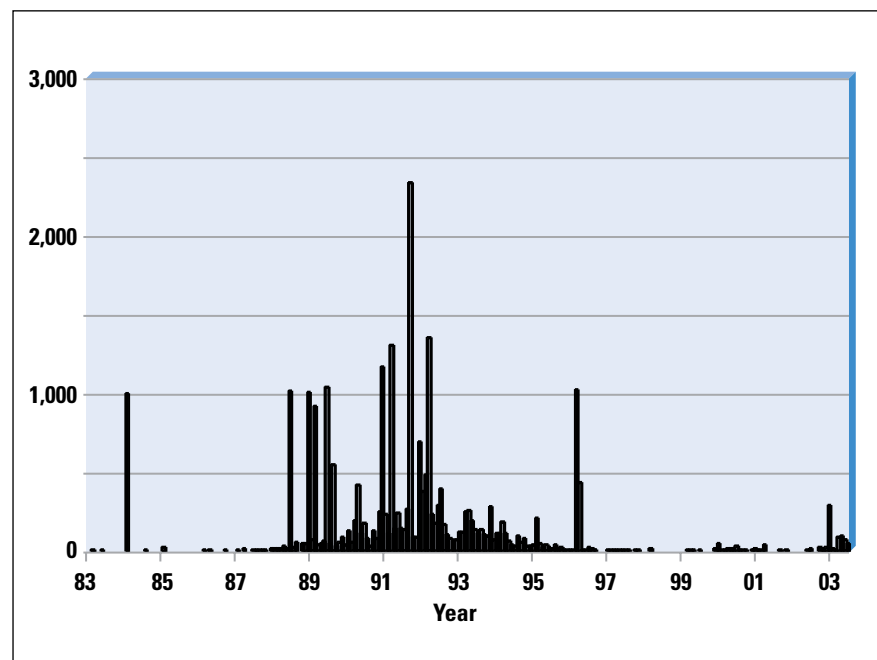
Combining the Three Views

One question remains before we complete the projections for IPv4 address space. There are 43.3 /8 blocks, or some 17 percent of the total IPv4 address space that has been allocated for use, but is not visible in the Internet routing table. This is a very significant amount of address space, and if it is growing at the same rate as the advertised space, then this will have a significant impact on any overall model of consumption of the use of address space.

The question here is whether this “invisible” address pool is a legacy of the address allocations policies in place before the RIR system came into operation in the mid 1990s, or some intrinsic inefficiency in the current system. If it is the latter, then it is likely that this pool of unannounced addresses will grow in direct proportion to the growth in the announced address space, whereas if it is the former, then the size of the pool will remain relatively constant in the future.

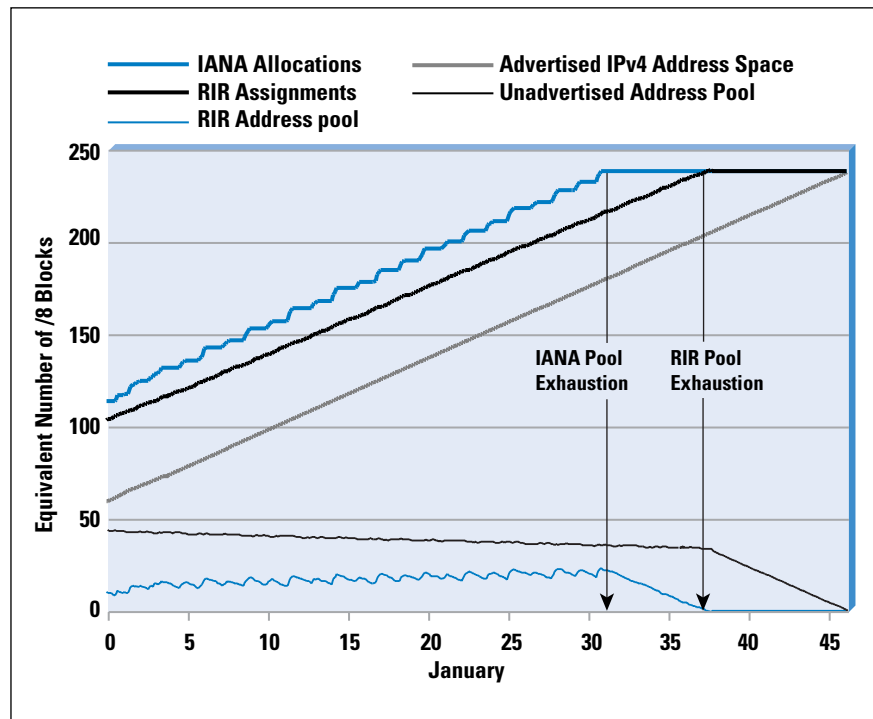
We can look back through the RIR allocation data and look at the allocation dates of unannounced address space (Figure 8). This view indicates that the bulk of the space is a legacy of earlier address allocation practices, and that since 1997, when the RIR operation was fully established, there is an almost complete mapping of RIR allocated address space to BGP routing announcements. The recent 2003 data indicates that there is some lag between recent allocations and BGP announcements, most probably due to the time lag between an LIR receiving an allocation and subsequent assignments to end users and advertisement in the routing table.

Figure 8: Age Distribution of Unadvertised Address Blocks (/8 Address Blocks)



This confirms that in recent years all the address space that has been assigned by the RIRs appears in the Internet routing table, implying that projections of the amount of address space advertised in the routing table is a good correlation to projections of address space consumption. With this in mind it is now possible to construct a model of the address distribution process, working backward from the BGP routing table address size. From the sum of the BGP table size and the LIR holding pool, we can derive the total RIR-managed address pool. To this number is added the RIR holding pool low size and its low threshold where a further IANA allocation is required. This allows a view of the entire system, projected forward over time, where the central driver for the projection is the growth in the network itself, as described by the size of the announced IPv4 address space. This is shown in Figure 9.

Figure 9: IPv4 Projections of Address Consumption



It would appear that the point of effective exhaustion is the point where the RIRs exhaust available address space to assign. In this model, RIR exhaustion of the unallocated address pool would occur in 2037.

Uncertainties

Of course such projections are based on the underlying assumption that tomorrow will be much like today, and the visible changes that have occurred in the past will smoothly translate to continued change the future. This assumption obviously has some weaknesses, and many events could disrupt this prediction.

Some disruptions could be found in technology evolution. An upward shift in address take-up rates could occur because of an inability of NAT devices to support emerging popular applications. Widespread deployment of peer-to-peer applications implies the need for persistent address presentation, which may imply greater levels of requirement for public address space. The use of personal mobile IP devices (such as PDAs in their various formats) using public IPv4 addresses would place a massive load on the address space, simply because of the very large volumes associated with deployment of this technology^[4].

Other disruptions have a social origin, such as the boom and bust cycle of Internet expansion in recent years. Another form of disruption in this category could be the adoption of a change in the distribution function. The current RIR and LIR distribution model has been very effective in limiting the amount of accumulation of address space in holding pools, and allocating addresses based on efficiency of utilization and conformance to the routing topology of the network.

Many other forms of global resource distribution use a geopolitical framework, where number blocks are passed to national entities, and further distribution is a matter of local policy^[5]. The disruptive nature of such a change would be to immediately increase the number of “holding” points in the distribution system, locking away larger pools of address space from being deployed and advertised and generating a significant upward change in the overall address consumption rates due to an increase in the inefficiency of the altered distribution function.

The other factor to be aware of is the steadily decreasing “buffer” of unallocated addresses that can be used to absorb the impacts of a disruptive change in address consumption rates. Although at present some 60 percent of the address space—or some 2.6 billion addresses—are available in the unallocated address pools or held in reserve, this pool will reduce over time. If a disruptive event is, for example, a requirement to directly address some 500 million devices, then such an event would reduce the expectancy of address space availability by some years, assuming it occurred within the period when sufficient address space remains to meet such a surge of demand.

The other source of uncertainty is that this form of predictive modeling assumes that the ratios of actual connected devices and the amount of address space deployed to service this device pool remain relatively constant.

This model also assumes some form of continuity of current address allocation policies. This is not a likely scenario, because it is likely that address policies will reflect some notion of balance between the level of current demand against future demands. As the unallocated address pool shrinks it is possible that policies will alter to express the increased level of competitive demand for the remaining resource. Consumption rates would be moderated by such a change in allocation policy. The commonly cited intended evolutionary path for the Internet is to a transition to ubiquitous use of IPv6, and at some point in that transition process it is reasonable to assume that further demands for IPv4 space will dwindle. It may be that at such a “crossover” time allocation policies may then be altered to reflect a drop in both current and future demands for IPv4 address space.

In attempting to assess the possible future path of address allocation policies, it is also evident that, from a market rationalist perspective, there is a certain contrivedness about the current address allocation process. The current address management system assumes a steady influx of new addresses to meet emerging demands, and the overall address utilization efficiency is not set by any form of market force, but by the outcomes of the application of RIR address allocation policies to new requests for address space. A market rationalist could well point to the use of market price as a means of determining the most economically efficient form of utilization of a commodity product. Such a position is based on the observation that the way that the consumer chooses between alternative substitutable services is by a market choice that is generally price sensitive.

If price is removed from an IPv4 address market, the choices made by market players are not necessarily the most efficient choices, and some would argue that the current situation underprices IPv4 at the expense of IPv6.

However, in venturing into these areas we are perhaps straying a little too far from exploring the degree of uncertainty in these predictive exercises. A discussion of the interaction between various forms of distribution frameworks and likely technology outcomes is perhaps a topic for another time.

So just how long does IPv4 have?

The assumptions used here include assuming that the trends in the growth in the advertised space are directly proportional to the future consumption rates for IP addresses, and that the constant growth model remains a best fit for this time series of data. It also assumes a continuation of the current utilization efficiency levels in the Internet, a continuing balance between public address utilization and the use of various forms of address compression, and continuity of current address allocation policies, as well as the absence of highly disruptive events. With all this in mind, then it would appear that the IPv4 world, in terms of address availability, could continue for up to another three decades or so without reaching any fixed boundary of exhaustion.

But it must be remembered that each of these assumptions is relatively sweeping, and to combine them as we have done here is pushing the predictive exercise to its limits, or possibly beyond them. Three decades out is way over the event horizon for any form of useful prediction for the Internet, so if we restrict the question to at most the next five to eight years, then we can answer with some level of confidence that, in the absence of any significant disruptions to the current deployment model of the Internet, there is really no visible evidence that IPv4 will exhaust its address pool by 2010, based on the available address consumption data.

Data Sources

IANA IPv4 Address Registry:

<http://www.iana.org/assignments/ipv4-address-space>

Registry “stats” report files:

APNIC: **<ftp://ftp.apnic.net/pub/apnic/stats>**

ARIN: **<ftp://ftp.arin.net/pub/stats>**

LACNIC: **<ftp://ftp.lacnic.net/pub/stats>**

RIPE NCC: **<ftp://ftp.ripe.net/ripe/stats>**

BGP Address Data: **<http://bgp.potaroo.net>**

Notes

- [1] “Tackling the net’s number shortage.” BBC News, World Edition, 26 October 2003. The item starts with the claim: “BBC ClickOnline’s Ian Hardy investigates what is going to happen when the number of net addresses—Internet Protocol numbers—runs out sometime in 2005.”
<http://news.bbc.co.uk/2/hi/technology/3211035.stm>
- [2] The work was undertaken in the *Address Lifetime Expectations* (ALE) Working Group of the IETF in 1993–1994. The final outcome from this effort was reported from the December 1994 meeting of this group: “Both models currently suggest that IPv4 addresses would be depleted around 2008, give or take three years.”
- [3] This registry is online at:
<http://www.iana.org/assignments/ipv4-address-space>
- [4] On the other hand, it is evident that the growth of the Internet in recent years has been fueled by the increasing prevalence of NAT devices. In order for applications to be accepted into common use in today’s Internet, they need to be able to function through various NAT-based constraints, and increasing sophistication of applications in operating across NAT devices is certainly evident today.
- [5] Such a geopolitical distribution system is used in the E.164 number space for telephony (“ENUM”).

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. Huston is currently the Chief Scientist in the Internet area for Telstra. He is also the Executive Director of the Internet Architecture Board, and is a member of the APNIC Executive Committee. He is author of *The ISP Survival Guide*, ISBN 0-471-31499-4, *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, and coauthor of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. All three books are published by John Wiley & Sons. E-mail: **gih@telstra.net**

Low-Tech Network Maintenance

by *Locum sysadmin*

In an ideal world, we all maintain networks composed of shiny, high-end equipment. Server rooms are stacked to the brim with racks of blinking lights. Neat bundles of cable wend their way through cable loops to orderly, labeled patch bays. When the occasional piece of equipment fails, a hot replacement is slotted in by trained technicians, often before users even notice the outage. Sleek, modern servers hum contentedly, offering their services all day, every day. All is well.

And then there are the other environments ...

Imagine, if you will, that you are a programmer, working for a small company. You are perhaps vaguely aware that all is not well with the small network that you use each day, but the system administrator (*sysadmin*, if there is one) is so busy with other duties that addressing your concerns seems to be last on the list. The occasional delay in CVS checkouts or e-mail that just never quite makes it seem like minor issues compared to... well, whatever it is that so occupies the sysadmin.

Or perhaps there is no sysadmin ... the network topology is neither ring, nor star, but more “accreted.” It is possible that the nephew of one of the managers was responsible for its setup. Like coral, successive waves of employees have washed over the network, leaving their small additions—a cheap 8-port hub here, some gaffer-taped wiring there.

You become aware that your LAN/WAN environment is a real-world test of how deeply Ethernet hubs may be cascaded. A trip to the server room (or server closet) reveals a mess of cabling that closely resembles blue spaghetti. Access to the outside world can take several forms, but it is not uncommon to find a couple of dialup modems lurking quietly in the mess, unnoticed until a failure in the regular link means a failover to the pleasures of 30 employees sharing a 33.6k modem. The concept of labeling cables never made it to this paleolithic theme park, so if you ever trip on one of the floor-dwelling blue vines, locating its original socket can be a challenging occupation.

The servers themselves seem to be an interactive museum display charting the history of computing up until the late 1990s. Old UNIX boxes spill a mess of cables and hard drives over the bench, generic white-box servers of unknown vintage litter the room, “Powered by Linux” or FreeBSD stickers adorning them. Discolored 15-inch monitors sometimes display a blue screen of death, letting you know that some people still love NT4. Assorted tape drives blink quietly away, backing up regularly, though no one seems quite sure what they are backing up, or how to recover them. An elderly Sun box whiles away its retirement transferring mail and playing host to the occasional crackers who exploit security holes in its ancient *sendmail*, then give up in disgust.

The spare parts for the network might occupy a shelf in the server room, or perhaps they nestle on top of a rack unit. A motley assortment of chewed-looking Category 5 cables, network cards so ancient that their manufacture date is in Roman numerals, and a sculpture of BNC connectors—the thought of turning here for help fills you with dread. A dead network adapter usually means a surreptitious raid of the petty cash and a trip to the local computer-parts store for a no-name Ethernet card.

Then—as it always does—disaster strikes. Somewhere, something goes wrong. One thing that you can be sure of is that it will happen at the worst possible time. It is likely that a crucial presentation will be under way, or perhaps a software release is due by close of business. Maybe you are hosting a server for a client, and the client has noticed its absence, and is on the phone, using words like “unscheduled outage” and “penalty clause.” If your clients are so inclined, words like “kneecap” and “sledgehammer” might also be heard. Another fact you can be reasonably sure of is that the sysadmin will not be present, and the next-most technical person will be called upon to work up a minor miracle to fix the ailing network.

Sound far-fetched? Believe it or not, I have been in this situation more than once. What follows are some hints that may help in fixing networks in suboptimal conditions, and as always, with the understanding that it must be done as cheaply as possible.

Many of the hints use features found on Linux boxes, beloved for its technical excellence (and its low cost). Most of the tips here can be adapted for whatever type of operating system you have.

Audible Ping

Ping is the venerable tool that we all know and love, and is the reigning king of the low-tech diagnostic tools. Linux (and other operating systems that use GNU tools) features an extension to *ping* that produces a beep on receipt of a response. The *audible ping* is designated by the **-a** command-line option.

Something as simple as **ping -a missinghost.your.net**, left running from a console in the server room, can alert you when you have finally reestablished network connectivity. It is like having a cable tester that can traverse routers.

Where Are You?

In a server room full of unlabeled generic boxes, it can sometimes be tricky to know which box is which. The following conversation is typical:

Hapless1: “Okay, I’ve logged into **srv7** by SSH [Secure Shell Protocol], and I think its second hard drive has died. Can you turn off its power switch when I shut it down?”

Hapless2: “Sure, which box is it?”

Hapless1: “Ummm... its hostname is **srv7**...”

Hapless2: “None of them are labeled!”

Hapless1: “Okay... [**cat /proc/cpuinfo**] it’s a Pentium 2.”

Hapless2: “That narrows it down to five boxes...”

This kind of guessing game can continue for quite some time. Following the ground-breaking research of Murphy, if you guess wrong, it is reasonably certain that you will pick a critical server to drop. My least-favourite twist on this is when the boxes have been labeled—but labeled wrong—or labeled with yellow post-it notes (which fall off as the temperature in the server room increases).

If you are using a Linux box, and it has a CD-ROM drive, why not try ejecting it? Using the **eject /dev/cdrom** (or other device name as appropriate) command will make the box spit out its CD tray. It is like telling the real **srv7** to put its hand up.

[Cautionary note: Be careful of doing this to machines where the CD-ROM tray is behind a closed door, such as with the Digital Prioris or the IBM NetVista. Like a tractor-pull for plastic components, you *will* find out whether the server door is stronger than the internal tray mechanism of the CD-ROM drive.]

[Disappointing note: Calling eject on a nonremovable drive does not cause the hard drive to eject its platters. Bummer! A hard drive that could unleash a couple of platters at 10,000 revolutions per minute would be an interesting sight.]

Change Default Passwords (and record them for your successor)

Sometimes in one of these computer ghettos, you will stumble across an unexpectedly nice piece of equipment, such as a managed switch or a decent router. The chances are strong that it will have been left in its default configuration, so that any devious member of staff can *telnet* to it, change its configuration, leaving the network even more fouled up.

Your natural inclination should be to change these passwords—even if people do not act maliciously, they can sometimes foul up equipment accidentally. However, because you have been pressed into service as the network admin, remember that the same fate will likely befall another hapless victim one day. As a mark of consideration, record the equipment description, location, serial number, and new password, on paper. If the company has a safe, store it there. If the company has a safety deposit box, store it there. Make sure someone (a manager or director) knows about it. The time you save may be your own.

Do-It-Yourself Router

Perhaps you have identified that the network really ought to be split up—maybe moving testing to its own segment so that the incessant load-testing does not choke the network for everyone. However, requests for budget allocation to buy a router might not actually be fulfilled. It is at times like this that an old Pentium, two network cards, and a copy of the *Linux Router Project* (LRP) can be pressed into service as a cheap router.

The throughput of such a lo-fi router may not match that of a dedicated unit, but it may suffice for a small organization.

For bonus points, you might also consider setting up some firewall rules on the router, so that the next virus-ridden e-mail opened by someone in marketing does not flood the entire network with excess traffic.

Nagios

Network monitoring tools can make a world of difference to your quality of life as a temporary network administrator. Rather than waiting for users to alert you to a downed Internet connection, you can detect and repair problems as they occur. The ability to maintain logs of link downtime can also help support arguments to replace unreliable links.

Nagios^[1] is a free network monitoring tool. It provides services such as:

- Monitor if a host is up
- Monitor if key services on a host are up
- Monitor if a host is running services it should not

A Web interface allows easy access to status reports. It can be configured to notify you when problems occur, for example, with an e-mail message. Of course, if the mail server is down, this notification method might not be so useful. Such a situation might be better handled by using the Nagios *Short Message Service* (SMS) messaging component.

Given that you might not have a dedicated *Global System for Mobile Communications* (GSM) modem available for sending these SMS notifications, you might like to investigate the Gnokii project^[2]. Ostensibly a project to assist the user in communicating with a mobile phone handset (over data-link cable or infrared), with a capable handset users can initiate sending SMS messages from their handset with Gnokii.

Snort

Intrusion detection might seem a luxury on a network that is struggling to stay operational, but when the price is right (free) and you can spare time to set it up, *Snort* offers a range of features that is surprisingly good. *Snort* can even run without an IP address, making its host computer a fairly difficult target for intruders. The documentation at the *Snort Website*^[3] is quite comprehensive, and I recommend it.

Squid

Squid^[4] is a popular, free HTTP and FTP proxy server. The simple act of caching banner and button graphics for frequently accessed sites can give an apparent increase in Internet bandwidth. The impression for the end user is that things just get faster, because all those pretty graphics load immediately. You may know it is just a nifty trick, but why let on?

Nmap

One characteristic of chaotic networks is that, like weeds after heavy rain, network services spring up everywhere. Programmers are prime offenders in this respect. But be wary—a service with a security flaw, running on an exposed server, can provide an easy beachhead for crackers (a lesson I learned the hard way).

Nmap^[6] is a free network scanner that can assist in finding servers that seem to be running more services than they ought to. It operates in several modes, and offers a range of switches to control its operation.

One of the features that seems more oriented toward people who are scanning networks they are not supposed to is the “Timing policy,” specified with the **-T** command-line switch. The options offered here are *Paranoid*, *Sneaky*, *Polite*, *Normal*, *Aggressive*, and *Insane*. This feature actually comes in handy if the target of your attentions is heavily laden, or lives at the end of a slow link. If you are in the process of tuning a firewall to detect port scans, *Nmap* offers an excellent test facility too.

Another feature that will likely be helpful is the *Nmap* OS fingerprinting facility. Using a combination of techniques^[5], it produces remarkably accurate results for most scans. Combine this result with a port scan and you can build a great picture of which machine has grabbed the wrong IP address (a favorite trick of laptop users: “I didn’t know what my IP address was supposed to be, so I picked one.”) You also can form a rough network map by OS-fingerprinting every active host on your network.

Immunization

It is a good idea to stay up-to-date on your tetanus shots because occasionally you will nick your hands on the sharp bits of metal found in computer equipment.

Traceroute

When licenses for your VisualRouteAnalyser2000 and TrafficGraphic tools have expired, remember that *traceroute* can be one of the most valuable tools to ascertain exactly where things are going wrong. The only (obvious) word of caution is to be aware that overzealous firewall rules can produce spurious results from *traceroute*.

Tag Cables

The desirability of labeling cables is so obvious that it seems silly to even mention it, but it might not have been standard practice for the sysadmin before you. All the more reason you should do the right thing. Sure, *you* know that the purple cable is the link from **gw-eng** to **gw-test**, but will the next person who has to diagnose network issues?

The other impediment to labeling cables is that the sheer volume of unmarked cables makes the task seem futile. Why bother labeling the new one you have just put in, when there are another 40 unknowns? Take heart—by gradually labeling a few here and there, the cables will gradually get less scary each time. Sometimes it can seem like the labor of Sisyphus, but every little bit helps.

Label Equipment

Post-it notes do not constitute an adequate label for network equipment or servers. You are strongly urged to preserve the sanity of other sysadmins by clearly labeling all equipment, using adhesive labels (in a pinch, the labels for a floppy disk will do).

At a minimum I would suggest that host name and operating system (where appropriate), IP address, and a dire warning against tampering with the unit be included. Bonus points are awarded to people who also maintain an equipment audit and record the details of the unit, plus a list of known services that it is running. Of course these will quickly become outdated, but with a known starting point confusion may be reduced.

Destroy Faulty Cables

After several hours of cable tracing, network-card replacement, checking switch link lights, and so on, it may be that you identify a network problem as being caused by a faulty network cable. It can happen anywhere, and is not necessarily a reflection on the skills of the [acting] sysadmin. (Although if the network cable has clearly been mangled and you should have spotted it with a quick visual inspection, you will probably feel a little silly if the time to locate the fault exceeded two hours).

So you whip a replacement cable out of your secret stash (you should have a secret stash of known-good cables) and voila! Network outage fixed. Now comes the most important duty of all—do not discard the damaged cable anywhere that subsequent admins might find it. On several occasions, damaged cables have been put back in operation, only to cause a repeat of the problem that caused them to be removed from service in the first place. It is not uncommon in server rooms to have an empty box that serves as a rubbish bin, but those unfortunates who come after you may not recognize its role as a waste repository in a time of crisis.

If waste is so abhorred that discarding cables is frowned upon, perhaps you can redo the ends of the cable and vigorously retest. Some even maintain that a long cable run can be split into several shorter runs and reused, because the cable fault is likely to be caused by a single break. I disagree—any cable that has broken in one place is likely to suffer further breaks. Demonstrating this principle to overly frugal managers is sometimes best achieved by ensuring the outcome of the demonstration. I suggest laying the cable through a close-fitting door frame and slamming the door on it a few times prior to testing.

Help Dying Equipment on Its Way

Sometimes it can be difficult to discard equipment. Combine this with the almost pathological frugality common in the small business owner, and you find the most decrepit network gear being nursed along. “I just know this old hub has another few years in it. Sure, a few of the Ethernet ports are stuffed, it overheats on warm days, and looks like it might have a mouse nest in the power supply, but that is no reason to discard it.” Nothing is going to convince the owner of this piece of gear that it is time to “redeploy” it in the rubbish bin.

Sometimes you have to be cruel to be kind. Without wanting to seem too much like the *Bastard Operator from Hell* (BOFH)^[7], you may have to help some of this equipment meet its end. It is difficult to identify any one method that fulfills this requirement. My best suggestion is to avoid solutions that leave any externally visible marks (unless they are carbonization marks caused by electrical fault).

You may find that some equipment shows a perverse ability to survive conditions well outside their “recommended operating environment,” and nothing short of a sledgehammer will cause those last two operational ports to die. My recommendation here is to do some network reorganization so that the people responsible for the retention of the equipment are directly affected by it. Nothing says “replace me” quite like frequent trips to the server room to toggle the power switch on an ailing hub. It is surprising how fast requisition orders get signed when managers can no longer browse their favorite Websites.

Conclusion

The crisis has passed. Your time as a sysadmin has passed, and you are free to return to your real job. You have acquitted yourself admirably as sysadmin, and you have learned something in the process.

Like the end of a horror movie, you know that it does not really end here. Somewhere, something is waiting to go wrong. Will you be ready the next time?

References

- [1] Nagios: <http://www.nagios.org/>
- [2] Gnokii project: <http://gnokii.org/>
- [3] Snort: <http://www.snort.org/>
- [4] Squid: <http://www.squid-cache.org/>
- [5] <http://www.insecure.org/nmap/nmap-fingerprinting-article.html>
- [6] Nmap: <http://www.insecure.org/nmap/>
- [7] BOFH: <http://bofh.ntk.net/Bastard.html>

LOCUM SYADMIN is the nom de guerre of a roving programmer who often seems to find himself in sysdamin roles. Operating in deep secrecy, this elusive creature may sometimes be seen tracing cables and cursing. E-mail: locum_sysad@yahoo.com

Letters to the Editor

Ole,

I just finished reading the article about Secure BGP [*Border Gateway Protocol*] by Stephen T Kent. It was very informative and educational with regard to the application and overhead of using the additional BGP attributes and IPSec [*IP Security*]. However, it should be noted that the reliance of a PKI [*Public Key Infrastructure*]-based system, although strong, may also present another possible exploit. If the PKI KDS (*Key Distribution System*) is attacked and subsequently knocked out, including redundant *Key Distribution Engine* (KDE) servers, this may cause serious ramifications to the operation of *Secure BGP* [S-BGP].

Here is a very informative link regarding S-BGP resources for your readers: <http://www.ir.bbn.com/projects/s-bgp>

Also, did you know that the *North American Operators' Group* (NANOG) in conjunction with Cisco engineers recently conducted a BGP vulnerability test? This test confirms that BGP implemented properly is pretty secure in and of itself, without the need for something like S-BGP. The article, titled "BGP Vulnerability Testing: Separating Fact from FUD," was written by Sean Convery and Matthew Franz, Cisco Systems. The article can provide a contrast to the one submitted by Kent and give the technical community both sides of the BGP security issues. Following is the link:

<http://www.nanog.org/mtg-0306/pdf/franz.pdf>

I thoroughly enjoy IPJ and look forward to each issue. Keep up the great work.

—Jeffrey J. Sicuranza, *Applied Methodologies Inc.*
jsicuran@optonline.net

The author responds:

Ole,

Jeffrey makes a few observations about S-BGP in his letter, and they merit responses.

First, I would hope that the discussion of the security features of S-BGP and their direct derivation from the semantics of BGP was as informative as the discussion of performance aspects of the system. After all, a system with good performance but questionable security is probably a poor candidate to S-BGP routing.

Jeffrey raises the question of whether the reliance of S-BGP on certificates, CRLs [*Certificate Revocation Lists*], and address attestations creates significant vulnerabilities that need to be addressed. This is a fair question, but one which I think we have addressed.

The data that S-BGP stores in repositories is data that changes slowly, and thus the system tolerates unavailability of these repositories fairly well. Note that no router ever accesses these repositories in order to verify a route attestation received in an UPDATE. Instead, each ISP [*Internet Service Provider*] or multihomed subscriber NOC [*Network Operations Center*] accesses the repositories to retrieve this data, process it, and distribute the extracted public keys and authorization data to the routers in its network. We anticipate that this process might occur roughly every 24 hours. Because the information represented by the signed objects in the repositories changes very slowly, this retrieval rate seems appropriate. One would expect that these repositories can be engineered to meet these availability requirements. In the worst case, network operators can choose to keep working with the last set of data that they have successfully retrieved. This works because operators process the data before distributing it to their network, and thus can override expired CRLs, etc. So, I think the answer to Jeffrey's cited concern is that S-BGP is not very vulnerable to attacks against these repositories.

I strongly disagree with the conclusions Jeffrey draws from the BGP vulnerability tests he cites. Numerous incidents of BGP security breaches have been reported over the last few years, so there is no question that BGP, as implemented, deployed, and operated, is insecure. Correct implementation of BGP and improved network operator management practices certainly can reduce BGP vulnerabilities. However, the article in question is hardly a refutation of the wide range of vulnerabilities that exist both in practice and in principle. Much of it focuses on a narrow range of attacks, not broader security concerns.

In addressing broader security concerns, for example, the article argues that proper filtering of routes will mitigate the impact of a compromised router. But we know that such filtering is not feasible for many transit network connections, and route filterers are prone to configuration errors. Reliance on transitive trust (for example, assuming that peers filter routes appropriately) makes BGP intrinsically insecure. Relying on *all* ISP operators to *never* make exploitable errors in configuring their route filters, where such filters can be used, is a fundamentally flawed security approach. S-BGP accounts for the reality that not every ISP will operate its network perfectly, and employs mechanisms to allow other ISPs to detect and reject a wide class of errors (or attacks) that may result from such imperfect operation. Thus I reassert that the security vulnerability characterizations that appear in the S-BGP publications are accurate, not overblown.

As a side note, I find it odd that some critics of S-BGP argue that it fails to account for operational reality, yet they offer alternatives that are based on unrealistic assumptions about network operators acting perfectly!

—Steve Kent, BBN Technologies
kent@bbn.com

Book Review

IP for 3G *IP for 3G, Networking Technologies for Mobile Communications*, by David Wisely, Philip Eardley, and Louise Burness, ISBN 0-471-48697-3, John Wiley & Sons, 2002.

I was looking for a book covering mobile communication issues from an IP perspective and IP issues from a mobile communications perspective in order to better clarify details of IP and *third-generation* (3G) convergence. The issue is becoming more and more concrete with the early implementations of 3G networks, so this is a timely book for networking professionals.

Organization

This well-organized textbook helps readers easily understand the “IP-for-3G” issues. It gives a clear vision of that convergence as well as the current snapshot of the recent developments about the subject within the research community. The book is more than an introductory textbook; but readers interested in more technical elaboration can refer to a detailed list of references and further readings given at the end of each chapter.

The book begins with a short chapter that explains the case for IP for 3G. The authors discuss in detail what the term means. They give possible interpretations of IP (Internet, IP Protocol, applications) and their consequent implications on the meaning of IP for 3G. Then they elaborate the IP case within first the “Engineering Reasons for IP for 3G” and then “Economic reasons for IP for 3G” sections.

The second chapter is an introduction to 3G networks. The chapter mostly concerns the core and the access part of 3G networks, skipping the air interface part, because core and access are where IP would make a real difference to the performance and architecture of a 3G network. The chapter reviews briefly the history of 3G developments, from conception to implementation. Then the architecture of *Universal Mobile Telecommunications Service* (UMTS) is introduced, followed by the section where elements of the core network and the architecture of the radio access part are examined. For each part, main functional components such as *Quality of Service* (QoS), mobility management, security, transport, and network management are discussed in detail.

The third chapter discusses the basics of IP and IP networks. Authors give excellent remarks about IP design principles, which are then compared to those of classical telecommunications. Subsequent short sections inform readers about IP addressing schemes, routing, layer behavior, etc. The final section covers the issue of application layer security, which is irrelevant to me for the content of this book. A note: Some of the following chapters require better IP know-how, especially about domain segmentation and intra- and interdomain routing issues. Readers with no prior information are encouraged to refer to other materials before examining the details of, for example, mobility management and QoS.

The fourth chapter is about the multimedia support and session management. First, the concept of session management is introduced. The chapter focuses mainly on the control plane functions of the session management, and the data plane functions are covered in detail in the sixth chapter. The concept of the *Virtual Home Environment* (VHE) is introduced, which forms one of the major requirements of the next-generation mobile system. The authors then review control plane session management protocols, namely H.323 and the *Session Initiation Protocol* (SIP). More discussion is given to SIP, because it is included in the next generation of UMTS standards as the major session management protocol.

The fifth chapter reviews a major problem of the IP-for-3G concept: mobility management. Other key issues of IP such as QoS, IPv6, and session management have always been subject to preceding studies, because those protocols have already been proposed for use in stationary networks. However, the issue of mobility management is a major subject to be investigated for any proper convergence scenario. Personally, I find that this is the biggest challenge of the “long-time-discussed” convergence of IP and mobile communications, and hard work is still ongoing in order to properly resolve the mobility problem. The chapter reviews the basics of mobility such as personal or terminal mobility. From there, macromobility (interdomain or global mobility) and micromobility (intradomain or local mobility) concepts are discussed, followed by proposed protocols for each type of mobility. Mobile IP is examined as the (unique) macromobility protocol. More attention is given to micromobility because it is the most sensitive part of the mobility, under the assumption that 3G BTSs (B nodes) will be simple routers with some extra capabilities. Two variants are discussed, mobile IP schemes, which are based on dynamic tunneling mechanisms, and “per-host forwarding” schemes based on dynamic routing functions. A comparison of major proposals for micromobility management protocols follows.

The sixth chapter considers current IP QoS mechanisms, their operation and capabilities. Those mechanisms created mostly for stationary IP networks may provide a bounded QoS for some “non-real-time” applications, but they are not enough to support any QoS request within the wireless or mobile environment. After giving details of current QoS mechanisms and discussing wireless implications for TCP QoS as well as mobility and wireless issues for *Real Time Protocol* (RTP) QoS, the chapter examines the key elements of QoS and generic features that any prospective QoS mechanism must have. Finally, the authors analyze recent Internet QoS mechanisms such as *Integrated Services* (IntServ), *Differentiated Services* (DiffServ), *Multiprotocol Label Switching* (MPLS), and *Resource Reservation Protocol* (RSVP). The closing section proposes a possible outline solution for how to provide IP QoS for 3G, based on previous work done during the EU BRAIN project.

In the final chapter, the authors summarize all previously given subjects to sketch out the vision of an “All-IP” mobile network. Principles, architecture, routing and mobility issues, QoS, security issues, and interfaces are all discussed to elaborate the generic vision of All-IP networks. Finally, 3G network evolution covering UMTS R4 and R5, and what is beyond 3G, are all discussed.

The book is perfect in the sense that it touches a very hot topic, most of the technical details of which are still in the process of evolving. The authors manage very well the level of details about each subject; they first discuss the overall material before examining details, so readers can obtain a generic but complete view before studying technical details. Each chapter is followed by a comprehensive list of references and further readings, each of them classified by topic. The only fault I find in the book is that SIP should be discussed in more detail.

Recommended

Overall, I would highly recommend this book to any network professional, especially one who is part of any IP-3G convergence process for mobile operators. Still, data network professionals can glean much from the book, because the aim is to carry—a little differently—the same old data, whether or not it contains multimedia, voice, or standard data information.

—Dr. K. Murat Eksiöglu, RT.NET, Turkey
`murat.eksioglu@o2.net.tr`

[Ed.: A version of this review was previously published in the October 2003 issue of *IEEE Communications Magazine* (Vol. 41, No. 10). Used with permission.]

Tim Berners-Lee Knighted by Queen Elizabeth

31 December 2003 — Tim Berners-Lee, the inventor of the World Wide Web and director of the *World Wide Web Consortium* (W3C), will be made a *Knight Commander, Order of the British Empire* (KBE) by Queen Elizabeth. This was announced earlier today by Buckingham Palace as part of the 2004 New Year's Honours list.

The rank of Knight Commander is the second most senior rank of the Order of the British Empire, one of the Orders of Chivalry awarded. Berners-Lee, 48, a British citizen who lives in the United States, is being knighted in recognition of his "services to the global development of the Internet" through the invention of the World Wide Web.

"This is an honor which applies to the whole Web development community, and to the inventors and developers of the Internet, whose work made the Web possible," stated Berners-Lee. "I accept this as an endorsement of the spirit of the Web; of building it in a decentralized way; of making best efforts to keep it open and fair; and of ensuring its fundamental technologies are available to all for broad use and innovation, and without having to pay licensing fees."

"By recognizing the Web in such a significant way, it also makes clear the responsibility its creators and users share," he continued. "Information technology changes the world, and as a result, its practitioners cannot be disconnected from its technical and societal impacts. Rather, we share a responsibility to make this work for the common good, and to take into account the diverse populations it serves." For more information see:

http://www.w3c.org/2003/12/timbl_knighted

SECSAC Publishes DNS Report

The *Security and Stability Advisory Committee* (SECSAC) has published a report entitled "DNS Infrastructure Recommendation." For details see:

<http://www.icann.org/committees/security/dns-recommendation-01nov03.htm>

Coordination, not Governance says ISOC re WSIS

The *Internet Society* (ISOC) published the following text at the *World Summit on the Information Society* (WSIS 2003) which was held in Geneva in early December, 2003:

ISOC is a global not-for-profit membership organisation founded in 1991 to provide leadership in Internet-related standards, education, and policy issues. We are dedicated to ensuring the open development, evolution and use of the Internet for the benefit of people throughout the world. Our education initiatives, for example, have helped bring Internet connectivity to virtually all developing countries over the last 12 years.

ISOC is the organisational home of the *Internet Engineering Task Force* (IETF)—an open consensus-based group responsible for defining Internet protocols and standards. Through our participation in WSIS 2003 we aim to increase understanding and awareness of what is important in order to develop and maintain the Internet’s stability, open nature and global reach.

The Internet has come of Age

In many countries, the Internet has become a mass medium. This has brought with it reflexive pressure on policy makers to regulate it as if it were radio, television, or other mass media. While Governments naturally seek to address their citizens’ interests regarding online privacy, spam, Internet security, intellectual property protection, the price of Internet access, and the digital divide, our position is that better use of technology, and broad participation in today’s Internet coordination processes, not Government regulation, are the most effective and appropriate ways to satisfy these concerns.

The biggest barrier to the Internet fulfilling its immense potential could turn out to be misinformed and inappropriate intervention in the way in which the Internet’s technologies, resources and policies are developed, deployed and coordinated. The Internet Society can help provide guidance here.

What is the nature of the Internet?

The Internet is a modern distributed communications medium. No one is in charge of the Internet and yet everyone is in charge. Unlike the antiquated system of national telephone network monopolies, the global Internet consists of tens of thousands of interconnected networks run by Internet Service Providers, individual companies, universities, Governments, and other institutions. Some of these are global in scope, others regional or local. Hundreds of different organisations and thousands of different companies make decisions every year that contribute to how the Internet develops.

These varied entities, together with the users of the Internet and the developers of Internet technologies and applications, have specific needs for coordination. Collaborative processes that are critical for the future stability and evolution of the Internet, and which should not be modified arbitrarily or abruptly, satisfy these needs.

Coordination, not Governance

It is misleading to use the term “Internet Governance” when the Internet is clearly not a single entity to govern. It is more useful to refer to “Internet Coordination.” The multiple facets of the Internet require different types of coordination, each calling for specific competencies and sensitivities to balance the needs of the Internet user community globally and locally. Specific Internet Coordination activities are taking place globally at three levels:

- Coordination of the definition of Internet standards
- Coordination of the availability and assignment of Internet resources
- Coordination of the policies preventing misuse of the Internet

This coordination is best performed by the existing set of organisations using proven processes. Because of the diverse nature of these activities, it is unrealistic to expect a single body— Government or otherwise—to take on all these roles effectively.

Coordinating Internet standards

The IETF under the umbrella of the Internet Society, is one of the oldest and most successful Internet coordination processes. Other organisations are also involved in Internet-related standards, including the IEEE, the W3C and the ITU.

Many of the protocols at the heart of today's Internet (for example, TCP, IP, HTTP, FTP, SMTP, Telnet, PPP, POP3, the DNS protocol etc.) were developed through IETF standards activities. The results of the IETF are well engineered and practical open protocol standards that are trusted and open to global implementation with little or no licensing restrictions—they are freely available on the Internet, without cost, to everyone.

The strength of the IETF process lies in its unique culture and talented global community of network designers, network operators, service providers, equipment vendors, and researchers. They all openly contribute their individual technical experience and engineering wisdom in an environment that fosters innovation and the open exchange of ideas. This process, which is open to anyone, helps quickly identify and articulate problems of common interest. It also helps build the trust required to make the further investments necessary for a protocol to be usefully implemented and deployed. Ultimately, however, it is the Internet users themselves that determine whether or not a protocol is valuable and useful enough for widespread use. Here the IETF track record of producing useful, widely deployed protocols is unrivaled.

Coordinating Internet resources: The Internet Registry System

There has always been a need to manage the allocation of Internet resources such as the unique addresses that identify devices connected to the Internet (IP addresses), generic top-level domain names (for example, **.org**), country code top-level domain names (for example, **.ch**), domain names (such as **www.isoc.org**), and the systems that translate domain names into IP addresses (for example, the *Domain Name System* or DNS).

This coordination activity has been handled by long-standing, not-for-profit membership organisations such as the *Regional Internet Registries* (RIRs) and *top-level domain* (TLD) registries.

More recently, coordination at a global level has been supported by the *Internet Corporation for Assigned Names and Numbers* (ICANN). Established in 1998, ICANN is also a not-for-profit organisation. Business, technical, non-commercial, academic, governmental and end-user communities participate in ICANN.

These organisations are a meeting point for bottom-up, consensual, industrial self-regulation by the groups and individuals that use their services and resources.

Coordinating policies preventing misuse of the Internet

As we have seen, organisations such as the RIRs, TLD registries, ICANN and the IETF all have very specific roles. It is neither within their charters, nor within their capabilities, to take on responsibility for all areas of Internet Coordination—particularly that of preventing inappropriate use of the Internet. For example, areas such as “cyber crime” (for example, fraud and child pornography) require coordinated global attention by lawmakers—and not by those responsible for the equitable coordination of the underlying Internet infrastructure. Security matters also need to be addressed by organisations providing Internet access (not only by standards developers), and intellectual property issues may best be handled by organisations such as the *World Intellectual Property Organization* (WIPO).

In discussions about these broader Internet policy issues there is cooperation between all the organisations mentioned above. ICANN for example works with WIPO to implement its *Uniform Domain Name Dispute Resolution Policy* (UDRP). And the Internet Society, with technical advice from the IETF, works with Governments and policy makers to explain the effects and possibilities of new Internet technologies.

The way forward: Make your voice heard

Existing consensus-based processes have given us the Internet and have successfully coordinated its phenomenal growth: thousands of new networks, new policy procedures, new top-level domain names, new protocols etc. All of them constantly balance the needs and stability of today’s Internet with future demands.

An open debate is now needed to move towards common, globally acceptable policies, processes and technologies to prevent misuse of the Internet. Governments have a vital role to play here as a concerted effort on the part of the Internet community, non-governmental organisations and Governments can help strengthen and extend today’s successful coordination processes.

The successful continued development of the Internet for the benefit of everyone can be ensured by participation in these proven processes rather than by attempting to create new untested mechanisms that are inappropriate to the unique characteristics of the Internet.

The Internet Society remains dedicated to providing information and orientation about Internet structures and processes. We encourage broad participation in the activities of each of the organisations involved in Internet coordination. For more information on ISOC, visit: **www.isoc.org**

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Technology Strategy
MCI, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.
Copyright © 2003 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PSNRT STD U.S. Postage PAID Cisco Systems, Inc.

The Internet Protocol Journal

March 2004

Volume 7, Number 1

A Quarterly Technical Publication for
Internet and Intranet Professionals

In This Issue

From the Editor 1

High Availability in
Routing..... 2

The Lures of Biometrics..... 15

Book Reviews 35

Fragments 38

FROM THE EDITOR

The operational stability of the global Internet (or any network based on TCP/IP technology) is in large part the result of a carefully configured routing system. Routing continues to be one of the most complex topics in Internet engineering. In our first article, Russ White describes some mechanisms for the design of large-scale, stable routing systems. The article is entitled “High Availability in Routing.”

Security continues to be a high-priority item in computer networks and in society in general. One aspect of security is the identification system by which an individual is given authorized access to a particular facility, be it physical or virtual. Edgar Danielyan gives us an overview of one key element of identification, namely *biometrics*.

The Internet is “going where no network has gone before.” The *National Aeronautics and Space Administration* (NASA) has been working on the *Interplanetary Internet Project* (<http://www.ipnsig.org/>). We hope to bring you an in-depth article about this project in a future issue. An important demonstration of this system took place recently. To quote from the press release:

“A pioneering demonstration of communications between NASA’s Mars Exploration Rover *Spirit* and the *European Space Agency* (ESA) *Mars Express* orbiter has succeeded. On February 6, 2004, while Mars Express was flying over the area Spirit was examining, the orbiter transferred commands from Earth to the rover and relayed data from the robotic explorer back to Earth. The commands for the rover were transferred from Spirit’s operations team at NASA’s *Jet Propulsion Laboratory* (JPL), in Pasadena, California, to ESA’s European Space Operations Centre in Darmstadt, Germany, where they were translated into commands for Mars Express. The translated commands were transmitted to Mars Express, which used them to successfully command Spirit. Spirit used its ultra-high frequency antenna to transit telemetry information to Mars Express. The orbiter relayed the data back to JPL, via the European Space Operations Centre.”

We often receive requests for back issues of IPJ. Although we cannot provide paper copies, all of our previously published editions are available in both PDF and HTML format from the IPJ Website: www.cisco.com/ipj.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

High Availability in Routing

by Russ White, Cisco Systems

A network is a complex system of interacting pieces, as anyone who has ever worked with a large-scale network “in the wild” can tell you. So, when businesses begin asking for a network that can converge in something under 1 second, especially in a large network, network engineers begin to scratch their heads, and wonder what their counterparts in the business world are thinking about. Just about everyone in the network engineering business knows scale and speed are, generally speaking, contradictory goals. The faster a network converges, the less stable it is likely to be; fast reactions to changes in the network topology tend to create positive feedback loops that result in a network that simply will not converge.

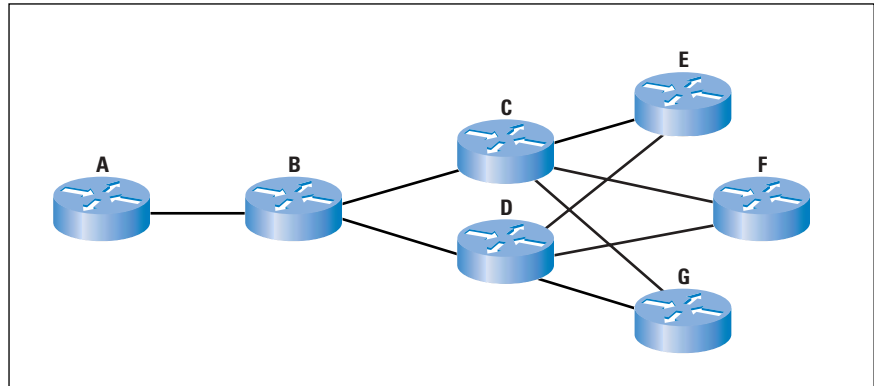
But recent experience has shown that subsecond convergence in a network—even a large network in the wild—is definitely possible. How do we go about building a large-scale network that can converge in times that were, before recently, considered impossible, or improbable, at best? We approach the problem the same way network systems, themselves, are approached. We break the problem down into smaller pieces, and try to solve each piece individually. When we have solved each of the smaller pieces, we recombine them, and see what needs to be adjusted to make it all work together properly.

What pieces of a network do we need to be concerned about when considering subsecond (fast) convergence? Generally, we are concerned with the physical layer (how fast can a down link be detected?), routing protocol convergence (how fast can a routing protocol react to the topology change?), and finally, forwarding (how fast can the forwarding engine on each router in the network adjust to the new paths calculated by the routing protocol?). This article focuses on routing protocols convergence, with some discussion of fast down detection as well, specifically the interior gateway protocols, *Enhanced Interior Gateway Routing Protocol* (EIGRP), *Intermediate System-to-Intermediate System* (IS-IS), and *Open Shortest Path First* (OSPF).

Network Meltdowns

Before beginning to work on a network so it will converge quickly, we need to set some realistic expectations. As mentioned previously, a routing protocol configured to react very quickly to changes in network topology tends to develop positive feedback loops, which result in a network that will not converge at all. Using the following example, consider how a single problem can produce feedback that causes a failure to cascade through the network.

Figure 1: Positive Feedback Loops in a Network



Suppose the link between routers D and G flaps, meaning that it cycles between the down and up states slow enough for a routing adjacency to be formed across the link, or for the new link to be advertised as part of the topology, but too quickly for the link to actually be used. In this situation, the adjacency (or neighbor relationship) between routers D and G forms and tears down as quickly as the routing protocol will allow.

While this is occurring, the routing information at routers E, F, and G is changing as quickly as the adjacency between D and G can form and tear down. This change in routing information is, in turn, passed on to C, which then must process it as fast as it possibly can. It is possible that the routing information presented to router C will overcome the ability of its processor to process the information, causing router C to fail, or drop its neighbor adjacencies.

At the same time, the constantly changing routing information at router B will also cause problems, possibly causing it to periodically drop its adjacencies, specifically with routers C and D. At this point, if the routers B, C, and D are all three consuming a large amount of memory and processing power adjusting to apparent topology changes because of changing adjacency states, the flapping link between routers D and G, which originally caused the problem, can be removed from the network, and the routing protocol will still not converge. This is what network engineers consider a classic *meltdown* in the routing system.

Solving the Meltdown

Typically, when a network engineer faces a network in this condition, the first step is to simply remove routing information from the system until the network “settles.” This typically involves removing parallel (redundant) links from the view that the routing protocol has of the topology until the routing protocol converges. At this point, the network would be examined, routers reloaded as needed, and the parallel links brought back up. The network design might then be reviewed, in an attempt to prevent recurrence of a meltdown.

Routing protocol designers and developers would also like to move the point at which a routing protocol “melts” as far along the curve of network design as possible.

Of course, it is impossible to prevent all network meltdowns through protocol design; there are limits in any system where the implementation steps outside the “state machine,” and the system will simply fail. But how would a routing protocol designer work around this sort of a problem in the protocol itself? The answer is actually very simple: Slow down.

The main problem here, from a protocol designer’s point of view, is that routers D and G are simply reacting too fast to the changing topology. If they were to react more slowly, the network would not fall into this positive feedback loop, and the network would not melt. And, in fact, slowing down is really quite simple. Various methods of slowing down include:

- Not reporting all interface transitions from the physical layer up to the routing protocol. This is called *debouncing* the interface; most interface types wait some number of milliseconds before reporting a change in the interface state.
- Slow neighbor down timers. For instance, the amount of time a router waits without hearing from a given neighbor before declaring that a neighbor has failed is generally on the order of tens of seconds in most routing protocols. The dead timer does not impact down-neighbor detection on point-to-point links, because when the interface fails, the neighbor is assumed to be down, but there are other “slow-down” timers here, as well.
- Slow down the distribution of information about topology changes.
- Slow down the time within which the routing protocol reacts to information about topology changes.

All four of these methods are typically used in routing protocols design and implementation to provide stability within a routing system. For instance:

- In IS-IS, a timer regulates how often an intermediate system (router) may originate new routing information, and how often a router may run the *shortest path first* (SPF) algorithm used to calculate the best paths through the network.
- In OSPF, similar timers regulate the rate at which topology information can be transmitted, and how often the shorter path first algorithm may be run.
- In EIGRP, the simple rule: “no route may be advertised until it is installed in the local routing table” dampens the speed at which routing information is propagated through the network, and routing information is also paced when being transmitted through the network based on the bandwidth between two routers.

It seems like the simplest place to look when trying to decrease the time a routing protocol requires to converge, then, is at these sorts of timers. Reduce the amount of time an interface waits before reporting the transition to a down state, reduce the amount of time a router must wait before advertising topology information, etc. But when we consider implementing such changes, we remove much of the stability we have all come to expect in routing systems—the size a network can be built without melting down decreases below an acceptable threshold, even with modern processors, more memory, and implementation improvements in place.

There is another place to attack this problem: the frequency of changes within the network. This is the same concept—speed—from a different angle. How does looking at it from a different angle help us? By allowing us to see that it is not the speed of the network changes that causes the positive feedback loop, but rather how often the changes take place. If we could report the changes quickly when they occur slowly, and report them more slowly when they occur quickly, or if we could just not report some events at all, routing could converge much faster, and still provide the stability we expect.

The two options we want to examine, then, are not reporting every event, and slowing down as the network speeds up. First we will discuss these two options, and then discuss speeding up the reporting of network events, which plays a large role in decreasing convergence times.

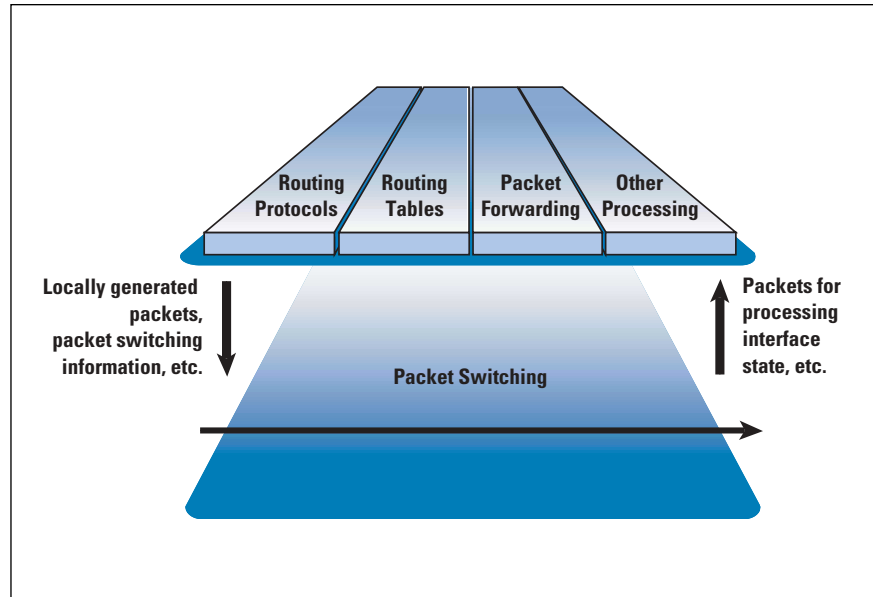
Do Not Report Everything You See (NSF and GR)

It sounds simple just to say that a router should not report every event within the network it is aware of, but it becomes more complicated as we consider the issues involved. What we need to do is sort out which events are important, in some sense, and which are not. For instance, if a router loses contact with an adjacent router because the adjacent router restarted for some reason, do not report the resulting change in topology until you are certain the neighbor is not coming back.

But the classic questions follow: How long do you wait before deciding the problem is real? And what happens to traffic you would normally forward to that neighbor while you are waiting? Finally, how do you reconnect in a way that allows the network to continue operating correctly? A technology recently incorporated in routing protocols called *Graceful Restart* (GR), combined with another technology called *Non-Stop Forwarding* (NSF), can combine to answer these questions.

Let's start at the bottom of the *Open Systems Interconnection* (OSI) model, at the physical and data link layers, and discuss the second question, what happens to traffic that would normally be forwarded while a router is restarting? Normally, this traffic would be dropped, and any applications impacted would need to retransmit lost data. How could we prevent this? We can take advantage of the separation between the control plane and the forwarding plane in a large number of modern routers.

Figure 2: Control and Data Plane Interaction in a Router



In some routers, such as the Cisco® 12000, 10000, 7600, and others, the actual switching, or forwarding, of packets is performed by different processors and physical circuitry than the control plane processes run on (such as routing protocol processes, routing table calculation, and other processes). Therefore, if the control plane fails or restarts for any reason, the data plane could continue forwarding traffic based on the last known good information.

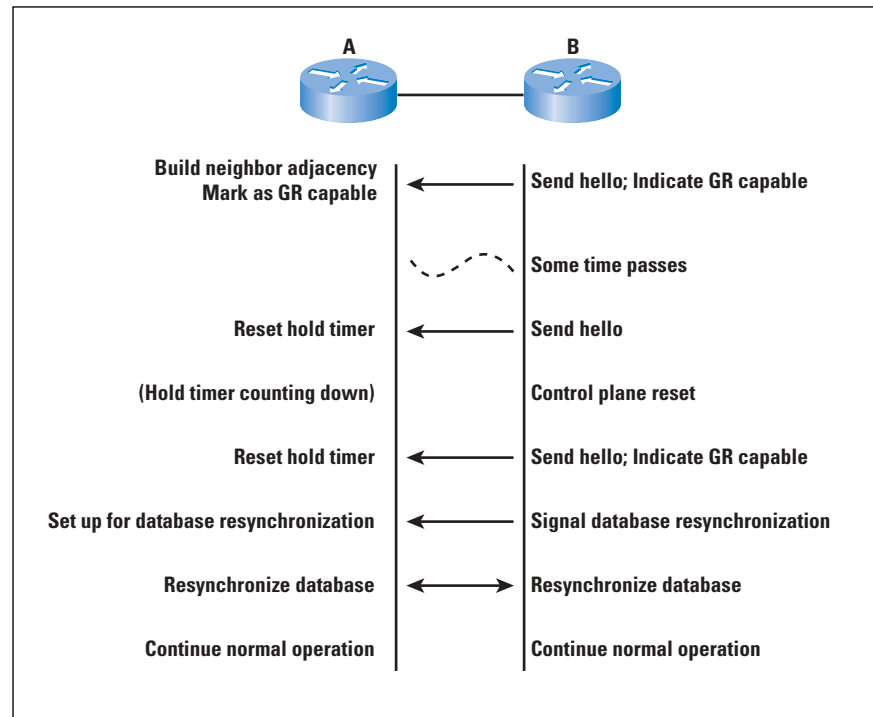
NSF, implemented through *Stateful Switchover* (SSO) and *Stateful Switchover+* (SSO+) in Cisco products, allows this continuous forwarding, regardless of the state of the control plane, to take place. Normally, when the control plane resets, it sends a signal to the data plane that it should clear its tables out, and reset, as well. With NSF enabled, this signal from the control plane simply acts as a signal to mark the current data as stale, and to begin aging the information out.

Now we need to be able to bring the control plane back up, resynchronize the routing protocol databases, and rebuild the routing table, all without disturbing the packets still being switched by the data plane on the router. This is accomplished through GR. GR starts by assuming two critical things:

- The normal hold times are acceptable, within this network environment, for reporting a network event or topology change. In other words, if a router's control plane fails, the event wouldn't be reported until the routing protocol's default hold or dead timers expire, whether or not GR is configured.
- The control plane on the router can reload and begin processing data within the hold or dead time of the routing protocol.

Let's examine how, in principle, GR works, so we can put these two requirements into context, and understand where GR is best deployed in a live network. Consider the following chart to understand how GR works between two peers of any generic routing protocol.

Figure 3: The Process of Graceful Restart



When two routers begin forming an adjacency (or neighbor relationship, or begin peering, depending on which routing protocol is being run between them), they exchange some form of signaling noting that they are capable of understanding GR signaling, and responding to it correctly.

[Note that this does not imply the router is GR-capable, only that it can support a neighboring router performing a GR. For instance, the Cisco 7200 supports switching modes only where the control and data planes are not cleanly separated, so it cannot fully support GR. It can, however, support the signaling necessary for a neighboring router to gracefully restart.]

Assume some time passes, and router B is transmitting Hello packets to router A normally, on a periodic basis. Each time router A receives one of these Hello (or *keepalive*) packets, it resets the hold, or dead, timer on router B, indicating that it should wait that amount of time before declaring router B down if it stops receiving Hellos. Now, at some point, after sending a Hello packet, the router B control plane resets. While the control plane is down, the router A hold timer is still counting down; the routing protocol does not reset the session. This is, in fact, normal routing protocol operation, which normally results in the packets forwarded by router A toward router B to be dropped. Because router B is NSF-capable, however, its data plane is still forwarding this traffic to the correct destination, even though the control plane is down.

If the router B control plane does not come back up within the dead or hold timer allowed by the routing protocol, router A declares the adjacency down, and begins routing around router B. This explains why the router B control plane must come back up within the hold interval of the routing protocol, one of the two assumptions we outlined GR as making at the beginning of this section. For this case, we assume that the router B control plane comes back up before the router A hold timer expires, and router B sends a Hello with no information other than indicating it is restarting.

When router A receives this Hello, it acts as though it has received a normal Hello, and simply keeps its adjacency with router B up. In other words, although router B may not know what the network it is connected to looks like at this point, router A does not report this failure to the rest of the network. Convergence time is, from a network standpoint, effectively reduced to 0.

When the router B control plane completes its reset, it then signals router A to begin resynchronizing their databases. The two routers then use some method specific to each protocol to resynchronize their databases, and begin operating normally, in a stable condition once again.

Slow Down When the Network Speeds Up

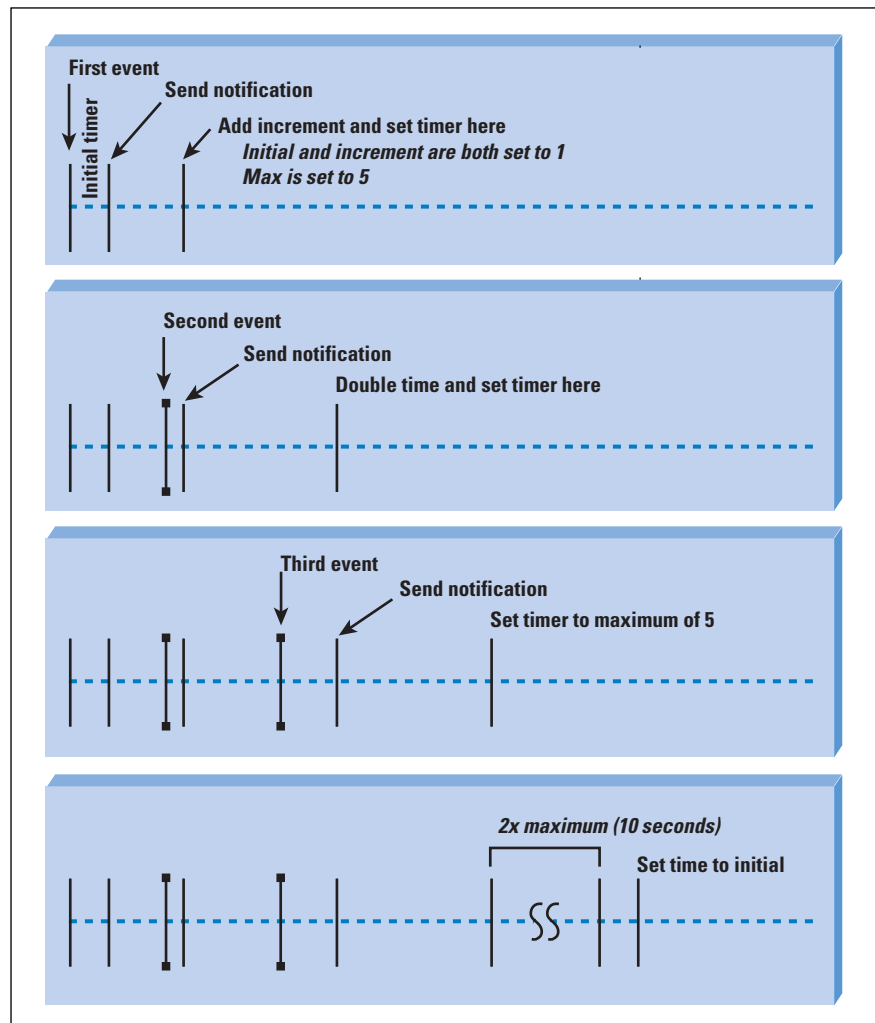
The second option we discussed originally was to attack the problem by reducing the frequency, rather than the number, of updates. What we want to do is to slow down the reporting of events when they occur more frequently (or when they occur rapidly), and speed up the reporting of events when they occur less frequently (or when they occur slowly). This is possible through a series of features built into Cisco IOS® Software within the last year or two, applying the concept of the *exponential timers*.

An exponential timer changes the amount of delay between an event occurring and the reporting of that event by the frequency at which the event occurs—possibly not reporting the event at all, in some situations. Two implementations of exponential timers are *exponential backoffs* and *dampening*. Let's examine each of these individually, and then consider where they are implemented in Cisco IOS Software.

Exponential Backoffs

Consider the following figure to examine how exponential backoff works.

Figure 4: Exponential Backoff



When the first event occurs, a timer is set to the initial time, 1 second in this case, meaning that the router waits for one second before notifying other routers in the network about the event. When the notification is sent, the router adds the initial timer to the increment, and sets a timer for this period. We call this timer the *backoff timer*.

When the second event occurs, the backoff timer is still running; the router waits until this timer expires to send the notification about this event occurring. When this notification is sent, the backoff timer is set to twice the previous setting or the maximum backoff time, whichever one is shorter. In this case, doubling the backoff timer results in 4 seconds, so it is set to 4 seconds.

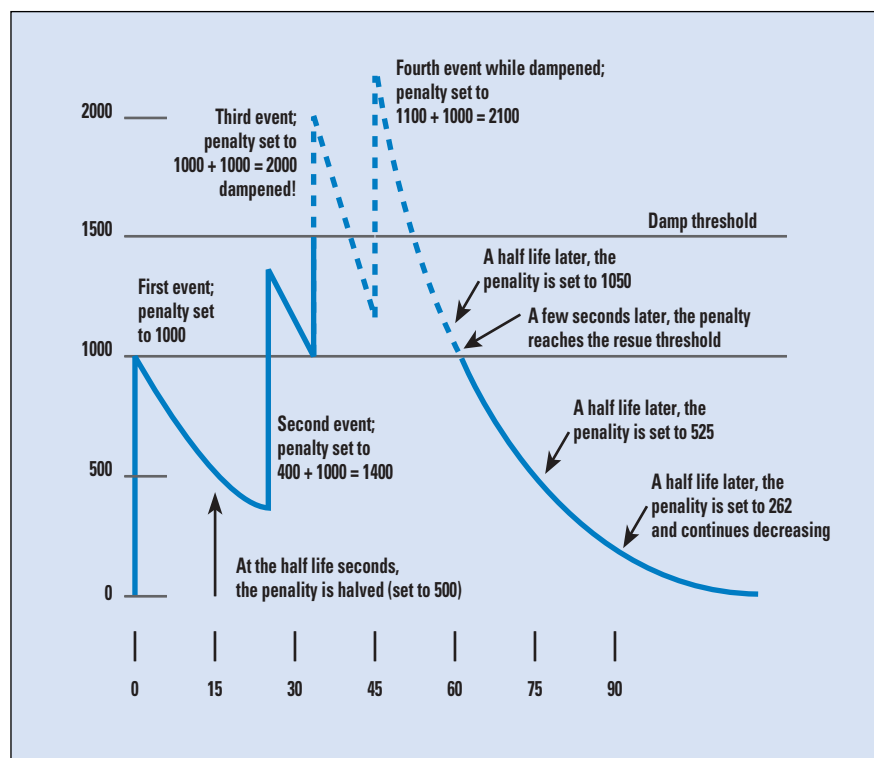
When the third event occurs, the backoff timer is still running; the router waits until the timer expires before sending any notification of the event occurring. Again, the timer is doubled, this time to 8 seconds, and compared to the maximum time, which is 5 seconds. The shorter of the two times is taken, so the backoff timer is now set for 5 seconds.

At this point, any future events will be reported only at 5-second intervals, as long as one event occurs at least every 5 seconds. If no events occur for an interval of 10 seconds, the timers are all reset to their initial condition, so the initial timer is set to 1 second, and the backoff timer is not set at all.

Dampening

Dampening, or damping, is also an exponential backoff mechanism similar to the exponential backoff algorithm we examined previously. The primary difference is that dampening is applied to events that have a Boolean component; a route that is either advertised or withdrawn, an interface that is either up or down, etc. Exponential backoff simply deals with events in general, whereas dampening adds value based on the type of event, as well as the frequency at which the event occurs. Consider the following figure to understand dampening.

Figure 5: Dampening Over Time



In dampening, the desirability of reporting an event is set using the *penalty*; the higher the penalty applied to a given item, such as a route or an interface, the less desirable it is to advertise changes in the state of that item. Dampening always leaves the item in the “off,” or “down,” state, when it stops reporting state changes; this is called the *dampened* state. A penalty is normally added when transitioning from “down” to “up” in most dampening systems.

Here, we start at time 0, with a penalty of 0; when the first event occurs, a penalty of 1000 is added, making the total penalty 1000. As time passes without another event occurring, the penalty is decreased, based on the *half life*. Each time the half life passes, in this case 15 seconds, the current penalty is halved, so after 15 seconds, the half life is set to 500.

A few seconds later, while the penalty is still decreasing, the second event occurs; 1000 is added to the current penalty, making the total penalty 1400. Again, as time passes, the penalty decays exponentially, reaching 1000 before the third event occurs. When the third event occurs, 1000 is again added to the total penalty, so it reaches 2000—which is above the *damp threshold*, so future events are dampened by simply leaving the interface or route in the down state.

Again, as time passes, the penalty is cut in half for each passing half life, reaching 1100 before the fourth event occurs. When the fourth event occurs, 1000 is again added, making the penalty 2100, and leaving us in the dampened state until the penalty can be reduced again. Over time, the penalty finally drops to 1000 (at around 60 seconds in the example), which is the *reuse threshold*. At this point, state changes in the item being tracked are once again reported as they occur, unless the penalty reaches the dampening threshold at some future point.

So, dampening reacts to events by simply not reporting events if they occur too frequently, whereas exponential backoff reacts to events by reporting each event that occurs, but slowing down the reporting of events as they occur more frequently.

Speeding Up the Reporting of Events

When we have some methods in place to prevent a network meltdown when events occur, we can consider ways to discover events faster. Primarily, these techniques are used in conjunction with exponential backoff and dampening.

There are two ways to detect a down neighbor or link: *polling* and *event driven*. We will briefly discuss each of these, and some various techniques available in both cases.

Polling

One method commonly used for detecting a link or adjacency failure is polling, or periodically sending Hello packets to the adjacent device, and expecting a periodic Hello packet in return. The speed at which Hello packets are transmitted and the number of Hello packets missed before declaring a link or adjacency as failed are the two determining factors in the speed at which polling can discover a failed link or device.

Normally, a neighbor or link is declared down if three Hello packets are lost, meaning that the hold time, or the dead time, will always be about three times the Hello time, or polling interval. Normally, for Layer 2 links and routing protocols, the Hello interval is measured in seconds. For instance:

- EIGRP running over a point-to-point link sends one Hello every 5 seconds, and declares a neighbor down if no Hellos are heard for 15 seconds.
- EIGRP running over a lower-speed point-to-multipoint link sends one Hello every 60 seconds, and declares a neighbor down if no Hellos are received in 180 seconds.

- OSPF normally sends a Hello every 10 seconds, and declares a neighbor down if no Hellos are heard for 40 seconds.
- *Frame Relay Local Management Interface* (LMI) messages, the equivalent of a Hello, are transmitted every 10 seconds. If an LMI is not received in 30 seconds, the circuit is assumed to have failed.
- *High-Level Data Link Control* (HDLC) keepalive messages are transmitted every 10 seconds. If a keepalive message is not received within 30 seconds, the circuit is assumed to have failed.

Fast Hellos can decrease these timers to Hello intervals on the order of 300 milliseconds, and dead timers of around 1 second.

The primary problem with fast Hellos is scaling, particularly in receiving and processing fast Hellos from a large number of neighboring routers. For instance, if a router has 1000 neighbors and is using a Hello interval of 330 milliseconds, the router has to be able to receive and process 3000 Hellos per second and send 1000 Hellos per second. Timers in this range leave little room for processes that consume a router processor for long periods of time, short-term packet loss on a link due to congestion, and other factors.

Event Driven

Rather than polling at a fast interval, event-driven notifications rely on devices within the network that can sense the state of a link through lower layers (electrical, electronic, or optical state) to notify the routers attached to the link when the link has failed. SONET links are probably the best example of media with built-in capabilities for sensing link failures and notifying attached devices. This Tech Note on Cisco Online:

http://www.cisco.com/en/US/tech/tk482/tk607/technologies_tech_note09186a0080094522.shtml

... provides information about SONET alarms. There are also techniques that can be used to speed up the reporting of failed links in Frame Relay circuits, and techniques are being developed for allowing switches to notify devices attached to an Ethernet VLAN about a loss of connection to an attached device.

Implementations

Now that we have discussed what exponential backoff and dampening are, we can consider how they are implemented, and how their implementation helps you build highly available networks (through fast convergence) without risking major network instability along the way. We start by examining where dampening is implemented, and then follow that with a discussion about where exponential backoff is implemented. These sections do not provide a great deal of detail on the implementation of these features; vendor documentation and other sources of information (such as the forthcoming book *Designing to Scale*) should be consulted for technical details.

Dampening

Dampening is currently implemented in two places:

- *Border Gateway Protocol* (BGP) route flap dampening
- Interface dampening

BGP route flap dampening is a well-known technology, deployed in the Internet on a wide scale to increase the stability of the Internet routing table.

Interface dampening allows the network engineer to prevent rapidly flapping interfaces from having a wide-ranging impact on the entire network. When an interface fails and comes back up numerous times within a short time period, the interface is placed in the down state from an IP perspective, and not advertised within routing protocols, or used for forwarding packets.

It is important to note that the interface is allowed to change states freely at Layer 2; an interface that continues to change state rapidly continues to accumulate penalties, and continues to show down to the IP subsystem.

Exponential Backoff

Exponential backoff is implemented in several places in link state protocols at this point, including:

- The ability to exponentially back off the amount of time between a change in the network topology being detected and the transmission of a link state packet being transmitted to report the change; exponential backoff has been applied to the link state generation timer.
- The ability to exponentially back off the amount of time between receiving a link state packet reporting a change in the network topology, and running SPF to recalculate the path to each reachable destination in the network; exponential backoff has been applied to the SPF timer.

Fast Hellos

Each routing protocol has a different limit on how Fast Hellos can be transmitted and how often they must be received for a neighbor to be considered alive. OSPF and IS-IS have both implemented the fastest Hellos, with a minimum of 330 millisecond Hellos, and a dead interval of 1 second.

EIGRP can run with Hellos as fast as one per second, with a 3-second dead time. BGP can use similar timers, with a keepalive interval of 1 second.

Caution should be used when configuring Fast Hellos on a network. Congestion, high processor use, and other problems can cause false down indications that may cause higher rates of network failure than would normally occur.

Deploying GR and Fast Convergence Technologies

We now have a full range of options we can use to improve network availability, including GR and NSF, event dampening, and fast convergence techniques. How can we deploy these in a real network to improve network uptime? Generally, the technologies can be placed in one of three categories:

- *Fast reaction to node or link failure, to route around the failure.* We use Layer 2 techniques and Fast Hellos to quickly determine when an adjacent node, or a link to that node, has failed.
- *Slow reaction to node or link failure, combined with routing through the failure.* We rely on moderate speed reactions to node failures to allow resynchronization of routing data while forwarding of traffic continues.
- *Fast recalculation of the best path when a topology change has been reported.*

As we can see, the first two are complementary; we could not deploy both of them in the same place in the network. The third one, fast recalculation, can be deployed with either (or both) fast reaction and slow reaction techniques to increase network availability. The primary question then becomes: which of these two techniques do you deploy in your network, and where?

The basic premise behind making this design decision follows:

- If there is a fast, online backup available to a node or a link, it probably makes more sense to route around any problems that occur as rapidly as possible.
- If any existing backup is going to take a good deal of time to bring online, or there is no backup path (such as a single homed remote office, or a remote office with dial backup), it probably makes more sense to route through any problems.

In general, then, we want to deploy techniques that improve network convergence time everywhere—techniques that bring down the time a network is down when a failure occurs, is detected, and a new path calculated. At the same time, we want to evaluate each point in the network we would like to protect from failure, and determine the best means to protect that point of failure: redundancy with fast down detection, GR, or NSF.

Fast, stable networks are possible with today's techniques in routing; some large networks, with several hundred routers, measure their convergence times in the milliseconds, with 1 second as their outside convergence goal.

RUSS WHITE is on the Cisco Systems Routing DNA Team in Research Triangle Park, North Carolina, specializing in the deployment and architecture of routing protocols. He has coauthored books on routing protocols, network design, and router architecture, regularly speaks at the Cisco Networkers conference, and is active in the Internet Engineering Task Force. Russ can be reached at riw@cisco.com. This article offers a high-level overview of material covered in depth in a forthcoming network design book, *Designing to Scale*, being published through Cisco Press.

The Lures of Biometrics

by Edgar Danielyan, Danielyan Consulting LLP

This article introduces biometrics and discusses some of the complex issues associated with use of biometrics for identification and authentication of individuals and its impact on both standalone and networked information systems, as well as on physical security. The agenda is not to show whether biometrics is your best investment or a useless thing—these two polar viewpoints share the same quality of being oversimplifications, to say the least. It also certainly does not purport or try to tell everything there is to tell about biometrics or its applications. Legal and social implications of biometrics are also not discussed in this article because these would differ considerably, depending on the legislation and cultural traditions of countries concerned; we also do not consider the complex performance, design, and implementation questions, because these are of too specialized nature—for more in-depth coverage of these topics a list of biometrics organizations and publications are provided at the end of this article, along with a list of references.

Before we continue, it would be useful to examine the current deployment of biometrics outside testing laboratories and the corporate perimeter. With the U.S. government fingerprinting and taking photographs of some of the visitors coming to the United States beginning January 5, 2004, under the US-VISIT program, biometrics and associated issues such as privacy and personal data protection are bound to get unprecedented levels of publicity^[1]. Although it is too early to judge whether this innovation will actually contribute to overall security of the country or rather increase the general confusion surrounding security procedures, it has already resulted in more questions asked than answered. To some of its proponents, biometrics is a magic technology that would contribute to the security of their societies, to others the same technology heralds the coming of a police state and erosion of personal privacy and liberties and discrimination against (potentially not only) foreign citizens. Indeed, that was the opinion of Julier Sebastiao da Silva, a federal judge in Mato Grosso state of Brazil, who ordered similar measures to be taken in the case of U.S. citizens visiting Brazil^[2]. Despite the announcement of Brazil's federal police that they may well seek to have this judgment overturned, this is a significant event because it illustrates that the use of biometrics is not only a technical procedure but also has its far-reaching social, legal, and international implications. It is immaterial whether this judgment will be upheld or overruled—it is the fact that introduction of the mandatory use of biometrics at borders resulted in such a response that is important.

Earlier announcement by the U.S. authorities that they expect the visa-waiver countries whose citizens currently may enter the U.S. without visas, simply upon presentation of their passports, to provide biometric data in newly issued passports also resulted in different reactions, ranging from support for the measure to outright condemnation^[3].

Aside from the huge technological and logistical work that must be done in order to introduce biometrics into passports, these requirements also pose considerable legal and social issues in countries with strong personal privacy and data protection legislation in place. However, one thing is clear—biometrics ceases to be an exotic and little-used technology and is bound to be increasingly used in one way or another.

This article is organized as follows. First biometrics and related concepts are introduced, along with descriptions of the most widely used and understood physiological and behavioral biometrics. We will also see how biometric systems fail when inadequately designed or implemented. Later we describe the system and design issues of biometrics, such as security, accuracy, speed, resilience, privacy, and cost of biometric identification and verification systems, as well as practical applications of biometrics in network authentication and international travel documents.

Definition of Biometrics

A *biometric* is a physiological or behavioral characteristic of a human being that can distinguish one person from another and that theoretically can be used for identification or verification of identity. For a biometric to be practically useful, ideally it should be unique, universal, permanent, recordable, and acceptable—more on these properties of practical biometrics later.

Authentication in General

Authentication is the second step in the identify-authenticate-authorize process, which is done countless times every day by humans and computers alike. When speaking about human authentication, basically we have three choices: using something we know (such as passwords and passphrases), something we have (such as access tokens, smart cards, and so on) or something we are (biometrics). There is no “best” authentication method; each has its pros and cons, depending on the application, the users, and the environment. Whatever authentication method we use, we can make it stronger by using one or both of the other methods. An example of strong authentication would be a system that requires possession of a smart card, knowledge of a password or *Personal Identification Number* (PIN), and biometric verification. Obviously to steal or fake all three would be much more difficult than to steal or fake any one of these—however, more expensive and laborious to operate as well. The other two factors—the time of access and the location of subject—may also be used for access control, but usually only as auxiliary factors.

What You Know

Unquestionably the most widely used method of authentication, passwords, passphrases, and PINs share both pros and cons with each other. Moreover, an advantage in one situation easily becomes a problem in another—an example being the ease of password sharing. Passwords are easy to change, but are also easy to intercept. Systems can force the use of strong passwords, but the user may respond by storing or transmitting them in such a way that the added security is effectively reduced to nil.

Unauthorized disclosure of a password is not usually detected until after unauthorized access has already taken place. Passwords are also vulnerable to guessing, dictionary, and brute-force attacks. On the other hand, they require no additional hardware, they are an accepted method of authentication, and they are well-understood—even by the most technologically challenged part of human species.

What You Have

Smart cards, access tokens (both challenge-response and time-based), and other “what you have” authentication methods solve some of the problems associated with “what you know” authentication, but they create a set of different problems. Unlike theft of a password, theft of a smart card or access token can, of course, be easily detected. Unlike passwords, smart cards usually cannot be used simultaneously by two or more parties in different places. However, “what you have” authentication devices may be lost, damaged, and stolen. They may also run out of power (if self-powered) or may be prone to power-, synchronization- and time-based attacks if externally powered. They may also be subjected to reverse engineering and other treatment, which may compromise their security.

What You Are: Biometric Authentication

There are two biometric authentication methods: biometric verification and biometric identification of identity. Biometric identification is also sometimes referred to as *pure biometrics* because it is based only on biometric data and is more difficult to design and operate—but alas, pure biometrics is not the most secure, useful, or efficient one. Also, both methods can not always be used with all biometrics—some biometrics can only be used in verification mode because of their intrinsic properties.

Verification

Biometric verification uses entity IDs and a biometric—in this case biometric merely serves to prove identity already declared by the entity—which may be done using something you know (a username) or something you have (a smart card). Biometric (something you are) works to actually complete the authentication process. Hence, the biometric database keeps a list of valid entity IDs (which may be said to serve as primary keys to the database) and corresponding biometric templates, and compares (“matches”) the stored template with the biometric provided. The result of this comparison is either an accept or reject decision based on a complex algorithm and system settings (refer to the section “Matching”).

Identification

Unlike biometric verification of identity, biometric *identification* is based solely on biometrics. The biometric serves as both the identifier and the authenticator. The biometric database contains the enrolled biometric templates, and they all are compared against the provided biometric to find a match. Biometric identification may be described as “putting all your eggs in one basket,” partly because somehow faking or stealing a biometric compromises both the ID and the authenticator.

A biometric identification system may operate in one of the two modes: positive identification or negative identification. In a positive identification biometric system, the provided biometric must be in the database and there must be only one match to positively identify the person. The risks present in a biometric system are false acceptance and false rejection, whereas unauthorized subjects are incorrectly accepted, or authorized ones are denied identification, resulting in a denial of service. A negative identification system, in contrast, works by determining whether the provided biometric is not in the database.

Enrollment

Regardless of the type of a biometric system, *enrollment* is a mandatory part of the process. Biometric enrollment is the registration of subjects’ biometrics in a biometric database. Positive enrollment results in a database of recognized persons’ biometric templates that may be later used for positive identification or verification. Negative enrollment results in a database of “excluded” persons, a black list if you wish. Security and reliability of the enrollment process and the biometric database are fundamental to the security of the entire system, but in practice they are difficult to achieve because of the myriad of issues that affect collection, transmission, storage, and usage of biometric data (see “Security” and “Privacy,” later in this article for an overview of just some of the risks).

Matching

After an individual is enrolled—that is, the individual’s biometrics are scanned and registered in the biometric database—*matching* is the next step. Biometric matching is essentially the comparison of the enrolled person’s known biometric data stored in the biometric database in the form of biometric templates—binary representation of biometric sample—with the biometric provided by the individual at the identification or verification time. However, biometric matching is a pattern-recognition problem and not a simple bit-by-bit comparison—representation of the same biometric taken by two input sensors or taken at two different points in time does not match bit by bit because of numerous factors such as sensor resolution, system noise, and so on. Therefore, a degree of likeness (usually referred to as the *matching score*) is used to express how like the stored biometric is to the provided biometric. A *threshold level* is used to decide whether the matching score is high enough to be considered a match—if the score is at or below the threshold level, matching fails. This threshold level is one of the many variables that affect the accuracy—and hence security—of biometric authentication systems.

For biometric identification applications, the provided biometric is compared against all entries in the database and should result in only one successful match to result in positive identification. In biometric verification systems, the provided biometric is compared only with the biometric template or templates corresponding to the specified identity. As a result of biometric matching, the following system errors may occur:

- *False match or acceptance*: This occurs when the system decides that the two biometrics (the one stored in the database and the one provided now) are the same, when in reality they are not. The rate of false matches is known as *False Matching Rate* (FMR) or *False Acceptance Rate* (FAR). False acceptance is a confidentiality and integrity risk.
- *False nonmatch or rejection*: This is expressed as *False Rejection Rate* (FRR), and *False Nonmatching Rate* (FNMR). False nonmatch is when the system erroneously decides that biometrics are from different identities while in reality they are from the same person. False rejection is an availability risk.

In practice, both FRR and FAR do not equal zero, and in different applications one of them may be more important than the other. In an application that requires higher security (and hence as low FAR as possible), users may be troubled with high false rejection rates; whereas in an application that can accept somewhat higher false acceptance rates (such as public transport), false rejection rate is of more concern because of convenience and manual processing concerns. When FAR and FRR meet, that is the *Cross-over Error Rate* (CER). The lower the CER, the better—hence it is frequently used to express accuracy of biometric systems (although it is not the infallible measure as some suppose). Additionally, *Failure to Acquire* (FTA) errors occur when an individual does not have the required biometric or the biometric cannot be read by the sensor; and *Failure to Enroll* (FTE) is when a part of the targeted population may not be enrolled for whatever reason (such as a FTA). These errors directly affect the practicality of biometrics and must be accounted for with regard to the projected population of users.

Practicality of Biometrics

Writing in the December 1994 issue of *Information Technology & People* (“Human identification in Information Systems: Management Challenges and Public Policy Issues”)^[4] ten years ago, Roger Clarke proposed some criteria that should be met in order for a biometric to be practically usable:

- *Universality*: Every relevant person should have an identifier.
- *Uniqueness*: Each relevant person should have only one identifier, and no two people should have the same identifier.
- *Permanence*: The identifier should not change, nor should it be changeable.

- *Indispensability*: The identifier should be one or more natural characteristics, which each person has and retains.
- *Collectibility*: The identifier should be collectible by anyone on any occasion.
- *Storability*: The identifier should be storable in manual and in automated systems.
- *Exclusivity*: No other form of identification should be necessary or used.
- *Precision*: Every identifier should be sufficiently different from every other identifier that mistakes are unlikely.
- *Simplicity*: Recording and transmission should be easy and not error-prone.
- *Cost*: Measuring and storing the identifier should not be unduly costly.
- *Convenience*: Measuring and storing the identifier should not be unduly inconvenient or time-consuming.
- *Acceptability*: Its use should conform to contemporary social standards.

Although some of these criteria may be argued over, this set is nevertheless a useful reference. An interesting point is that no known biometric completely satisfies all of these criteria, perhaps proving that these are not strict “must haves” but instead guidelines to be accounted for.

Types of Biometrics

Two broad categories of biometrics exist: *physiological* biometrics (such as fingerprints, hand geometry, iris recognition) and *behavioral* biometrics (such as signature and voice biometrics). Physiological biometrics is based on direct measurements and data derived from measurements of a part of the human body, whereas behavioral biometrics is based on measurements and data derived from human actions, and indirectly measures characteristics of the human body over a period of time.

Physiological Biometrics

Relatively widely understood and used physiological biometrics are fingerprint recognition, face recognition, hand geometry, and iris recognition. These methods are introduced in the following sections.

Fingerprint Recognition

It is believed that no two persons share the same fingerprints—not even identical twins—because the fingerprint patterns are part of a person’s phenotype and do not apparently depend on genetics^[5]. Fingerprints have been used to identify humans for a long time—there is some evidence that thousands of years ago ancient Chinese were aware of the uniqueness of fingerprints^[6], not speaking about their current use in forensic science and law enforcement. The traditional fingerprint acquisition mechanism—finger into ink and then on to paper—obviously is not usable in many—if not most—noncriminal applications.

Currently there are four known inkless fingerprint acquisition mechanisms considered suitable for use in practical biometrics.

Optical Sensing

Optical fingerprint sensing works by acquiring light reflected from the finger surface through a special prism. The result is an image of the finger surface. The downside of this method is that wet, dirty, or dry finger skin may result in a bad image.^[7]

Thermal Sensing

With the thermal sensing method, a thermogram of the finger surface is taken and the resulting image is used.^[8]

Capacitance Sensing

Because of differing capacitance of the ridges and valleys of fingers, a *Complementary Metal-Oxide Semiconductor* (CMOS) capacitance sensor can obtain an image of the finger when it is touched. However, like optical sensing, capacitance sensing may be negatively affected by dry, dirty, or wet skin.^[9]

Ultrasound Sensing

Ultrasound sensing works by using an ultrasound beam to scan the skin surface. Ultrasound sensing is not affected much by dry, dirty, or wet skin but takes longer to perform and the ultrasound sensing equipment is usually not compact and consequently not widespread.^[10]

In addition to the mentioned issues of wet, dry, or dirty skin, numerous other factors may also affect the quality or the very possibility of taking a fingerprint. For example, although the absolute majority of people have at least one finger, many people may also have damaged skin or skin illnesses that may degrade the quality of fingerprints or render them unusable. Fingerprint matching approaches may be broadly categorized into three classes: feature techniques, imaging techniques, and hybrids of the two. In feature-based fingerprint matching techniques, a symbolic representation of the fingerprint, defined by so-called *minutiae*, is created from the fingerprint image, and it is this representation that is later stored and used to match fingerprints—not the raw fingerprint image itself^[11]. Imaging techniques use the fingerprint images directly—image correlation algorithms are then used to compare the fingerprints^[12].

The Mighty Fingers

If the defending technology is expensive and complex, it does not mean the attacking technology will also be complex and expensive—this has been proven by many successful security attacks. Tsutomu Matsumoto of the Yokohama National University successfully fooled numerous fingerprint readers into accepting fake fingers made of gelatin with a 80-percent success rate, sending a shock wave among biometrics proponents^[13].

In a paper ambiguously entitled “Impact of Artificial Gummy Fingers on Fingerprint Systems,” co-authored with H. Matsumoto, K. Yamada, and S. Hoshino and presented at the Optical Security and Counterfeit Deterrence Techniques IV conference (Proceedings of the *International Society for Optical Engineering*, 2002), Matsumoto describes relatively easy ways to create artificial clones of fingers using cheap and freely available materials such as gelatin, free molding plastic, and photosensitive printed circuit boards.

Not only was he able to create a copy of a live finger that was good enough to fool most fingerprint readers used in the experiment, he also created an artificial finger using a latent fingerprint left on a glass, which was also accepted as genuine. In addition, Matsumoto mentions several other attack vectors against fingerprint systems, including instances where the registered finger is presented by an armed criminal, under duress, or on a sleeping drug; a severed fingertip of the registered finger; or a genetic clone of the registered finger.

Even if we disregard the last possibility as too expensive and unlikely, the others are indeed very real and must be disturbing to current users of fingerprint-based identification or verification systems. After this research was published, Bruce Schneier wrote in the May 2002 issue of his monthly newsletter CRYPTO-GRAM^[14]:

“There’s both a specific and a general moral to take away from this result. Matsumoto is not a professional fake-finger scientist; he’s a mathematician. He didn’t use expensive equipment or a specialized laboratory. He used \$10 of ingredients you could buy, and whipped up his gummy fingers in the equivalent of a home kitchen. And he defeated eleven different commercial fingerprint readers, with both optical and capacitive sensors, and some with “live finger detection” features. (Moistening the gummy finger helps defeat sensors that measure moisture or electrical resistance; it takes some practice to get it right.) If he could do this, then any semi-professional can almost certainly do much much more. More generally, be very careful before believing claims from security companies. All the fingerprint companies have claimed for years that this kind of thing is impossible. When they read Matsumoto’s results, they’re going to claim that they don’t really work, or that they don’t apply to them, or that they’ve fixed the problem. Think twice before believing them.”

Face Recognition

One of the most powerful drivers behind the use of face recognition is the fact that we all use face recognition every day to recognize people—so it seems to be one of the most acceptable biometrics we have (unlike, for example, fingerprints, which are often associated with criminal prosecution), not speaking about photographs that have been used for identification for many years^[15]. However, despite progress in this area of biometrics, face recognition is still not accurate and dependable enough, and factors such as aging, changing hairstyles, beards, and moustaches only make reliable face recognition more difficult. Bruce Schneier, in his recent book *Beyond Fear*, had the following to say about the usefulness of face recognition systems^[16]:

“I’ll start by creating a wildly optimistic example of the system. Assume that some hypothetical face-scanning software is magically effective (much better than is possible today)—99.9% accurate. That is, if someone is a terrorist, there is a 1-in-1000 chance that the software fails to indicate “terrorist” and if someone is not a terrorist, there is a 1-in-1000 chance that the software falsely indicates “terrorist.” In other words, the defensive-failure rate and the usage-failure rate are both 0.1%. Assume additionally that 1 in 10 million stadium attendees, on average, is a known terrorist (this system won’t catch any unknown terrorists who are not in the photo database). Despite the high (99.9%) level of accuracy, because of the very small percentage of terrorists in the general population of stadium attendees, the hypothetical system will generate 10,000 false alarms for every one real terrorist. This would translate to 75 false alarms per Tampa Bay football game and one real terrorist every 133 or so games.”

Of course these issues do not apply exclusively to face recognition systems, but we get the idea—a system that generates so many false alarms and catches so few terrorists is not going to be successful. This was proven on several occasions. First at the Palm Beach International Airport, where a face recognition system failed by providing less than 50-percent recognition rate and generating a large number of false positives, resulting in a decision by the airport not to use the system at all^[17]. Almost the same happened in the second case, at a face recognition system trial at the Boston Logan International Airport^[18].

Hand Geometry

Features measured and used by hand geometry biometrics typically include length and width of fingers, different aspect ratios of palm and fingers, thickness and width of the palm, and so on^[19]. Existing hand geometry systems mostly use images of the hand. Like face recognition, hand geometry is a user-friendly technology that scores higher on the acceptability test than, for example, fingerprints. It is also relatively more easily measurable and recordable than some other biometrics. Several patents have been issued for hand geometry systems, but there is not as much research as on fingerprints^[20]. However, because of its biometric properties, hand geometry is not suitable for use in the identification mode.

Iris Recognition

Iris recognition-based biometric systems are believed to be very reliable and accurate^[21]. Like fingerprints, the iris image is a part of human phenotype and is believed to be unique in every individual. Perhaps one of the most known cases of deployment of the iris recognition system is the Privium at Amsterdam’s Schiphol International Airport. Frequent travelers may enroll in the system to enjoy fast border crossing by simply looking at the iris scanner, which authenticates the person and opens the gate^[22]. In February 2004, an iris recognition system will also be piloted at the Frankfurt International Airport, and if the six-months-long trial concludes successfully, the system may be installed and deployed in 18 European countries^[33]. Obviously, iris recognition would not work for people who are missing both eyes or who have serious eye illnesses that affect the iris.

Behavioral Biometrics

Two of the most used behavioral biometrics are signature- and voice-based systems. Another behavioral biometric, keystrokes (where the timing between successive key pressings is used), seems to receive increasing attention and use.

Signature

In use for centuries, signatures enjoy a high degree of acceptance, largely because of their everyday use and familiarity, but as a behavioral biometric, signatures lack permanence: they may change at the will of a person, or under influence from such factors as illness, mental state, medicines, emotions, or age. For these and other reasons, signature-based biometric systems function in the verification and not in the identification mode.

Two subtypes of signature verification systems exist: static signature verification systems, where only the graphical representation (image) of the signature is used, and dynamic signatures, where the dynamics, pressure, and speed of the movement of a special pen are used for verification. Although the first method does not require any special hardware, the dynamic signature verification requires the use of special electronic signature readers or high-quality tablets. It is understood that dynamic signature verification is more secure and reliable than static signatures^[23]. However, some people do not have consistent signatures, resulting in increased false rejection rates to unacceptable levels and severely affecting the practical use of signature-based biometric systems.

Voice

Voice recognition systems (not to be confused with speech recognition systems, which are concerned with the actual words said and not the identity of the speaker) depend on numerous characteristics of a human voice to identify the speaker. Voice recognition holds much potential because it is acceptable and it does not require expensive input devices, unlike some other biometrics. Like face recognition, voice recognition is something we humans do many times a day; additionally, voice recognition is ideal for many practical and widespread telephony applications, and in theory voice recognition systems may even function in the background without forcing the users to go through a separate identification and verification process, saving us from another password to remember. But as usual, voice recognition systems also have their fair share of potential problems. As we all know, some people with exceptional vocal abilities may skillfully imitate others' voices, potentially defying such systems. Another issue is the ease of sound recording and replay, so any voice recognition system must be designed to withstand "record and replay" attacks.

Voice recognition also is influenced by the usual suspects—illness, mental state, emotions, age—which may substantially modify an enrolled subject's voice to a degree that it does not match the stored templates anymore. Several voice recognition models varying in accuracy and complexity exist.

The *fixed-text* model involves a person saying a word or phrase previously recorded and enrolled in the biometric database. The verification process is the simple comparison, possibly accounting for some allowable differences. However, if this word or phrase can be recorded, the entire system fails, because it is fairly easy to reproduce words and phrases.

Another model is *text-dependent*, meaning the system instructs the person to speak words or phrases—naturally this system is less prone to replay attacks because supposedly the person does not know in advance what words or phrases the system will ask for. A hybrid system, also known as *conversational voice verification*, combines something you are—your voice—and something you know—such as a password—to provide a higher degree of verification accuracy and reliability, and this system may well be the best choice in practice^[24], so multimodal biometrics may hold the key to more accurate and practical biometric authentication. Again, we should keep in mind that some people cannot use this biometric for one reason or another.

System and Design Issues

The following is a quick overview of only some of the most important biometric system design and implementation considerations:

Security

Biometrics is invariably associated with security, hence the biometric system itself should be reasonably secure and trustworthy. Not only should the system provide the required functionality, but we also should have a degree of security assurance. Keeping in mind our track record of creating secure complex systems (almost an oxymoron), we should not really have high expectations this time either. If we have learned a lesson, it is that systems fail and malfunction, so recovery and compensating mechanisms should be in place from the beginning, and even the most sophisticated system should be expected to fail sooner or later, one way or another. Some of the biometrics security issues are discussed in the following section.

Rogue Sensors and Unauthorized Acquisition (theft) of Biometric Samples

One of the risks associated with the use of biometrics for identification or verification is that a biometric cannot be changed by definition—your fingerprint is your fingerprint and there is no easy way to change it—so if it is stolen and used to create a fake finger to impersonate you, there is not much you can do about yours. Therefore, the issue of mutual authentication of the individual and the sensor is of much importance. In practice, however, as illustrated by numerous stories about rogue *Automated Teller Machines* (ATM) harvesting unsuspecting victims' card and PINs, this would prove to be a difficult task. Unlike, for example, smart cards, which may use cryptographic protocols to establish with whom they are communicating, we humans have no secure way to ascertain whether the biometric reader attached to a computer somewhere is indeed under control of (let's say) a genuine Internet banking application and will not relay or store our biometric template without authorization.

In contrast, bank customers asked to authenticate themselves at a bank counter may have a reasonable expectation that their biometric will be used by the same bank for lawful purposes only—because of their and the sensor’s physical location (so called location-based authentication). Still, unauthorized acquisition and use of biometrics remains one of the issues to be considered in any practical implementation.

The fact that not all biometrics require placing your finger on a fingerprint reader (such as face recognition systems) and that some biometric samples may be obtained without any action on part of the subject is further food for thought because one’s biometrics may be acquired without knowledge or authorization.

Communications Security Between Sensors, Matchers, & Biometric Database(s)

Although as important as the previous issue, communications security between sensors, matchers, and biometric databases is easier to provide than to solve the problems of mutual authentication of humans and biometric sensors. Well-designed and well-implemented secure cryptographic protocols may provide the required security for sensitive data exchange between parts of a biometric identification or verification system, and they are unlikely to be the weak link in the biometrics chain.

Accuracy

A biometric system must be reasonably accurate—otherwise why would we need it? The widely used FAR and FRR, and their product, CER, are not really exact measures but often estimates made using assumptions—and these assumptions may not be reasonable in all circumstances.

Speed

Although the question of how fast the system works may not be a pressing issue in, say, a nuclear reactor access control system, it will be a crucial factor at installations such as airports or border crossing points where a large number of people needs to be reliably and quickly identified and authenticated.

Scalability

Biometric verification systems are significantly and inherently more scalable than biometric identification systems particularly because only one-to-one matching is required. A distributed, combined system using smart cards that store the owner’s biometric template and compare the provided biometric in card is an example of a scalable distributed biometric verification system. However, as the previously described face recognition system experiences at airports show, system properties such as FRR must be considered in context—one false rejection a month may be acceptable, but a hundred false rejections a day clearly would not. Another scalability issue is the nature of biometrics. A scalable biometric—such as the iris—can theoretically be deployed on a large scale (with thousands or millions of enrolled users), but a biometric with weak scalability could provide acceptable error rates and performance only in small installations. Therefore, scalability is directly linked with the particular type of biometric used, and this seems to be accounted for by the International Civil Aviation Organization (see the section “Biometrics and Passports”).

Resilience

A biometric system should be able to handle exceptions. An exception in this context might be a person without the required biometric or a person whose biometric may not be usable for some reason. In many cases exception handling means resorting to a manual process, which of course brings all the issues of human intervention (speed and social engineering, to name only two) with it and may mean life or death for a particular system or application.

Cost

Because laws of economics apply to almost every human activity, a biometric system should be reasonable in cost. Of course reasonableness of cost is a very subjective concept and would vary greatly between different environments and different uses.

Privacy

As mentioned in the beginning of this article, biometrics is argued to be one of the threats to privacy and anonymity in the modern age. The *Electronic Frontier Foundation* (EFF) lists the following as being the most important privacy concerns:

- Biometric technology is inherently focused on individuals and interfaces easily to database technology, making privacy violations easier and more damaging.
- Biometric systems are useless without a well-considered threat model.
- Biometrics are no substitute for quality data about potential risks.
- Biometric identification is only as good as the initial ID.
- Biometric identification is often overkill for the task at hand.
- Some biometric technologies are discriminatory.
- Biometric systems accuracy is impossible to assess before deployment.
- The cost of failure is high.

Indeed it is very depressing to imagine a society—or even worse, a world order—where everyone is forced into a biometric database and total control over all your actions and whereabouts during your entire life is maintained—and where you can never “change your username” or “log out.” One cannot help but remember Benjamin Franklin’s immortal statement that those who are willing to trade liberty for security deserve neither. However depressing, this image hopefully will not materialize—and to achieve that, biometric systems should provide reasonable privacy and specific use guarantees to the enrolled subjects; in addition, they must have effective systems of checks and balances to audit and assure conformance with these guarantees.

Standards in Biometrics

As Andrew Tanenbaum once supposedly said, the good thing about standards is that there are so many to choose from—regardless of whether he did or not, this statement perhaps does not yet seem to apply to biometrics standards.

- The *Common Biometric Exchange File Format* (CBEFF) describes a set of data elements necessary to support biometric technologies in a unified way, and provides for the exchange of security, processing, and biometric data in a single file. The U.S. *National Institute for Standards and Technology* (NIST) describes CBEFF as facilitating interoperability between different systems or system components, forward compatibility for technology improvement, and software/hardware integration^[26].
- *BioAPI and Human Authentication API*. BioAPI and HA-API efforts merged in 1999 under the umbrella of the BioAPI Consortium. The current version of the BioAPI Specification is Version 1.1, which aims to provide a “standardized *Application Programming Interface* (API) that will be compatible with a wide range of biometric applications and a broad spectrum of biometrics technologies”^[27].
- The Open Group’s *Human Recognition Services* (HRS) is a module of the *Common Data Security Architecture* (CDSA), which in particular is used in Apple’s Mac OS X. HRS is compatible with the CBEFF and, thanks to the CDSA modular and layered approach, can use services provided by other CDSA modules^[28].
- *Biometrics Management and Security for the Financial Services Industry* (ANSI X9.84-2000) specifies minimum security requirements for effective use of biometrics data in the U.S. financial services industry, including collection, distribution, and processing of biometrics data. In particular, it specifies the security of the physical hardware used throughout the biometric life cycle; the management of the biometric data across its life cycle; the use of biometric technology for verification or identification of bank clients and employees; and other aspects. The data objects specified in X9.84 are compatible with CBEFF^[29].
- The *American Association of Motor Vehicle Administrations* (AAMVA) *Driver’s License and Identification* (DL/ID) standard provides a uniform way to identify holders of driver license cards within the United States and Canada. This standard specifies identification information on drivers’ license and ID card applications, provides for inclusion of fingerprint data, and is compatible with BioAPI and CBEFF^[30].
- *ANSI/NIST Data Format for the Interchange of Fingerprint, Facial, Scar Mark, and Tattoo Information* (ANSI/NIST-ITL 1-2000). This standard defines the content, format, and measurement units for the exchange of the specified information that may be used for identification of persons, and it is mainly directed at U.S. law enforcement agencies and government.^[31]

Additionally, one of the groups of the *International Organization for Standardization* (ISO) is working toward inclusion of biometrics specifications in the widely used ISO 7816 standard for smart cards (Part 11: personal verification through biometric methods)^[32].

Practical Uses of Biometrics

Because there may be as many practical uses of biometrics as users, we address just two of them: the use of biometrics for network authentication and the use of biometrics in international travel documents.

Biometrics for Network Authentication

As we saw earlier in this article, the accepted and widely used *what you know* and *what you have* authentication methods are not always—nor are they necessarily—secure or convenient, and they have their share of weaknesses.

The additional challenge of using biometrics for network authentication is the fact that the subject and the object of access are separated by a (usually uncontrolled, untrusted, and possibly hostile) network, which does not add to the simplicity or security of the system as a whole. As illustrated by the case of gelatin fingers described earlier, the question of whether a live person provided the biometric to a remote biometric sensor is even more important in network authentication applications when there are no preventive or detective controls, such as a watching guard, in place.

Although we have relied mostly on passwords to serve as the only or the main authentication mechanism until today, it has been clear for a while that passwords do not provide strong authentication. Keeping this lesson in mind, a biometric network authentication system should not depend solely on biometrics but should use one of the other authentication methods (what you know or what you have) as well.

The remote biometric sensors required in any biometric network authentication system are one of the most vital parts of the entire system, yet they are most vulnerable ones as well. For our purposes, we define the remote part of a centralized network authentication system as including a human user who needs to be authenticated as being physically present at the site and time of authentication, a general-purpose computer running a general-purpose operating system, and a special-purpose biometric sensor device directly connected to the general-purpose computer. This setup, therefore, includes the following high-level potential points of attack:

1. User
2. Path from the user to the sensor
3. Biometric sensor
4. Path from sensor to the general-purpose computer
5. Network
6. The central database

Even if the central authentication database is left out of the picture, the most simple risk assessment would reveal, among others, the following issues:

1. The user should be accurately identified or the declared identity should be verified; the sensor should be able to differentiate between a live human being providing live biometric and a biometric replica, such as an iris photograph or a gelatin finger. This includes, *inter alia*, reasonable assurance of the physical presence of the whole individual and not just the particular biometric at a particular point in time (hence, in part, the need for multimodal authentication involving not only what you are but also what you know or what you have).
2. The sensor should be sufficiently tamper-proof to withstand a defined set of attacks by a defined class of attackers, which would of course differ from environment to environment.
3. The communication protocol used between the sensor and the general-purpose computer should be simple, well-defined, and verified.
4. The role of the (untrusted) general-purpose computer and its software in such a system should be kept to a minimum. The biometric data acquired by the sensor should be cryptographically protected (encrypted and signed with the device key, for instance) inside the same sensor, without any dependence on action or inaction of the general-purpose computer. Their only role in this play should be to relay the bits from the sensor to the central authentication server for verification. Confidentiality and integrity of the biometric data should not be affected by a malicious, general-purpose computer or its software; the worst that can happen is the nondelivery of such data to the central authentication database.

An example of this approach would be a tamper-resistant fingerprint reader able to accurately recognize live human fingers (and reject fake ones), extract the required information, append a time stamp from an internal independent time source, encrypt and sign the resulting minutiae + time stamp data block using some digital signature algorithm, and send the resulting information through, for example, a *Universal Serial Bus* (USB) connection to the general-purpose computer. The general-purpose computer may then use the provided token to seek authentication from the central authentication database, provided all other requirements have been met.

Today a variety of network authentication systems that use or can use biometrics are available from numerous vendors. Aside from the objectively subjective information provided by vendors of such systems, little evidence of assurance exists that could enable potential users to evaluate them for their particular environments. The fact that most of these systems run as applications on the most widespread and arguably the least secure of operating systems perhaps speaks for itself.

Biometrics and Passports

For many years now more than 110 nations have issued machine-readable travel documents (mainly passports and visas) that conform to the *International Civil Aviation Organization* (ICAO) standard 9303. ICAO, a United Nations specialized agency, in addition to being responsible for international civil aviation matters, is also mandated to develop and adopt international standards on customs and immigration documents and procedures under the Chicago Convention. These machine-readable travel documents include a two-line area printed in *Optical Character Recognition* (OCR) B format, which contains information usually required for international travel (such as a person's name, date of birth, citizenship, document validity dates, and other information). These documents have greatly reduced the time necessary to check passports and visas by border officials, and have contributed to smoother international travel. In May 2003, the ICAO adopted a set of documents on integration of biometrics into machine-readable passports, choosing three most suitable for these purposes^[25]. The main biometric chosen was a digitized face image, followed by two optional biometrics: fingerprints and irises. The ICAO also selected high-capacity, contactless smart cards as the storage method for this biometric data and gave other recommendations related to integration and use of biometrics in passports and other documents. It remains to be seen if or how and when 188 member states of the ICAO will integrate biometrics into their passports.

New Biometrics

It would be unreasonable to assume that we are aware of all possible biometrics. It may very well be the case that new biometrics are discovered and possibly, in the fullness of time, considered fit for practical use. An example would be a behavioral biometric proposed by Ross Anderson of Cambridge University, author of the already classic *Security Engineering*:

“Are there any completely new biometrics that might be useful in some circumstances? One I thought up while writing this chapter, in a conversation with William Clocksin and Alan Blackwell, was instrumenting a car so as to identify a driver by the way in which he or she operated the gears and the clutch.”

Summary

Biometrics is a promising and exciting area, where different disciplines meet and provide an opportunity for a more secure and responsible world. However, the same biometrics, if misused or poorly engineered, may instead bring many hassles—if not troubles. Some biometrics are less usable than others, and different environments warrant different biometrics and design considerations. The best advice would be to differentiate between market-ready biometric technologies and technologies that are not yet (if ever) ready for deployment outside testing grounds. However much fervent proponents and keen vendors of biometric solutions market their wares, the guiding factor should be proven reliability and appropriateness of these solutions to specific uses, not marketing hype, which seems at times to dominate this arena.

Organizations and Publications

The following organizations and publications may be useful sources of further information on biometrics and biometric applications:

The International Biometric Society: www.tibs.org

Biometric Consortium: www.biometrics.org

BioAPI Consortium: www.bioapi.org

International Biometrics Industry Association: www.ibia.org

International Association for Identification: www.theiai.org

Journal of the International Biometric Society:
stat.tamu.edu/Biometrics

Biometric Digest: www.biodigest.com

Biometric Technology Today: www.biometrics-today.com

Additionally, the following books may serve as good introductions to biometrics:

Guide to Biometrics, by Bolle, Connell, Pankanti, Ratha, Senior, ISBN 0-387-40089-3, Springer Verlag, 2003

Practical Biometrics, Julian Ashbourn, Springer Verlag, 2003

One of the best publicly available works on security engineering is *Security Engineering: A Guide to Building Dependable Distributed Systems*, by Ross Anderson (Wiley, 2001).

References

- [1] http://www.dhs.gov/dhspublic/interapp/editorial/editorial_0333.xml
- [2] <http://news.bbc.co.uk/2/hi/americas/3358627.stm>
- [3] http://www.usatoday.com/tech/news/techpolicy/2003-08-24-biometrics-travel_x.htm
- [4] <http://www.anu.edu.au/people/Roger.Clarke/DV/HumanID.html>
- [5] "On the individuality of Fingerprints. Pankanti," Prabhakar, Jain; *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, December 2001.
- [6] "The History and Development of Fingerprinting," Lee, Gaensslen; *Advances in Fingerprint Technology*, CRC Press, 1994.
- [7] *Guide to Biometrics*, Bolle et al., Springer Verlag, 2003.
- [8] "Fingerchip: Thermal Imaging and Finger Sweeping in a Silicon Fingerprint Sensor," Mainguet, Pegulu, Harris; *Proceedings of AutoID 99*, October 1999.

- [9] “Low-power and high-performance CMOS Fingerprint Sensing and Encoding Architecture,” Jung, Thewes, Scheiter, Gooser, Weber; *IEEE Journal of Solid-State Circuits*, July 1999.
- [10] “Ultrasound Sensor for Fingerprint Recognition,” Biez, Gurnienny, Pluta; *Proceedings of SPIE—Optoelectronic and Electronic Sensors*, June 1995.
- [11] “A Tree System Approach for Fingerprint Pattern Recognition. Moayer,” Fu; *IEEE Transactions on Computers*, C-25(3).
- [12] *Guide to Biometrics*, Bolle et al., Springer Verlag, 2003
- [13] <http://www.itu.int/itudoc/itu-t/workshop/security/present/s5p4.pdf>
- [14] <http://www.schneier.com/crypto-gram-0205.html#5>
- [15] “Face Recognition: Features versus Templates,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 12(10), October 1993.
- [16] *Beyond Fear: Thinking Sensibly about Security in an Uncertain World*, Bruce Schneier; Copernicus Books, 2003.
- [17] <http://www.theregister.co.uk/content/archive/25444.html>
- [18] <http://www.theregister.co.uk/content/archive/26298.html>
- [19] “A Hand Shape Identification System,” Biometric Systems Lab, <http://bias.csr.unibo.it/research/biolab/hand.html>
- [20] U.S. Patent 3,576,537; U.S. Patent 3,648,240
- [21] “Iris Recognition: An Emerging Biometric Technology,” Wildes; *Proceedings of the IEEE*, 85(9), September 1997.
- [22] http://www.schiphol.nl/schiphol/privium/privium_home.jsp
- [23] “Automatic On-line Signature Verification,” Nalwa; *Proceedings of the IEEE*, 85(2), February 1997.
- [24] “Speaker Recognition,” Campbell, in *Biometrics: Personal Identification in Networked Society*, by Jain, Bolle, Pankanti, ISBN 0-7923-8345-1, Kluwer Academic Publishers, 1999.
- [25] <http://www.icao.int/mrtd/download/technical.cfm>
- [26] <http://www.itl.nist.gov/div895/isis/bc/cbeff/CBEFF010301web.PDF>

- [27] <http://www.bioapi.org/>
- [28] <http://www.opengroup.org/security/cdsa.htm>
- [29] <http://www.ansi.org>
- [30] <http://www.aamva.org>
- [31] ftp://sequoyah.nist.gov/pub/nist_internal_reports/sp500-245-a16.pdf
- [32] <http://www.iso.org>
- [33] http://news.com.com/2100-7348_3-5158973.html

The author of this article does not work for, is not affiliated with, and has no financial interest or shareholding in any vendor of any biometric technology at the time of submission of this article for publication.

EDGAR DANIELYAN, CISSP, is a self-employed consultant, published author, editor, and instructor specializing in information security, UNIX, and internetworking. He is the principal partner at Danielyan Consulting LLP (www.danielyan.com), an information security assurance consultancy, and a member of ACM, IEEE, ISACA, USENIX, and the British Computer Society's Information Security Specialist Group.
E-mail: edd@danielyan.com

Book Reviews

The Unicode Standard

The Unicode Standard, Version 4.0, by The Unicode Consortium, ISBN: 0-321-18578-1, Addison Wesley Professional, 2003.

The Unicode 4.0 book is a thick, heavy one, but it is good. If you work with the Unicode character set, you should have this book on your bookshelf.

This book consists of four parts:

- Background and explanation of terms (103 pages)
- Implementation guidelines (29 pages)
- Technical specifications (60 pages)
- The Unicode Character Tables (1150 pages)

A review must describe each of these sections by itself, because they are important for different reasons. Unfortunately, however, the sections in the book are not clearly divided into sections as I outlined, so you don't necessarily know where to start. You don't need to read the characters section—just the sections you are interested in.

You should read the “Preface” (Section 0), because this section describes the rest of the book. It starts on page xxxi (before chapter 1).

You can then immediately go to the section you are interested in. Each section more or less stands by itself, and the book is easy to read. If something is not clear, you should look for text in another section that describes the subject. Reading from start to finish is possible, but I use this book as reference material, like an encyclopedia (except for the characters).

The background material is easy to read. It covers basic concepts such as differences between *characters* and *glyphs*, definition of terms such as *equivalence*, character encoding schemes and implication of things such as bidirectional text (mixed right-to-left and left-to-right text). Knowing how these things work is essential for anyone who either implements text engines of any kind or works on developing protocols or standards. This background material is easier understood read on paper and not electronically. It also is the part of the book I return to most often.

The second very good part concerns implementation guidelines. Even though it is (relatively) short, it is very important material. It discusses selection algorithms and other user interface guidelines, as well as other algorithms needed for, for example, comparison (what is called “Normalization”). I like this section as well, because it really describes the details you need to know when implementing anything Unicode related.

Unicode is a *large* character set. You see that in the more-than-1000 pages of “just characters.” Of course, the tables themselves can be found on the Unicode Consortium Web site, but this book gives you a good overview. Part of this overview is a description of the *scripts* that Unicode covers, one at a time before the *codepoints* that come from those scripts. Still, this is the part that makes this book heavy, and a version without the codepoints would have been interesting by itself.

The book ends with more technical material, consisting mostly of references to, for example, *Unicode Technical Notes* and other standards documents that the Unicode Consortium produces, in addition to the Unicode Standard itself.

Useful reference

In summary, the first 130 pages (well, starting at page 40) in the book are very good. If you work at all with Unicode, you should read those pages. The rest of the book is good reference material.

Even though I have been working with Unicode for almost 10 years now, and for the last 8 years have weekly reviewed Unicode-related standards in the Internet Engineering Task Force, I see myself opening this book now and then. There is always something I need to check, and to be honest, I like encyclopedias on paper.

As reference material, this is a must-have item. If you want to read only the 140 interesting pages once, well, the book is possibly overkill.

—Patrik Fältström, Cisco Systems
paf@cisco.com

iSCSI: The Universal Storage Connection

iSCSI: The Universal Storage Connection, by John L. Hufferd, ISBN 0-201-78149-X, Addison-Wesley, 2003.

I have to come clean straightaway and say that when I received this book to review I had never even heard of *Internet Small Computer System Interface* (iSCSI) and, to be honest, I have never heard it mentioned by anyone again since the day the book arrived. This is, of course, not a criticism of this book, just a comment on the current state of penetration of iSCSI into everyday computing discourse. In fact, if you search Google for “iscsi,” you get only 465,000 hits—very few indeed these days, though this does have the decided advantage that the links you get are generally pretty useful. I’m sure that this will change because there are lots of big names behind the protocol, and certainly when vendors start really selling kit that uses it. *Storage Area Networks* (SANs) are important (though also not yet at the forefront of most computing people’s minds)—and iSCSI will probably make them bigger.

However, to the book. And, really, if you want to know pretty well everything about iSCSI and don't want to read lots of Web sites, then this book is for you. It covers everything from the background behind the protocol, to how and where it might be applied, to all the low-level information that most of us hope that we never need to see. I'm not going to list it all and go into detail: the whole thing is here, from soup to nuts.

As to the presentation of the material, it is excellent—clear diagrams and useful tables. The layout is spacious without huge amounts of wasted white space on every page—making a change from many textbooks you see today.

The writing is clear too, though I did find myself becoming a bit bogged down in all the abbreviations (no, not acronyms—most of them are not words), which seem to pile up in the sentences. I got a bit tired of seeing iSCSI everywhere after a while too. I wasn't keen on the end-of-chapter summaries, finding them a bit redundant.

Good Reference

All in all, if you are in a position where you need to know about iSCSI and may have to be involved in working with it at a low level, this book is a good reference. I doubt that there is anything more comprehensive or better written at the present time.

—*Lindsay Marshall, University of Newcastle upon Tyne*
`lindsay.marshall@newcastle.ac.uk`

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases we may be able to get a publisher to send you a book for review if you don't have access to it. Contact us at **`ipj@cisco.com`** for more information.

NRO Comments Concerning ICANN and WSIS

The *Number Resource Organization* (NRO) is the coalition of *Regional Internet Registries* (RIRs) which operate in the world today. The NRO is an organization representing the collective experience of individual RIRs and their communities. While the prime subject of its work are matters of joint interest relating to Internet numbering resources, the NRO provides an efficient interface to other parties interested in these issues. As the Internet continues to evolve, the NRO will ensure continuity of the operational infrastructure of Internet number resource allocation.

The RIRs are responsible for distribution of *Internet Number Resources* [IPv4 and IPv6 addresses and *Autonomous System Numbers*]. These number resources are the most fundamental of the identifiers on which the Internet relies: the Internet can operate without domain names; but it cannot operate without numbers. The RIRs have carried the responsibilities associated with managing these critical resources collectively for over 10 years, since well before the start of ICANN. This has been done very effectively through the entire “modern history” of today’s Internet which includes both the “dot com boom” and the “dot com bust.”

The RIRs have participated in the *World Summit on the Information Society* (WSIS) processes for over a year, including regional Prepcoms and the Summit itself. This is probably longer than any other Internet organization. The RIRs have attended as observers, and as subject matter experts with a genuine aim to assist in debates and discussions around issues related to Internet Number Resources in general and to IP addresses in particular.

The RIRs participated in the WSIS Phase I process as full supporters of ICANN as the model which represents not only the fundamental and critical aspects of Internet development to date, but also the means of community self-regulation to administer and manage Internet Number Resources. It must be understood that this is not given by the RIRs as mere components of ICANN, dependent upon it for support; but rather as independent components of the broader Internet administrative framework which ICANN itself is intended to support.

In the second round of WSIS, the NRO speaking for the collective RIRs will assert an active role vis-à-vis ICANN in order to aid that organization to address the genuine questions that it faces. The principle of these issues within the WSIS context is that of the independence and genuine internationalization of ICANN.

Therefore the NRO calls on ICANN to continue its work in this area, not by building a multinational organization, but rather by including and gaining the genuine support of its significant base of core stakeholders, namely those in the DNS, IP address, and protocol communities. Furthermore, the NRO calls on ICANN to work with the US Government to demonstrate a genuine and unambiguous plan for its independence and to commit to this plan before the conclusion of the second phase of the WSIS.

Finally, the NRO rejects any concept of an alternative Internet administrative model located within any governmental or intergovernmental structure. The NRO acknowledges that there is a valid role for governments in the administration of the Internet but this must be in the context of the current model. There is a need for the continual improvement of the current model of industry self-regulation to the extent that the ultimate solution may look little like today's ICANN.

<http://www.apnic.net/index.html>

<http://www.arin.net/index.html>

<http://www.lacnic.net/>

<http://www.ripe.net/index.html>

Upcoming Events

INET/IGC 2004 will be held in Barcelona, Spain, May 10–14, 2004. INET, which is the annual conference of the *Internet Society* (ISOC), will this time be held jointly with Spain's *Internet Global Congress* (IGC). For more information, visit: **<http://www.isoc.org/inet04/>**

The *North American Network Operators' Group* (NANOG) will meet in San Francisco, May 23–25, 2004. For more information see:

<http://nanog.org/>

The *South Asian Network Operators Group* (SANOG) will meet 23–30 July, 2004 in Kathmandu, Nepal. More info at:

<http://www.sanog.org/>

The *Internet Corporation for Assigned Names and Numbers* (ICANN) will meet in Kuala Lumpur, Malaysia, July 19–23, 2004, and in Cape Town, South Africa, December 1–5, 2004. For more information see:

<http://www.icann.org>

The *Internet Engineering Task Force* (IETF) will meet in San Diego, CA, August 1–6, 2004 and in Washington, DC, November 7–12, 2004. For more information, visit: **<http://ietf.org>**

The *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) will be held February 16–25, 2005 in Kyoto, Japan and February 15–24, 2006 in Bangalore, India. For more information visit: **<http://www.apricot.net/>**

This publication is distributed on an "as-is" basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Technology Strategy
MCI, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.
Copyright © 2004 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PSRST STD
U.S. Postage
PAID
Cisco Systems, Inc.

The Internet Protocol Journal

June 2004

Volume 7, Number 2

A Quarterly Technical Publication for
Internet and Intranet Professionals

In This Issue

From the Editor	1
Content Networks	2
IPv6 Autoconfiguration	12
DNSSEC	17
Book Review	29
Fragments	31

FROM THE EDITOR

The Internet Protocol Journal continues to be a forum for discussion of current and emerging technologies. In this issue, we first look at *content networking*. One can describe the Internet as a system of interconnected devices, but equally as a collection of information, called *content*, that resides on a distributed set of *servers* and is accessed by numerous *clients*. Our first article is by Christophe Deleuze.

Engineers are hard at work planning for an eventual transition to the next version of IP — IPv6. We've published several articles about IPv6 in previous editions. This time, François Donzé describes the automatic address configuration feature of IPv6. Of note is also the increasing global support for IPv6 deployment, (refer to "Fragments" on page 31).

Our final article returns to our recurring theme: adding security to existing Internet protocols. Because many malicious attacks on the Internet are perpetrated by "spoofing" information in one form or another, it makes sense to look at the *Domain Name System* (DNS), a critical component of the Internet infrastructure. Today, it is possible to create systems which provide fake answers to DNS queries. Miek Gieben explains what is being done to address this issue in his tutorial on DNSSEC, the secure version of the DNS protocols.

Please take a moment to renew or update your subscription to this journal. You can do so by visiting www.cisco.com/ipj and clicking on the "Subscription Information" link on the left. You will need to supply your subscription ID and e-mail address in order to gain access to your database record. If you have any questions, please send a note to ipj@cisco.com.

This is the 25th edition of IPJ. The journal now has more than 32,000 subscribers world-wide, and is available on paper and electronically on our Website in PDF and HTML format. The Website, located at www.cisco.com/ipj, contains all our back issues, and will soon offer a cumulative index in ASCII format that will make it easier to find particular articles. As always, we welcome your feedback.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Content Networks

by *Christophe Deleuze*

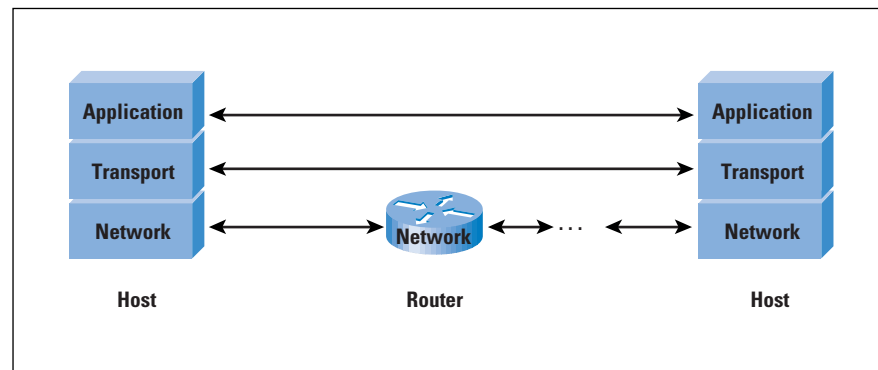
The Internet is constantly evolving, in both usage patterns and underlying technologies. In the last few years, there has been a growing interest in *content-networking* technologies. Various differing systems can be labelled under this name, but they all share the ability to access objects in a location-independent manner. Doing so implies a shift in the way communications take place on the Internet.

The Classic Internet Model

The Internet protocol stack comprises three layers, shown in Figure 1. The network layer is implemented by IP and various routing protocols. Its job is to bring datagrams hop by hop to their destination host, as identified by the destination IP address. IP is “best effort,” meaning that no guarantee is made about the correct delivery of datagrams to the destination host.

The transport layer provides an end-to-end communication service to applications. Currently two services are available: a reliable ordered byte stream transport, implemented by the *Transmission Control Protocol* (TCP), and an unreliable message transport, implemented by the *User Datagram Protocol* (UDP).

Figure 1: The Three Layers of the Internet Protocol Stack



Above the transport layer lies the application layer, which defines application message formats and communication semantics. The Web uses a client-server application protocol called *Hypertext Transfer Protocol* (HTTP)^[10].

A design principle of the Internet architecture is the “end-to-end principle,” which states that everything that can be done in the end hosts should be done there, and not in the network itself^[8]. That is why IP service is so crude, and transport and application layer protocols are implemented only in the end hosts.

Application objects, such as Web pages, files, etc. (we will simply call those “objects”) are identified by URLs. (Actually URLs identify “resources” that can be mapped to different objects called “variants.” A variant is identified by a URL and a set of request header values, but in order to keep things simple, we will not consider this in the following.) URLs for Web objects have the form `http://host:port/path`. This means that the server application lives on a host with *hostname* (or possibly IP address) on port *N* (with default value of 80), and knows the object under the name *path*. Thus URLs, as their name implies, tell where the object can be found. To access such an object, a TCP connection is open to the server running on the specified host and port and the object named *path* is requested.

Content Networks

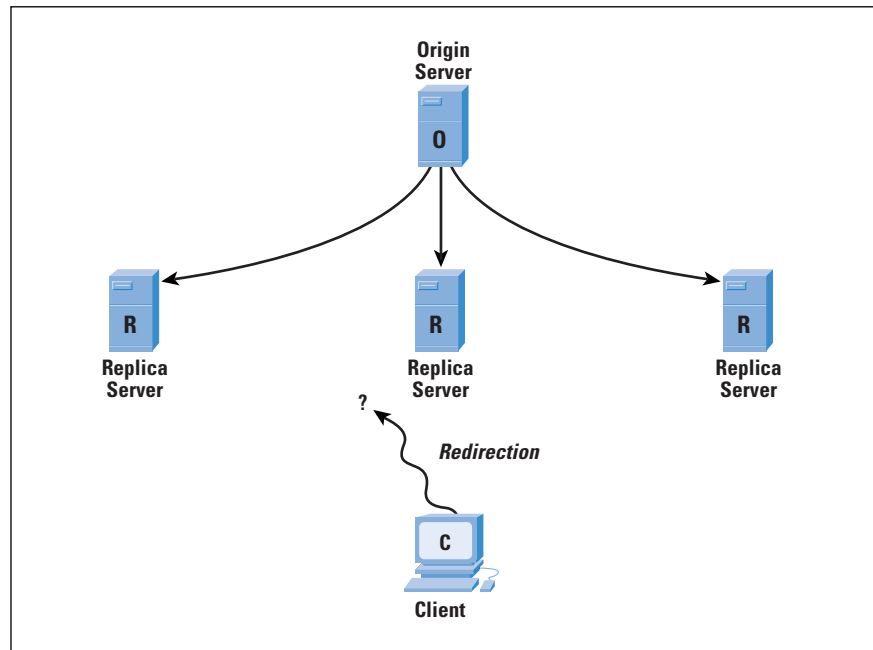
Content networks aim to provide location-independent access to an object, most commonly because they handle some kind of (possibly dynamic) replication of the objects. By design, URLs are not suited to identify objects available on several places on the network.

Handling such replication and location-independent access usually involves breaking the end-to-end principle at some point. Communication is no more managed end to end: intermediate network elements operating at the application layer (whose most common types are “proxies”) are involved in the communication. (Content networks are not the only case where this principle is violated.)

In the same way that IP routers relay IP datagrams (that is, network layer protocol data units), routing them to their destination according to network layer information, those application layer nodes relay application messages, using application layer information (such as content URLs) to decide where to send them. This is often called *content routing*.

So the goal of a content network is to manage replication, handling two different tasks: *distribution* ensures the copying and synchronization of the instances of an object from an *origin server* to various *replica servers*, and *redirection* allows users to find on instance of the object (possibly the one closest to them.) (By “replica,” we mean any server of any kind other than the origin that is able to serve an instance of the object. This term often has a narrower meaning, not applying, for example, to caching proxies.) This is illustrated in Figure 2.

Figure 2: Elements of a Content Network



Various kinds of content networks exist, differing in the extent to which they handle these tasks and in the mechanisms they use to do so. There are many possible ways to classify them. In this article, we use a classification based on who owns and administers the content network. We thus find three categories: content networks owned by network operators, content providers, and users.

Network Operators' Content Networks

Network operators (also called *Internet Service Providers*, or ISPs) often install caching proxies in order to save bandwidth^[11]. Clients send their requests for objects to the proxy instead of the origin server. The proxy keeps copies of popular objects in its cache and can answer directly if it has the requested object in cache. (To be precise, such a caching proxy does not cache objects, but server responses.) If this is not the case, it gets the object from the origin server, possibly stores a copy in its cache, and sends it back to the client.

This caching proxy scheme can be used recursively, making those proxies contact parent proxies for requests they cannot fulfill from their local store. Such hierarchies of caching proxies actually lead to constructing content-distribution trees. This makes sense if the network topology is tree-like, although there are some drawbacks, including the fact that less popular objects (those not found in any cache) experience delays, which increase with the depth of the tree. Another problem is with origin servers whose closest tree node is not the root.

The Squid caching proxy^[5] can be configured to choose the parent proxy to query for a request based on the domain name of the requested URL (or to get the object directly for the origin server). This allows setting up multiple logical trees on the set of proxies, a limited form of content routing.

Such manual configuration is cumbersome, especially because domain names do not necessarily (and actually most do not) match network topology. Thus the administrator must know where origin servers are in the network to use this feature effectively.

The same effects can be achieved, to some extent, in an automatic and dynamic fashion using ICP, the *Internet Cache Protocol* [16, 15]. ICP allows a mesh of caching proxies to cooperate by exchanging hints about the objects they have in cache, so that a proxy missing an object can find a close proxy that has it. One advanced feature of ICP allows you to select among a mesh of proxies the one that has the smallest *Round-Trip Time* (RTT) to the origin server.

One design flaw of ICP is that it identifies objects with URLs. We mentioned previously that a URL actually identifies a resource that can be mapped to several different objects called variants. Thus information provided by ICP is of little use for resources that have multiple variants. However, in practice most resources have only one variant, so this weakness does little harm.

Users normally configure their browsers to use a proxy, but automatic configuration is sometimes possible. Multiple proxies can be used by a client with protocols such as the *Cache Array Routing Protocol* (CARP)[14]. To avoid configuration issues, a common trend is for ISPs to deploy *interception proxies*. Network elements such as routers running the *Cisco Web Cache Communication Protocol* (WCCP)[6,7] redirect HTTP traffic to the proxy, without the users knowing. The proxy then answers client requests pretending to be the origin server. This poses numerous problems, as discussed in [12].

Caching proxies have limited support for ensuring object consistency. Either the origin server gives an expiration date or the proxy estimates the object lifetime based on the last modification time, using an heuristic known as *adaptive TTL* (time to live).

Content Providers' Content Networks

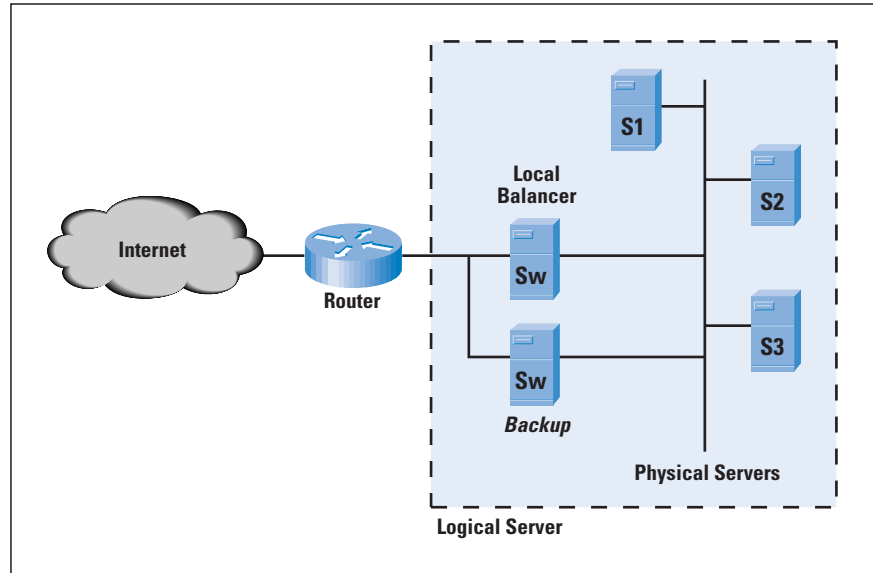
Contrary to ISPs whose main goal is to save bandwidth, content providers want to make their content widely available to users, while staying in control of the delivery (including ensuring that users are not delivered stale objects). We can again roughly classify such content networks in three subcategories:

- *Server farms*: Locally deployed content networks aimed at providing more delivery capacity and high availability of content
- *Mirror sites*: Distributed content networks making content available in different places, thus allowing users to get the content from a close mirror
- *Content-Delivery Networks* (CDNs): Mutualized content networks operated for the benefit of numerous content providers, allowing them to get their content replicated to a large number of servers around the world at lower cost.

Server Farms

Server farms are made of a load-balancing device (we will call it a *switch*) receiving client requests and dispatching them to a series of servers (the *physical* servers). The whole system appears to the outside world as a single *logical* server. The goal of a server farm is to provide scalable and highly available service. The switch monitors the physical servers and uses various load metrics in its dispatching algorithm. Because the switch is a single point of failure, a second switch is usually set up in a hot failover standby mode, as shown in Figure 3.

Figure 3: Server Farm



Some switches are called *Layer 4 switches* (4 is the number of the transport layer in the *OSI Reference Model*), meaning they look at network and transport layer information in the first packet of a connection to decide to which physical server the incoming connection should be handed. They establish a state associating the connection with the chosen physical server and use it to relay all packets of the connection. The exact way the packets are sent to the physical servers varies. It usually involves some form of manipulation of IP and TCP headers in the packets (like *Network Address Translation* [NAT] does) or IP encapsulation. These tricks are not necessary if all the physical servers live on the same LAN.

More complex *Layer 7 switches* (7 is the number of the application layer in the *OSI Reference Model*) look at application layer information, such as URL and HTTP request headers. They are sometimes called *content switches*. On a TCP connection, application data is available only after the connection has been opened. A proxy application on the switch must thus accept the connection from the client, receive the request, and then open another connection with the selected physical server and forward the request. When the response comes back, it must copy the bytes from the server connection to the client connection.

Such a splice of TCP connections consumes much more resources in the switch than the simple packet manipulation occurring in Layer 4 switches. Bytes arrive at one connection and are handed to the proxy application, which copies them to the other connection—all of this involving multiple kernel mode-to-user mode memory copy operations and CPU context switches. Various optimizations are implemented in commercial products. The simplest one is to put the splice in kernel mode. After it has sent the request to the physical server, the proxy application asks the kernel to splice the two connections, and forgets about them. Bytes are then copied between the connections directly by the kernel, instead of being given to the proxy application and back to the kernel.

It is even possible to actually merge the two TCP connections, that is, simply relay packets at the network layer to establish a direct TCP connection between the client and the physical server. This requires manipulating TCP sequence numbers (in addition to addresses and ports) when relaying packets, because the two connections will not have used the same initial sequence numbers. This can be much more complex (or even impossible) to perform if TCP options differ in the two connections.

Mirror Sites

In such a content network, a set of servers are installed in various places in the Internet, and they are defined as *mirrors* of the master server. Synchronization is most commonly performed periodically (often every night), using FTP or specialized tools such as *rsync*^[4].

Redirection is performed by the users themselves for most sites. The master server, to which the user initially connects, displays a list of mirrors with geographic information and suggests that users choose a mirror close to themselves, by simply clicking on the associated link.

This process can be automated sometimes. One trick is to store the user's choice in a *cookie*, such that the next time the user connects to the master site, the information provided in the cookie will be used to issue an *HTTP redirect* (an HTTP server response asking the client to retry the request on a new URL) to the previously selected site.

Other schemes involve trying to find which of the mirrors is closest to the user based on information provided in the user request (such as preferred language) or indicated by network metrics. Such schemes were not very common for simple mirror sites, but today many commercial products allowing for this kind of “global load balancing” are available.

In any case (except if redirection is automatic and *Domain Name System* [DNS] based—this is discussed in the next section) the URLs of objects change across mirrors.

CDNs

Most content providers cannot afford to own numerous mirror sites. Having servers in different places around the world costs lots of money. Operators of CDNs own a large replication infrastructure (Akamai, the biggest one, claims to have 15,000 servers) and get paid by content providers to distribute their content. By mutualizing the infrastructure, CDNs are able to provide very large reach at affordable costs.

CDN servers do not store entire sites of all the content providers, but rather cache a subset according to local client demand. Such servers are called *surrogates*. They manage their disk store like proxies do, and serve content to clients like mirrors do (that is, contrary to proxies, they act as the authoritative source for the content they deliver).

Because the number of surrogates can be so large, and because of the argument that “no user configuration is necessary,” CDNs typically include complex redirection systems that allow them to perform automatic and user-transparent redirection to the selected surrogate. The selection is based on information about surrogate loads and on network metrics collected by various ways such as routing protocol information, RTTs measured by network probes, etc. The client is made to connect to the selected surrogate either by sending it an HTTP redirect message, or by using the DNS system: when the client tries to resolve the host name of the URL in an IP address to connect to, it is given back the address of the selected surrogate instead. Using the DNS ensures that the URL is the same for all object copies. In this case, CDNs actually turn URLs into location-independent identifiers.

In addition to proxy-like on-demand distribution, content can also be “pushed” in surrogates in a proactive way. Synchronization can be performed by sending invalidation messages (or updated objects) to surrogates.

CDN principles are also being used in private intranets for building *Enterprise CDNs* (ECDNs).

Users' Content Networks

User-operated content networks are better known as *Peer-to-Peer* (P2P) networks. In these networks, the costly replication infrastructure of other content networks is replaced by the users, who make some of their storage and processing capacities available to the P2P network. Thus, no big money is needed, and no one has control over the content network.

One advantage P2P networks have over other content networks is that they are usually built as overlay networks and do not strive for transparent integration with the current Web. Thus they are free to build new distribution (some of them allow downloading files from multiple servers in parallel) and redirection mechanisms from scratch, and even to use their own namespace instead of being stuck with HTTP and URLs.

P2P networks basically handle the distribution part of replication in a straightforward way: the more popular an object is, the more users will have a copy of it, thus the more copies of the object will be available on the network. More complex mechanisms can be involved, but this is the basic idea.

The redirection part of replication is more problematic with most current P2P networks. It can be handled by a central directory as in *Napster*: every user first connects to a central server, updates the directory for locally available objects, and then looks up the directory for locations of objects the user wants to access. Of course, such a central directory poses a major scalability and robustness problem.

Gnutella and *Freenet*, for example, use a distributed searching strategy instead of a centralized directory. A node queries neighbors that themselves query neighbors, and so on until either one node with the requested object is found or a limit on the resources consumed by the search has been hit. Although there is no single point of failure, such a scheme is no more scalable than the central directory. It seems easy to perform denial-of-service attacks by flooding the network with requests. Additionally, you can never be sure you have found the object even if someone has it.

These examples are primitive and have serious flaws, but much research work is being performed on this topic; refer to [13] for a summary.

Although they are currently used mainly for very specific file-sharing applications, P2P networks do provide new and valuable concepts and techniques. For example, *Edge Delivery Network* is a commercially available software-based ECDN inspired by Freenet. Various projects use a *scatter/gather* distribution scheme, useful for very large files: users download several file chunks in parallel from other currently downloading users, thus refraining from using server resources for long periods of time.

Some projects attempt to integrate P2P principles in the current Web architecture and protocols. Examples are [3] and [1].

Conclusion

Current networks have been designed and deployed as ad-hoc solutions of specific problems occurring in the current architecture of the network. Caching proxies lack proper means to ensure consistency, but CDNs trick the DNS to turn URLs into location-independent identifiers. P2P networks are mostly limited to file-sharing applications.

Content networks implement mechanisms to ensure distribution of content to various locations, and redirection of users to a close copy. They often have to break the end-to-end principle in order to do so, mainly because current protocols assume each object is available in only one statically defined location.

Probably the first step in building efficient distribution and redirection mechanisms for providing an effective replication architecture is the setting up of a proper replication-aware namespace. Applications would pass an object name to a name resolution service and be given back one or more locations for this object. The need for such a location-independent namespace was anticipated a long time ago. URLs are actually defined as one kind of *Uniform Resource Identifier* (URI), another one being *Uniform Resource Names* (URNs) intended to provide such namespaces. A URN IETF working group [2] has been active for a long time, and recently published a set of RFCs (3401 to 3406).

Work on the topic of content networking has also been performed by the now closed *Web Replication and Caching* (WREC) IETF working group, which issued a taxonomy in [9]. An interesting survey of current work on advanced content networks is [13].

References

- [1] BitTorrent: <http://bitconjurer.org/BitTorrent/>
- [2] IETF URN Working Group:
<http://www.ietf.org/html.charters/urn-charter.html>
- [3] Open Content Network: <http://www.open-content.net>
- [4] Rsync: <http://rsync.samba.org>
- [5] Squid Internet Object Cache: <http://www.squid-cache.org>
- [6] M. Cieslak and D. Forster, “Web cache coordination protocol v1.0,” Expired Internet Draft, **draft-forster-wrec-wccp-v1-00.txt**, Cisco Systems, July 2000.
- [7] M. Cieslak, D. Forster, G. Tiwana, and R. Wilson, “Web cache Coordination Protocol v2.0,” Expired Internet Draft, **draft-wilson-wrec-wccp-v2-00.txt**, Cisco Systems, July 2000.
- [8] David D. Clark, “The design philosophy of the DARPA Internet protocols,” *Computer Communication Review*, Volume 18, No. 4, August 1988. Originally published in Proceedings of SIGCOMM’88.
- [9] Ian Cooper, Ingrid Melve, and Gary Tomlinson, “Internet Web Replication and Caching Taxonomy,” RFC 3040, January 2001.
- [10] R. Fielding, J. Gettys, J. Mogul, H. Frystyk, L. Masinter, P. Leach, and T. Berners-Lee, “Hypertext Transfer Protocol — HTTP/1.1,” RFC 2616, June 1999.

- [11] Geoff Huston, “Web Caching,” *The Internet Protocol Journal*, Volume 2, No. 3, September 1999.
- [12] Geoff Huston, “The Middleware Muddle,” *The Internet Protocol Journal*, Volume 4, No. 2, June 2001.
- [13] H. T. Kung and C. H. Wu, “Content Networks: Taxonomy and New Approaches,” 2002. <http://www.eecs.harvard.edu/htk/publication/2002-santa-fe-kung-wu.pdf>.
- [14] Vinod Valloppillil and Keith W. Ross, “Cache array routing protocol v1.0,” Expired Internet Draft, **draft-vinod-carp-v1-03.txt**, February 1998.
- [15] D. Wessels and K. Claffy, “Application of Internet Cache Protocol (ICP), Version 2,” RFC 2187, September 1997.
- [16] D. Wessels and K. Claffy, “Internet Cache Protocol (ICP), Version 2,” RFC 2186, September 1997.

CHRISTOPHE DELEUZE holds a Ph.D. degree in computer science from Université Pierre et Marie Curie, Paris. He worked on quality-of-service architectures in packet networks, and then spent three years in a start-up company designing CDN systems. He has also been a teacher. E-mail: christophe.deleuze@free.fr

IPv6 Address Autoconfiguration

by François Donzé, HP

Since 1993 the *Dynamic Host Configuration Protocol* (DHCP)^[1] has allowed systems to obtain an IPv4 address as well as other information such as the default router or *Domain Name System* (DNS) server. A similar protocol called DHCPv6^[2] has been published for IPv6, the next version of the IP protocol. However, IPv6 also has a stateless autoconfiguration protocol^[3], which has no equivalent in IPv4.

DHCP and DHCPv6 are known as *stateful* protocols because they maintain tables within dedicated servers. However, the stateless autoconfiguration protocol does not need any server or relay because there is no state to maintain.

This article explains the IPv6 stateless autoconfiguration mechanism and depicts its different phases.

Scope of IPv6 Addresses

Every IPv6 system (other than routers) is able to build its own unicast global address. A *unicast* address refers to a unique interface. A packet sent to such an address is treated by the corresponding interface—and *only* by this interface. This type of address is directly opposed to the multicast address type that designates a group of interfaces. Most of this article deals with unicast addresses. For simplicity, we will omit the unicast qualifier when there is no ambiguity.

Address types have well-defined destination scopes: *global*, *site-local* and *link-local*. Packets with a link-local destination must stay on the link where they have been generated. Routers that could forward them to other links are not allowed to do so because there has been no verification of uniqueness outside the context of the origin link.

Similarly, border-site routers cannot forward packets containing site-local addresses to other sites or other organizations. The IETF is currently working on a way to remove or replace site-local addresses. Hence, this article will refrain from any other reference to this address type. Finally, a global address has an unlimited scope on the worldwide Internet. In other words, packets with global source and destination addresses are routed to their target destination by the routers on the Internet. A fundamental feature of IPv6 is that all *Network Interface Cards* (NICs) can be associated with several addresses.

At minimum, a NIC is associated with a single link-local address. But in the most common case a NIC is assigned a link-local and at least one global address. The following command displays the configuration of network interface `eth1` on a Red Hat system. This interface is associated with two IPv6 addresses. One of them starts with `fe80::` and the other with `3ffe::`. The scope of the first one is the link and the second has a global scope.


```

root# ip address list eth1
3: eth0: <BROADCAST,MULTICAST,UP mtu 1500 qdisc pfifo_fast qlen 100
link/ether 00:0c:29:c2:52:ff brd ff:ff:ff:ff:ff:ff
inet6 fe80::20c:29ff:fec2:52ff/10 scope link
inet6 3ffe:1200:4260:f:20c:29ff:fec2:52ff/64 scope global

```

Creation of the Link-Local Address

An IPv6 address is 128 bits long. It has two parts: a *subnet prefix* representing the network to which the interface is connected and a *local identifier*, sometime called token. In the simple case of an Ethernet medium, this identifier is usually derived from the EUI-48 *Media Access Control* (MAC) address using an algorithm described later in this article. The subnet prefix is a fixed 64-bit length for all current definitions. Because IPv4 manual configuration is a well-known pain, one could hardly imagine manipulating IPv6 addresses that are four times longer. Moreover, a DHCP server is not always necessary or desired; in the case of a remote control finding the DVD player, a DHCP environment is not always suitable.

Because the prefix length is fixed and well-known, during the initialization phase of IPv6 NICs, the system builds automatically a link-local address. After a uniqueness verification, this system can communicate with other IPv6 hosts on that link without any other manual operation.

For a system connected to an Ethernet link, the build and the validation of the link-local address is the following:

1. An identifier is generated, supposedly unique on the link.
2. A tentative address is built.
3. The uniqueness of this address on the link is verified.
4. If unique, the address from phase 2 is assigned to the interface. If not unique, a manual operation is necessary.

Although a local policy can decide to use a specific token, the most common method to obtain a unique identifier on an Ethernet link is by using the EUI-48 MAC address and applying the modified IEEE EUI-64 standard algorithm. A MAC address (IEEE 802) is 48 bits long. The space for the local identifier in an IPv6 address is 64 bits. The EUI-64 standard explains how to stretch IEEE 802 addresses from 48 to 64 bits, by inserting the 16 bits **0xFFFE** at the 24th bit of the IEEE 802.

By doing so, transforming MAC address **00-0C-29-C2-52-FF** using the EUI-64 standards leads to **00-0C-29-FF-FE-C2-52-FF**. Using IPv6 notation, we get **000C:29FF:FEC2:52FF**. Recall that the notation of IPv6 addresses requires 16-bit pieces to be separated by the character “:”. Then, it is necessary (RFC 3513) to invert the universal bit (“u” bit) in the 6th position of the first octet. Thus the result is:
020c:29ff:fec2:52ff.

Universal uniqueness of IEEE 802 and EUI-64 is given by a “u” bit set to 0. This global uniqueness is assured by IEEE, which delivers those addresses for the entire planet. Inverting the “u” bit allows ignoring it for short values in the manual configuration case, as explained in paragraph 2.5.1 of RFC 3513^[4].

The second phase of creating automatically a link-local address is to prepend the well-known prefix **fe80::/64** to the identifier resulting from phase one. In our case we obtain **fe80::20c:29ff:fec2:52ff**. This address is associated with the interface and tagged “tentative.” Before final association, it is necessary to verify its uniqueness on the link. The probability of having a duplicate address on the same link is not null, because it is recognized that some vendors have shipped batches of cards with the same MAC addresses.

This is the goal of the third phase, called *Duplicate Address Detection* (DAD). The system sends ICMPv6 packets on the link where this detection has to occur. Those packets contain *Neighbor Solicitation* messages. Their source address is the undefined address “::” and the target address is the tentative address. A node already using this tentative address replies with a *Neighbor Advertisement* message. In that case, the address cannot be assigned to the interface. If there is no response, it is assumed that the address is unique and can be assigned to the interface.

We are reaching the last step of the automatic generation of a link-local address. This phase removes the “tentative” tag and formally assigns the address to the network interface. The system can now communicate with its neighbors on the link.

Global Prefixes

In order to exchange information with arbitrary systems on the global Internet, it is necessary to obtain a global prefix. Usually (but not necessarily), the identifier built during the first step of the automatic link-local autoconfiguration process is appended to this global prefix.

However, before assigning this global address, the system verifies again that no duplicate address exists on the link. DAD is performed for all addresses before they are assigned to an interface, because uniqueness in one prefix does not automatically assure uniqueness in any other available prefixes.

Generally, global prefixes are distributed to the companies or to end users by *Internet Service Providers* (ISPs).

Random Identifiers

The EUI-48-to-EUI-64 transform process is attractive because it is simple to implement. However, it generates a privacy problem. Global unicast as well as link-local addresses may be built with an identifier derived from the MAC address. A Website tracking where a node frequently attaches can collect private information such as the time spent by employees in the enterprise or at home.

Because a MAC address follows the interface it is attached to, the identifier of an IPv6 address does not change with the physical location of the Internet connection. Hence it is possible to trace the movements of a portable laptop or *Personal Digital Assistant* (PDA) or other mobile IPv6 device.

RFC 3041^[5] allows the generation of a *random* identifier with a limited lifetime. Because IPv6 architecture permits multiple suffixes per interface, a single network interface is assigned two global addresses, one derived from the MAC address and one from a random identifier. A typical policy for use of these two addresses would be to keep the MAC-derived global address for inbound connections and the random address for outbound connections. A reason for not using it for inbound connections is the need to update the DNS just as frequently as it is changes.

Such a system, with two different global addresses—one of which changes regularly—becomes very difficult to trace.

By default, Microsoft enables this feature on Windows XP and Windows Server 2003. The random-identifier-based global addresses of Microsoft systems have the address type “temporary.” EUI-64 global addresses have type “public.” Those types as well as other information can be displayed in a **cmd.exe** DOS-box with the command line:

```
netsh interface ipv6 show address
```

IPv6 Routers

By definition, a router is a node that forwards IP packets not explicitly addressed to it. IPv6 routers are certainly compliant with this definition but, in addition, they regularly advertise information on the links to which they are connected—provided they are configured to do so. These advertisements are *Internet Control Message Protocol Version 6* (ICMPv6) *Router Advertisement* (RA) messages, sent to the multicast group **ff02::1**. All the systems on a link must belong to this group, and nodes configured for autoconfiguration, among other things, analyze the option(s) of those messages. They might contain any routing prefix(es) for this segment.

Router Solicitation

Upon reception of one of those RA messages and according to local algorithm policy, an autoconfiguring node not already configured with the corresponding global address will prepend the advertised prefix to the unique identifier built previously.

However, the advertisement frequency, which is usually about ten seconds or more, may seem too long for the end user. In order to reduce this potential wait time, nodes can send *Router Solicitation* (RS) messages to all the routers on the link. Nodes that have not configured an address yet use the unspecified address “::”. In response, the routers must answer immediately with a RA message containing a global prefix. This router solicitation corresponds to ICMPv6 messages of type RS, sent to the all-router multicast group: **ff02::2**. All routers on the link must join this group.

Thus, a node soliciting on-link routers in such a way is able to extract a prefix and build its global address. Note that this method using an advertised prefix is possible only for end nodes. Today IPv6 routers are usually manually configured. The reason is obvious: a stateless automatic configuration requires the advertisement of a prefix. This prefix is sent by a router. The router sending the prefix must be fully configured to do so. The easiest way to break this seemingly unsolvable problem is to manually configure IPv6 routers. However, some automatic methods are being developed^[6].

Conclusion

Stateless address autoconfiguration is a new concept with IPv6. It gives an intermediate alternative between a purely manual configuration and stateful autoconfiguration. In addition to ease of use with no dedicated server or relay, this mechanism removes problems that have not been discussed here, such as the mismatch between the DHCP server and the router (prefix topology) or the IPv4 need to readdress subnets that have outgrown their prefix. Moreover, automatic renumbering (prefix change) is also possible on nodes using stateless autoconfiguration.

References

RFCs can be found at <http://www.ietf.org/rfc/>

- [1] Droms, R., "Dynamic Host Configuration Protocol," RFC 1531, October 1993.
- [2] Droms, R., Ed., Bound, J., Volz, B., Lemon, T., Perkins, C., Carney, M., "Dynamic Host Configuration Protocol for IPv6 (DHCPv6)," RFC 3315, July 2003.
- [3] Thomson, S., Narten, T., "IPv6 Stateless Address Autoconfiguration," RFC 2462, December 1998.
- [4] Hinden, R., Deering, S., "Internet Protocol Version 6 (IPv6) Addressing Architecture," RFC 3513, April 2003.
- [5] Narten, T., Draves, R., "Privacy Extensions for Stateless Address Autoconfiguration in IPv6," RFC 3041, January 2001.
- [6] Prefix delegation:
<http://www.ietf.org/internet-drafts/draft-ietf-dhc-dhcpv6-opt-prefix-delegation-06.txt>

FRANÇOIS DONZÉ studied at the University of Utah in Salt Lake City. In 1989 he joined Digital Equipment Corporation as a UNIX and network teacher. He is now a technical consultant at HP, based in Sophia-Antipolis, France, promoting IPv6 and other leading-edge technologies. The author of several internal articles, he also publishes in French magazines. E-mail: francois.donze@hp.com

DNSSEC: The Protocol, Deployment, and a Bit of Development

by Miek Gieben, NLnet Labs

“One Key to rule them all,
one Key to find them,
one Key to bring them all
and in the Resolver bind them.”

—Modified from *Lord of the Rings*.

The *Domain Name System* (DNS) (RFCs 1034 and 1035) is a highly successful and critical part of the Internet infrastructure. Without it the Internet would not function. It is a globally distributed database, whose performance critically depends on the use of caching.

Unfortunately the current DNS is vulnerable to so-called *spoofing attacks* whereby an attacker can fool a cache into accepting false DNS data. Also various man-in-the-middle attacks are possible. The *Domain Name System Security Extension* (DNSSEC) is not designed to end these attacks, but to make them detectable by the end user. Or more technically correct: detectable by the security-aware *resolver* doing the work for the end user. This saves users from doing online banking on the wrong server even if a secured connection is used and the address in the browser looks correct.

DNSSEC is about protecting the end user from DNS protocol attacks. In order to make it work, zone owners (such as *.com*, *.net*, *.nl*, etc.) need to deploy DNSSEC in their zones. End users then need to update their resolvers to become security-aware (that is, understand DNSSEC) and add some trusted keys. These keys are called *anchored keys*; they are configured in the resolver and cannot be changed or updated very easily. If this is all configured, the end user will (finally) be able to detect attacks.

DNSSEC, as defined in (hopefully soon-to-be-obsolete) RFC 2535, adds data origin authentication and data integrity protection to the DNS. The *Public Key Infrastructure* (PKI) in DNSSEC may be used as a means of public key distribution, which may be used by other protocols. *IP Security* (IPSec) and the *Secure Shell* (SSH) protocol, for example, are already considering the use of DNSSEC to carry their keying material.

In the course of early-deployment experiments carried out by various organizations, it became evident that RFC 2535 introduced an administrative key-handling and maintenance nightmare. This in turn would mean the DNSSEC deployment would never start (or be successful, for that matter).

The IETF DNSEXT working group decided to fix this problem, and to incorporate all drafts and RFCs written since RFC 2535 into a new DNSSEC specification.

This (still ongoing) effort became known as the *RFC 2535bis* DNSSEC specification. This work has resulted in three drafts, each handling a specific part of the new specification. These drafts follow:

1. dnssec-intro^[1] provides an introduction into DNSSEC.
2. dnssec-records^[2] introduces the new records for use in DNSSEC.
3. dnssec-protocol^[3] is the main document, which details all the protocol changes.

The documents are now almost ready (July 2004) to be submitted to the *Internet Engineering Steering Group* (IESG) for review. It is hoped that soon after this is done the drafts will become RFCs. It could be that 2004 will be the year of DNSSEC.

In this article I use the terms *domain* and *zone*. These are important concepts in the DNS and in DNSSEC. The difference between a zone and a domain is worth highlighting. A domain is a part of the DNS tree. A zone contains the domain names and data that that domain contains *except* for the domain names and data that are delegated elsewhere. Also refer to [4].

Consider, for instance, the *.com domain*, which includes everything that ends in *.com*. *CNN.com* is in the *.com domain*. The *.com zone*, however, is the entity handled by VeriSign.

One other important concept in DNS is the *Resource Record* (RR) and the *Resource Record Set* (RRset). An RR in DNS is, for instance:

```
www.example.org.  IN   A    127.0.0.1
```

... where *www.example.org* is the “ownername” or “name.” *IN* is the class (IN stands for Internet). *A 127.0.0.1* is the type (together with its rdata). *A* stands for “address.” This 3-tuple (name, class, type) together make up the resource record. RRset are all the RRs that have an identical name, class and type. Only the rdata is different. Thus:

```
www.example.org.  IN   A    127.0.0.1
www.example.org.  IN   A    192.168.0.1
```

... together form a RRset, but:

```
www.example.org.  IN   A    127.0.0.1
www.example.org.  IN   MX   mail.example.org.
```

... do not (their type is different). In the DNS an RRset is considered *atomic* and the smallest data item. In DNSSEC each RRset gets a signature.

What Is DNSSEC?

DNSSEC adds data origin authentication and data integrity to the DNS. To achieve this, DNSSEC uses public key cryptography; (almost) everything in DNSSEC is digitally signed.

Public key cryptography uses a single key split in two parts: a private and a public component. The *private* component, also known as the *private key*, must be kept secret. The *public* component (the *public key*) can be made public. Both these keys can be used for cryptographic operations, albeit with different goals.

If a message is scrambled with the public key, it can be decrypted only with the private key. This is called *encryption* of the message and it ensures that only the holder of the private key can read the original message. When the private key is used to scramble a message, everybody can use the available public key to decipher the message. This last operation is called (digitally) *signing* a message (for increased speed usually a hash of the message is signed). In this case you know where the message comes from (*authenticated data origin* in cryptographic jargon). An added benefit of signing messages is that when the data is mangled during transport the signature is no longer valid. This last property is called *authenticated data integrity*. A more lengthy introduction on public key cryptography can be found at [10]. In DNSSEC only digital signatures (signing) are used, and nothing is ever encrypted.

For every secure zone there must be a public key in the DNS for use by DNSSEC. Each zone administrator generates a key to be used for securing a zone. The private key is (of course) kept private and is used in the “signing process” to create the signatures. The public key is published in DNSSEC as a DNSKEY record, which is the zone key. The generated signatures are published as RRSIG records.

If RRsets in DNSSEC do not have a valid signature, they are labeled bogus by the resolver. Bogus data should not be trusted, because probably somebody is trying to conduct a spoof attack. DNSSEC further distinguishes between:

- Verifiable secure—The data has signatures that are valid.
- Verifiable unsecure*—The data has no signatures.
- Old-style DNS—A non-DNSSEC lookup is done.

* Yes, Unsecure. This word has somehow evolved from “insecure.”

Verifiable secure data is data that has valid signatures, and the key used to create those signatures is trusted (anchored in the resolver). Verifiable unsecure data is data for which we know for sure we do not need to do signature validation. Old-style DNS is the current (insecure) method of getting DNS data.

The signing of data in DNSSEC is comparable to the *Gnu Privacy Guard* (GPG) signing of e-mail. If I trust a public key from someone, I can use that key to verify the GPG signature and authenticate the origin of the e-mail.

The problem with both DNSSEC and GPG lies in the “...If I trust the public key from someone.” GPG solves this with public key servers, key signing parties at various events and thus the creation of a web of trust. For DNSSEC such solutions are impractical. DNSSEC uses a different, but very elegant mechanism called the *chain of trust*.

The chain of trust makes it possible to start with a root zone key, the highest possible key in the DNS tree, and following cryptographic pointers to lower zones. Each pointer is validated with the previous validated zone key. (The root key is the key used in the root zone of the Internet; it is the key used in the . (dot) zone. It could take a while before the root is signed.)

By using this mechanism only the root key is needed to validate *all* DNSSEC keys on the Internet. With these DNSSEC keys the DNS data in each zone can then be validated. So, unlike GPG, we need to distribute only one key. This can be done by publishing it on the World Wide Web or in a newspaper or putting an ad on TV, etc.

One of the current items in the DNSSEC community is to outline procedures and guidelines on how to update this root and other keys.

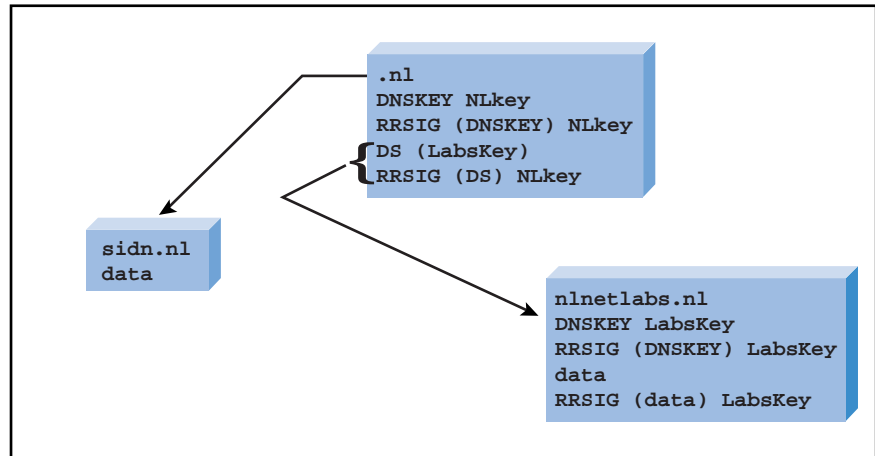
Chain of Trust

To start securely resolving in DNSSEC, a root key must be anchored in the resolver at your local computer or nameserver. Only when a resolver knows and trusts a zone key can it validate the signatures belonging to that zone. Because of the chain of trust, a resolver has to carry only a few zone keys to be able to validate DNSSEC data on the Internet.

The chain of trust works by following “secured pointers,” which are called *secured delegation* in DNSSEC. A special, new record called the *Delegation Signer* (DS) record delegates trust from a parental key to a child’s zone key.

The DS record holds a hash (*Secure Hash Algorithm 1* [SHA-1]) of a child’s zone key. This DS record is signed with the zone key from the parent. By checking the signature of the DS record, a resolver can validate the hash of the child’s zone key. If this is successful, the resolver can compare this (validated) hash with the (yet-to-be-validated) hash of the child’s zone key. If these two hashes match, the child’s real zone key can be used for validation of data in the child’s zone. Note: by successfully following a secured delegation, the amount of trust a resolver has in the parental key is transferred to a child’s key. This is the crux of the chain of trust.

Figure 1: **nlnetlabs.nl** is a secured delegation under **.nl**.
RTSIG(x)y denotes that a signature over a data *x* is created with key *y*.



In Figure 1 the following takes place.

The **.nl** zone contains the following:

```

nl.      IN      SOA (soa-parameters)
; the zone key
nl.      IN      DNSKEY NLkey
nl.      IN      RRSIG(DNSKEY)NLkey
nl.      IN      RRSIG(SOA)NLkey

nl.      IN      NS ns5.domain-registry.nl.
; this NS is authoratitive
nl.      IN      RRSIG(NS) NLkey

nlnetlabs.nl.  IN  NS open.nlnetlabs.nl.
; no RRSIG here (nonauthoritative data is not signed)

; DS record with a hash of the child's zone key
nlnetlabs.nl.  DS  hash(LabsKey)
; The signature of the parent
nlnetlabs.nl.  RRSIG(DS)NLkey
  
```

Note: It is important to see that we now have linked a parental signature to something that is *almost* the key of the child.

And the **nlnetlabs.nl** zone has the following:

```

nlnetlabs.nl.  IN      SOA (soa-parameters)
; The zone key
nlnetlabs.nl.  IN      DNSKEY LabsKey
nlnetlabs.nl.  IN      RRSIG(SOA)Labskey
; The (self) signature of the zone key
nlnetlabs.nl.  IN      RRSIG(DNSKEY)Labskey
nlnetlabs.nl.  IN      NS open.nlnetlabs.nl.
nlnetlabs.nl.  IN      RRSIG(NS)LabsKey
  
```

So the chain of trust looks like the following:

```
.nl DNSKEY -> nlnetlabs.nl DS -> nlnetlabs.nl DNSKEY
```

... and with that last key we can validate the data in the **nlnetlabs.nl** zone.

With this “trick” all keys from all the secure **.nl** zones can be chained from the **.nl** “master” key. So instead of one million (the number of zones in **.nl** currently) we need to configure only one key.

As you might have guessed, getting the root zone signed as soon as possible will make it possible to have one key that validates all other keys on the Internet.

We can also look at it from the resolver side. A resolver wants to get an answer. With DNSSEC it has to deal with signatures, keys, and DS records, but those are “side issues”; it still wants an answer.

Suppose **.nl** is secured and a secure delegation to **nlnetlabs.nl** exists. Our resolver has the key of **.nl** anchored. The nameservers of the root zone are also known to the resolver. We further assume the root is not signed. The resolver wants to resolve the address (A record) of **www.nlnetlabs.nl**. What does the actual resolving process look like in DNSSEC? Numerous steps need to be performed:

1. Go to a root server and ask our question.
2. The root server does not know anything about **www.nlnetlabs.nl**, but it *does* know something about **.nl**. The root nameserver refers us to the **.nl** nameservers. This kind of answer is called a *referral*.
- 3a. Notice that we have a key for **.nl** anchored.
- 3b. Go to the **.nl** nameserver and ask the **.nl** DNSKEY.
- 4a. Compare the two DNSKEYs. Continue with the secure lookup only if they match.
The **.nl** DNSKEY is now validated.
- 4b. Optionally, the RRSIG on the DNSKEY also can be checked.
5. Ask a **.nl** nameserver our question.
6. The **.nl** nameserver is also oblivious about **www.nlnetlabs.nl**, but it does know something about **nlnetlabs.nl**. It returns a secure referral consisting of a DS record plus the RRSIG and some nameservers.
7. The resolver now checks the signature on the DS record. If the signature is valid, the hash of the **nlnetlabs.nl** zone key is ok. The nameservers in the referral do not have any signatures on them.
The hash of the **nlnetlabs.nl** DNSKEY is validated with the **.nl** DNSKEY.
8. Go to the nameserver as specified in the referral and ask for the **nlnetlabs.nl** DNSKEY.
9. Hash the DNSKEY of **nlnetlabs.nl** and compare this hash with the hash in the DS record. If they match continue with the secure lookup.
The **nlnetlabs.nl** DNSKEY is now validated.
10. Ask the nameserver of **nlnetlabs.nl** our question.

11. The nameserver now responds with an answer consisting of the A record of **www.nlnetlabs.nl** and an RRSIG made with the **nlnetlabs.nl** DNSKEY.
12. The resolver now uses the already validated **nlnetlabs.nl** DNSKEY to check the RRSIG. If that signature is valid the RR with the answer is ok and can be given to the application.
13. After these steps we find out that the address of **www.nlnetlabs.nl** is 213.154.224.1. We also know it is not a spoofed answer.

This looks like a lot of work and it is—a recursive resolver is a complicated piece of software. Keep in mind, though, that only steps 3ab, 4ab, 7, 8, 9, and 12 are needed for DNSSEC; the rest is how resolving is done in the DNS today.

Deployment

As mentioned earlier, each zone owner generates its own key. To make the secure delegation actually work, this key must somehow be securely transferred to the parent, which is usually the local registry. The registry must have procedures in place to determine whether or not the uploaded key really belongs to the domain it claims to come from. During the *Secure Registry* (SECREG) experiment^[5] NLnet Labs has researched the impact DNSSEC has on registries.

But even before the key can be actually uploaded to the parent, a zone administrator still has to do some work; the DNS zone must be signed. This process, called *zone signing*, turns a DNS zone into a DNSSEC zone.

The signing is done offline; first you sign, and then you load the zone. This setup was chosen because at the time (late 1990) computers were not fast enough to generate the signature in real time. Currently it would be possible to do this, but having a server sign every answer it gives is a *Denial-of-Service* (DoS) attack waiting to happen. Especially root servers will be unable to do this.

In DNSSEC a zone can have multiple keys. The signed zone then has multiple signatures per RRset (one for each key). There is no protocol limit on the number of keys. Here we sign with only one zone key. Also signatures in DNSSEC have a start and end date, that is, before and after a certain date interval the signature can no longer be used for validation.

If you use DNSSEC, you must re-sign your zone to generate new signatures with a new validity interval.

The signing of a zone consists of the following steps:

1. The zone key is added to the zone file.
2. The zone file is sorted.

3. Each owner name (for example, a host name) in the zone gets a *Next SEcure* (NSEC) record. (Refer to the section “Authenticated Denial of Existence.”)
4. For each secured delegation, a DS record is added.
5. The entire zone is then signed with the private key of the zone. Each authoritative RRset gets a signature, including the newly generated NSEC records.

Berkeley Internet Name Domain (BIND)^[6] version 9—a popular implementation of the DNS protocols—contains a tool *dnssec-signzone*, which does steps 2 through 5 automatically; we only (manually) need to add the zone key to the zone file. The net result is that we have a bigger, signed, DNSSEC zone. A typical DNSSEC zone is 7 to 10 times larger than its DNS equivalent.

Experiments have shown that this does not pose much of a problem, even for such so-called country code *Top Level Domains* (ccTLDs) as **.nl**. The signed **.nl** zone was 350 megabytes, slightly more than a half a CD-ROM. And even if scaling problems are occurring, 64-bit machines would certainly help.

A few years ago there was much concern about the signing time. There was fear that it would be impossible to sign large zones, such as **.com**.

Experiments disproved this fear. Furthermore, a zone can be split up in pieces and each piece can be signed on a different machine. Later all the signed pieces can be put back together. Signing DNS zones is a highly parallel process.

After signing the zone, it can be loaded in the nameserver. If a resolver is DNSSEC-aware and has been configured with a trusted key that has a chain of trust to the zone key, it can validate the answers. If an answer does not validate, something is wrong and the DNS data must not be used.

The actual Internet-wide deployment of DNSSEC can happen incrementally. Each zone can decide to join independently. It is expected that initially DNSSEC is deployed in subsections of the Internet. These so-called *Islands of Trust* can appear anywhere on the Internet or even in intranets. The only requirement is that the key of the island of trust is distributed to the resolver. Resolvers configured with the key of a certain island of trust are called the *resolvers of interest*. Of course when DNSSEC is widely deployed on the Internet all resolvers are resolvers of interest and will have that key preconfigured.

Authenticated Denial of Existence

As mentioned previously, all records are signed offline. When a nameserver receives a query it looks up the answer plus the signature and returns the two (RRSIG + RRset) to the resolver. The signature is thus not created in real time. How can a secure-aware nameserver then respond to a query for something it does not know (that is, give an NX-DOMAIN answer)? The only way to have offline signing and NXDOMAIN answers work together is to somehow sign the data you do not have.

In DNSSEC this is accomplished by the *Next SECure* (NSEC) record. This NSEC record holds information about the next record; it spans the nonexistence gaps in a zone, so to say. For this to work, a DNSSEC zone must be sorted (this is where that requirement stems from). To clarify this, consider an example.

We have a DNS zone, with (for the sake of clarity only the NSEC records are shown):

```
a.nl  
d.nl  
e.nl
```

Next we generate (with the signer) our DNSSEC zone:

```
a.nl  
a.nl NSEC d.nl (span from a.nl to d.nl)  
  
d.nl  
d.nl NSEC e.nl (span from d.nl to e.nl)  
  
e.nl  
e.nl NSEC a.nl (loop back to a.nl)
```

1. If a resolver asks information about `b.nl`, the nameserver tries to look up the record fails. Instead it finds `a.nl`. It must then return: `a.nl NSEC d.nl` together with the signature. The resolver must then be smart enough to process this information and conclude that `b.nl` does not exist. If the signature is valid, we have an *authenticated denial of existence*. These NSEC records together with their signatures are the major cause of the zone size increase in DNSSEC.

Road to the DS Record

This section briefly considers the history of DNSSEC and, in particular, why the DNSEXT working group has invented this peculiar DS record, which can only exist at the parent side of a zone cut.

In RFC 2535 the DS record did not exist, and this is the reason that the key management in RFC 2535-DNSSEC is very, very cumbersome. In 2000 NLnet Labs ran its first experiment to test deployment of DNSSEC in the Netherlands. Because `.nl.nl` was chosen as the zone under which the secure tree would grow, this experiment became known as the *nl-nl-experiment*. With this experiment it was shown that the current DNSSEC standard (the soon-to-be-obsolete RFC 2535) was difficult to deploy^[7].

An update of a zone key in a child zone required up to 11 (coordinated and sequential) steps with the parent zone. The `.nl` zone now has more than 1 million delegations, so updating all the child zones would require more than 11 million steps. Because these updates could be quite frequent (once a month is typical), this is clearly an administrative nightmare.

Worse yet, if **.nl** lost its private key, all child-zone administrators would have to be notified and they would have to resubmit their public key for re-signing with the new **.nl** key. And because under these conditions the DNS may have been hacked and is thus untrusted, **.nl** is limited in its communication through the Internet; e-mail may not be the preferred method. A telephone call would be more safe, but what kind of organization can make up to one million phone calls in a few days ..?

After various failed attempts (sig@parent^[8]) to fix this behavior, the DS record was introduced^[1,3]. With this record the administration nightmare is solved, because DS introduces an indirection from the parent zone to a child's zone key.

If **.nl** loses its private key, it can easily resign its own zone, *without* contacting all its children. The DS to child key indirection is still valid, and only the signature of the DS record needs to be updated. This is a local operation.

To test this new DNSSEC specification, a new experiment was set up, which would build a shadow DNSSEC tree in the **.nl** zone. This experiment, called *SECREG*, was to test the new procedures in DNSSEC and, of course, the new DS record. Detailing the conclusions of this experiment is beyond the scope of this article, but in short the conclusion was that the new DNSSEC procedures do not pose much difficulty. At some point, more than 15,000 zones were delegated from the secure tree. A writeup of the experiment and the conclusions can be found in "DNSSEC in NL"^[5].

Settings and Parameters in DNSSEC

DNSSEC brings many new parameters to the DNS, including cryptographic ones such as key sizes, algorithm choices, and key and signature lifetimes. Because DNS never has involved cryptography, the best values for these parameters are still open for debate. There is, however, some documentation and knowledge available on this topic (refer to [9] for instance).

One of the major issues is how large (bit length) to make a zone key and how often to re-sign a zone file. The current view is that a parent zone should use larger keys and re-sign more often than a child zone. Also the signature lifetime should be shorter in a parent zone.

Because a parent zone has a DS record (and signature) of a child's zone key, it can decide how long this DS RRSIG must be valid. The shorter this validity interval is, the better protected the child. If a cracker steals a child's zone key, it can forge DNS data. This data looks genuine because the cracker has access to the private key. As long as there is a valid chain of trust to this hijacked key, the child is vulnerable. This chain of trust is broken as soon as the RRSIG of the DS record expires. This argues in favor of a very short parental RRSIG over the DS record.

However, making this interval too short opens the door for accidental mishaps. If a child zone makes an error and somehow the chain of trust is broken, it has until the RRSIG expires to fix the problem. This would recommend a longer signature lifetime. In DNSSEC these and other trade-offs have to be made.

The IETF DNSOP working group is currently addressing these parameters and their trade-offs. The current data came (and comes) from workshops and early test deployments.

Outlook and Prospects

Because DNSSEC requires some additions to the (cc/g)TLD registration process, it could be a while before ccTLDs are capable of deploying DNSSEC. If the protocol is completed this year (2004), it will probably take a few years before registries can advertise DNSSEC domain names.

It is important to consider what DNSSEC actually wants to accomplish; it makes spoofing attacks in the DNS visible—and nothing more. It is not a PKI with all the extra features because key revocation is, for instance, not implemented in DNSSEC. Seen in this light, the protection of private keys in DNSSEC is important, but when a private key is compromised we are just back to plain old DNS.

On the other hand, because DNSSEC does introduce cryptographic material in the DNS and allows for the addition of other (non-DNS) keys, some interesting possibilities emerge. Many technologies on the Internet want to have some kind of simple key distribution mechanism in place; for example: SSH and IPSec. What DNSSEC promises is a system in which we can validate the SSH key from an unknown host with only one key. If the validation is successful, we are quite certain the SSH host key comes from the host from which it claims to come. We get this without any extra effort or cost (from a client's perspective at least). The possibilities are probably endless.

References

- [1] Roy Arends, Rob Austein, Dan Massey, Matt Larson, and Scott Rose, "DNS Security Introduction and Requirements," Work In Progress,
<http://www.ietf.org/internet-drafts/draft-ietf-dnsext-dnssec-intro-10.txt>
- [2] Roy Arends, Rob Austein, Dan Massey, Matt Larson, and Scott Rose, "Resource Records for the DNS Security Extensions," Work In Progress,
<http://www.ietf.org/internet-drafts/draft-ietf-dnsext-dnssec-records-08.txt>
- [3] Roy Arends, Rob Austein, Dan Massey, Matt Larson, and Scott Rose, "Protocol Modifications for the DNS Security Extensions," Work In Progress,
<http://www.ietf.org/internet-drafts/draft-ietf-dnsext-dnssec-protocol-06.txt>

- [4] DNS and BIND Talk Notes:
<http://www.tfug.org/helpdesk/general/dnsnotes.html>
- [5] R. Gieben, "DNSSEC in NL,"
<http://www.miek.nl/publications/dnssecnl/index.html>
- [6] BIND9, Berkeley Internet Name Domain, Version 9:
<http://www.isc.org/sw/bind/>
- [7] R. Gieben, "Chain of Trust: The parent-child and keyholder-keysigner relations and their communication in DNSSEC," NIII report CSI-R0111:
<http://www.cs.kun.nl/research/reports/info/CSI-R0111.html>
<http://www.miek.nl/publications/thesis/CSI-report.ps>
- [8] R. Gieben and T. Lindgreen, "Parent's SIG over Child's KEY,"
<http://www.nlnetlabs.nl/dnssec/dnssec-parent-sig-01.txt>
- [9] O. Kolkman and R. Gieben, "DNSSEC Operational Practices," Work In Progress,
<http://www.ietf.org/internet-drafts/draft-ietf-dnsop-dnssec-operational-practices-01.txt>
- [10] Netscape Communications Corporation, "Introduction to Public-Key Cryptography,"
<http://developer.netscape.com/docs/manuals/security/pkin/contents.htm>

MIEK GIEBEN graduated in Computer Science in 2001 from the University of Nijmegen (Netherlands) on the subject of DNSSEC. He has been employed by NLnet Labs since that time. He has been using Linux and the Internet since 1995. Currently he is involved in DNSSEC deployment and has co-written parts of NSD2 (which is now fully DNSSEC aware). His personal home page can be found at <http://www.miek.nl/>. The home page of NLnet Labs can be found at <http://www.nlnetlabs.nl/>.
E-mail: miekg@atoom.net

Book Review

Network Management *Network Management, MIBs and MPLS* by Stephen B. Morris, ISBN 0131011138, Prentice Hall, June 2003.

Few people would question the need for good network management, and books about the *Simple Network Management Protocol* (SNMP) have been circulating for more than ten years now. But the key differentiator of this book is well recognized in its title—it's about SNMP in the context of a *Multiprotocol Label Switching* (MPLS) network. MPLS is now recognized as the convergence technology, and an increasing number of mission-critical services are being deployed over it. World-class network management is vital to keep these services running to the “five nines” level we've all come to expect.

Organization

In this book, Stephen Morris offers a very approachable and comprehensive look at SNMP and the methodology behind the all-important *Management Information Base* (MIB). The first chapter gives the obligatory justification for network management and sets the scene nicely for the rest of the book.

It's amazing to think that SNMP has been around since the late 1980s, and yet if you ask any MPLS operations person, the odds are that person is still using a *Command-Line Interface* (CLI) to actually configure boxes. CLI is a man-machine interface, not a machine-machine interface like SNMP. Even centralized provisioning platforms, such as the former Orchestream (now Metasolve) VPN Manager, simply created a friendly *Graphical User Interface* (GUI) front end for the provisioning procedure, and then ran CLI scripts frantically in the background. The drawbacks of CLI configuration are too numerous to list here, but the basic solution to the problem is to create a scalable and secure machine-to-machine interface. In the IP world the candidate technology for this is SNMPv3, and Morris discusses both the MIB structure (the key to scalability) and the security model in Chapter 2. Because premium MPLS-based services demand secure and robust provisioning, SNMPv3 is the technology of choice.

Chapter 3 describes what Morris calls the “Network Management Problem,” although in fact this is described as a whole set of problems, some of which are caused by deficiencies in the SNMP architecture, whereas others are caused by the scale and pace of operations in a modern network. A specific problem that Morris addresses very sensibly is the way that the rapid pace of network technology development impacts the ability to manage these networks. In other words, new technologies tend to appear too quickly for management mechanisms to be optimized for these protocols. To solve this problem, Morris (a software engineer by training) presents a series of “Linked Overviews” (these describe the properties of a given network technology—MPLS, *Asynchronous Transfer Mode* (ATM), etc.—in a procedural framework. In essence this is a kind of recipe for the software developer. In addition, the text is liberally sprinkled with “Developers Notes” that I'm sure will provide invaluable help for people trying to write management system code.

Chapter 4 then takes the approach of solving the “Network Management Problem” to a higher, and perhaps longer-term level, with the proposed development of smarter network management components and more integrated data frameworks. This culminates in a description of *Directory Enabled Networking*, a technology that seemed to flower briefly in the context of network management a few years ago, but then was buried when the telecom recession hit the industry. My own feeling is that the time is right for a rebirth of this approach in modern, converged networks.

Chapter 5 looks at some real *Network Management System* (NMS) issues, using the HP OpenView Network Node Manager as a worked example. Morris is quick to point out that this is not an endorsement of the product, but because it is the most well-known and widely used product in this class, it is the logical choice.

Chapters 6 and 7 look at software components, and Morris’s background in software development shines through here in the level of detail, coupled with well-structured explanations.

Chapter 8 describes a very useful case study of using SNMP to provision a tunnel through an MPLS network—a task that is typically performed today using crude CLI techniques.

Chapter 9 contrasts theory and practice in network management, and deals with the loose ends of various topics such as end-to-end security and the integration of a third-party *Open Source Software* (OSS) using standardized northbound *Element Management System* (EMS) interfaces.

Recommended

Overall this is an excellent book that really does deliver what it claims—a comprehensive and practical look at the latest SNMP technologies and techniques. In this regard it stays highly focused, and doesn’t waste time with irrelevant discussion on other topics. For example, at first I was disappointed to note that only a page or two of brief explanation is devoted to topics such as *Common Object Request Broker Architecture* (CORBA) and *Extensible Markup Language* (XML). But in the context of what this book is trying to tell us, it makes perfect sense. Each of these topics really needs its own book to cover the topic in similar detail to Morris’s work.

Similarly, if you’re expecting a description of emerging IP/MPLS *Operations, Administration, and Maintenance* (OA&M), then this book is not for you. Again, I would defend Morris’s use of Occam’s Razor because OA&M protocols are usually demanded by network staff, and not by OSS operatives. In my own opinion, this situation will gradually change in the next few years, as OA&M is recognized as the “eyes and ears” of the OSS. Perhaps this would be a good place for Mr. Morris to start his next book.

—Geoff Bennett, *Heavy Reading*
bennett@heavyreading.com

Cooperative Support for Global IPv6 Deployment

The *Regional Internet Registries* (RIRs), the *IPv6 Task Forces* and the *IPv6 Forum* are working in cooperation to support global IPv6 deployment.

The four RIRs, APNIC, ARIN, LACNIC and the RIPE NCC, are responsible for the management of global Internet numbering resources, including IPv4 and IPv6 address space, throughout the world. The RIRs confirm their commitment and continued support towards the deployment of IPv6 in cooperation with the IPv6 Task Forces and with the support of the IPv6 Forum.

The IPv6 Task Forces are focused on rapid IPv6 deployment. They see the adoption of IPv6 by industry, governments, schools and universities is particularly important. The extra address space offered by IPv6 will facilitate the deployment of widespread “always-on” Internet services including broadband access for all. In addition, IPv6’s built-in encryption will help improve Internet security and is promoted by many government institutions globally.

The cooperation among the RIRs and the IPv6 Task Forces includes key aspects such as:

- Supporting awareness, education and deployment of IPv6;
- Disseminating information on the progress of IPv6 deployment;
- Encouraging dialogue and ensuring the necessary cooperation between all involved parties;
- Benchmarking IPv6 deployment progress;
- Supporting the adoption of Domain Name Service infrastructure necessary for IPv6;
- Encouraging the participation of all those who are interested in the IPv6 policy development process.

This cooperative effort between the RIRs and the IPv6 Task Forces recognises that while IPv4 address space will be available for many years, new users and usages of the Internet have the potential to rapidly increase the utilisation of IPv4 address space. With the advent of multiple always-on devices, wireless handhelds and 3G mobile handsets, the Internet community needs to prepare for a sharp increase in IP address space utilisation. In order to prevent future operational problems, the global rollout of IPv6 is essential for enabling the development and adoption of new applications and services.

The rollout of IPv6 on this scale requires significant preparation, particularly in terms of training and planning. The RIRs and the IPv6 Task Forces encourage early evaluation by network operators and industry players, in order to promote the necessary technical dialogue and to facilitate widespread adoption. *Internet Service Providers* (ISPs) can already deploy IPv6 in non-disruptive ways that do not require additional investment while providing added value to their customers.

“The RIPE NCC has supported IPv6 from an early stage. We are committed to ensuring that IPv6 resources are provided to RIPE NCC members whenever they are required. We will continue to use the long-established system of address distribution where IP addresses are allocated according to demonstrated need wherever that need is demonstrated,” stated Axel Pawlik, Managing Director of the RIPE NCC. “The RIPE NCC is already providing IPv6 training to our members and other tools required to facilitate IPv6 deployment,” he added.

Jordi Palet, Founding Member of the EU IPv6 Task Force and co-chair of the IPv6 Forum’s Awareness and Education Working Group, sees the formalisation of this cooperative support of IPv6 deployment as an important development. “This cooperative effort ensures the global recognition of the strategic importance of IPv6 in enabling the continued development of the Internet and the worldwide information society. This ongoing coordination will have a positive global benefit for end users and the industry, by reinforcing the resilience of the Internet while allowing for the development of ever-improving applications and services,” he said.

Paul Wilson, APNIC Director General, noted that significant advances have been taking place in all the RIR regions with respect to IPv6 allocation and policy. “The RIRs are already working with the IANA and large ISPs to facilitate the delegation of large blocks of IPv6 address space,” he stated. “In the Asia Pacific region, a number of countries are taking the lead in terms of IPv6 deployment, and APNIC will continue to offer its support in these areas, and elsewhere, to allow the entire region to benefit from IPv6.”

“In the ARIN region, we have received clear direction from the community to make all necessary preparations for IPv6 deployment. This includes work on the allocation policies and procedures, as well as making our own services available via IPv6,” stated John Curran, Acting President of ARIN

“LACNIC is involved in the formation of the Latin American and Caribbean IPv6 Task Force and is active in encouraging the participation of its members and the community in IPv6 deployment and policy, and our services are already available over IPv6,” said Raúl Echeberría, CEO of LACNIC.

“This global cooperation signals another historic milestone to further accelerate take-up of IPv6 for the global good,” applauded Latif Ladid, President of the IPv6 Forum.

“The North American IPv6 Task Force supports the worldwide collaboration with the RIRs to further support the deployment of IPv6 and the next generation Internet mobile society using IPv6,” stated Jim Bound, Chair NAv6TF and IPv6 Forum CTO.

As an IPv6 Forum Board member and an ICANN Address Council member, Takashi Arano of the Asia Pacific IPv6 Task Force steering committee supports this collaboration. “Address management, which the RIRs are in charge of, is one of the crucial components for the commercial deployment of IPv6 and its stable operation.”

“I hope collaboration between IPv6 Task Forces and the RIRs will result in the advent of an IPv6-powered ‘everything-everywhere-every time’ networking world,” he stated.

IPv6 is a new version of the data networking protocols on which the Internet is based. The *Internet Engineering Task Force* (IETF) developed the basic specifications during the 1990s. The primary motivation for the design and deployment of IPv6 was to expand the available “address space” of the Internet, thereby enabling billions of new devices (PDAs, cellular phones, appliances, etc.), new users and “always-on” technologies (xDSL, cable, Ethernet-to-the-home, fibre-to-the-home, Power Line Communications, etc.).

The existing IPv4 protocol has a 32-bit address space providing for a theoretical 2^{32} (approximately 4 billion) unique globally addressable network interfaces. IPv6 has a 128-bit address space that can uniquely address 2^{128} (340,282,366,920,938,463,374,607,431,768,211,456) network interfaces.

The *European IPv6 Task Force* is a volunteer organisation, with over 500 members, open to all the interested parties in advancing the IPv6 deployment in the European region, in cooperation with the rest of the world and other related entities. Further information is available on the IPv6 Task Forces website: <http://www.ipv6tf.org>

Four RIRs exist today. They provide number resource allocation and registration services that support the operation of the Internet globally. The RIRs are independent, not-for-profit organisations that work together to meet the needs of the global Internet community. They facilitate direct participation by all interested parties and ensure that the policies for allocating Internet number resources (such as IP addresses and *Autonomous System Numbers*) are defined by those who require them for their operations.

The RIRs ensure that number resource policies are consensus-based and that they are applied fairly and consistently. The RIR framework provides a well-established combination of bottom-up decision-making and global cooperation that has created a stable, open, transparent and documented process for developing number resource policies.

The RIR framework contributes to the common RIR goal and purpose of ensuring fair distribution, responsible management and effective utilisation of number resources necessary to maintain the stability of the Internet. The RIRs currently consist of:

APNIC: *Asia Pacific Network Information Centre*
<http://www.apnic.net>

ARIN: *American Registry for Internet Numbers*
<http://www.arin.net>

LACNIC: *Latin American and Caribbean Internet Addresses Registry*
<http://www.lacnic.net>

RIPE NCC: *RIPE Network Coordination Centre*
<http://www.ripe.net>

The *IPv6 Forum* is a world-wide consortium of over 160 leading Internet service vendors, National Research & Education Networks and international ISPs, with a clear mission to promote IPv6 by improving market and user awareness, creating a quality and secure New Generation Internet and allowing world-wide equitable access to knowledge and technology. The key focus of the IPv6 Forum today is to provide technical guidance for the deployment of IPv6. IPv6 Summits are hosted by the IPv6 Forum and staged in various locations around the world to provide industry and market with the best available information on this rapidly advancing technology. <http://www.ipv6forum.org>

The *North American IPv6 Task Force* is an all-volunteer non-vendor/service/provider or other entity interest with the IPv6 mission of assisting the North American geography as sub task force of the IPv6 Forum for deployment, education, awareness, technical analysis/direction, transition analysis, political/business/economic/social analysis support and other efforts as required. The members see IPv6 as more important than their own self-interests. <http://www.nav6tf.org>

Upcoming Events

The *Internet Corporation for Assigned Names and Numbers* (ICANN) will meet in Kuala Lumpur, Malaysia, July 19–23, 2004, and in Cape Town, South Africa, December 1–5, 2004. For more information see: <http://www.icann.org>

ICANN and *The International Telecommunications Union* (ITU) will be jointly hosting a workshop on *country code Top Level Domains* (ccTLDs), in Kuala Lumpur on 24 July. The purpose of this joint ICANN/ITU-T open workshop is to focus on the operation and practical operational issues facing the ccTLDs and to give the opportunity for ccTLD operators and ITU Member States to share their experiences. The Workshop is not a policy meeting, but rather it is intended as a forum for the exchange of views and discussions. Written presentations are encouraged, but not required. Written presentations can be submitted to ICANN-ITU-T-Workshop@icann.org. Additional information can be found at the ITU-T website: <http://www.itu.int/ITU-T/worksem/cctld/kualalumpur0704/index.html>

The IETF will meet in San Diego, CA, August 1–6, 2004 and in Washington, DC, November 7–12, 2004. For more information, visit: <http://ietf.org>

Useful Links

The following is a list of Web addresses that we hope you will find relevant to the material typically published in the IPJ.

- The *Internet Engineering Task Force* (IETF). The primary standards-setting body for Internet technologies. <http://www.ietf.org>
- *Internet-Drafts* are working documents of the IETF, its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts. Internet-Drafts are not an archival document series.

These documents should *not* be cited or quoted in any formal document. Unrevised documents placed in the Internet-Drafts directories have a maximum life of six months. After that time, they must be updated, or they will be deleted. Some Internet-Drafts become RFCs (see below). <http://www.ietf.org/ID.html>

- The *Request for Comments* (RFC) document series. The RFCs form a series of notes, started in 1969, about the Internet (originally the ARPANET). The notes discuss many aspects of computer communication, focusing on networking protocols, procedures, programs, and concepts but also including meeting notes, opinion, and sometimes humor. The specification documents of the Internet protocol suite, as defined by IETF and its steering group the IESG, are published as RFCs. Thus, the RFC publication process plays an important role in the Internet standards process. <http://www.rfc-editor.org/>
- The *Internet Society* (ISOC) is a non-profit, non-governmental, international, professional membership organization. <http://www.isoc.org>
- The *Internet Corporation for Assigned Names and Numbers* (ICANN) "...is the non-profit corporation that was formed to assume responsibility for the IP address space allocation, protocol parameter assignment, domain name system management, and root server system management functions." <http://www.icann.org>
- The *North American Network Operators' Group* (NANOG) "...provides a forum for the exchange of technical information, and promotes discussion of implementation issues that require community cooperation." <http://www.nanog.org>
- The *Regional Internet Registries* (RIR) provides IP address block assignments for Internet Service Providers and others. See page 33 for links to APNIC, ARIN, LACNIC and RIPE NCC.
- The *World Wide Web Consortium* (W3C) "...develops interoperable technologies (specifications, guidelines, software, and tools) to lead the Web to its full potential as a forum for information, commerce, communication, and collective understanding." <http://www.w3.org>
- The *International Telecommunication Union* (ITU) "... is an international organization within which governments and the private sector coordinate global telecom networks and services." <http://www.itu.int>

This publication is distributed on an "as-is" basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Technology Strategy
MCI, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.
Copyright © 2004 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSR STD U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol *Journal*

September 2004

Volume 7, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Anatomy	2
Letters to the Editor	33
Fragments	36

FROM THE EDITOR

Network Address Translators (NATs) were designed to allow multiple devices in a private address realm to dynamically share a single public IP address. NATs are widely deployed in today's Internet. They provide an effective way of IPv4 address conservation while simultaneously offering some level of security because individual IP addresses on the "inside" are hidden from the "outside," or global Internet. But NATs also present a challenge to existing Internet applications that may depend on globally unique IP addressing for proper operation. To further complicate matters, not all NATs are created equal, leading to unpredictable behavior. This edition of IPJ is almost entirely devoted to an in-depth look at NATs. Geoff Huston looks inside the NAT, and explains the complexities behind each variation of NAT implementation. It seemed only natural that he would name such an exposé "Anatomy."

Many IPJ subscriptions had an official expiration date of September 30, 2004, but I am pleased to report that all these subscriptions have been extended for another year. You should still make sure your delivery address and e-mail is up-to-date in our database by using the link at www.cisco.com/ipj or sending e-mail to ipj@cisco.com with your updated information.

If you're hungry for even more networking-related reading material, look at the Internet Society's publication page at <http://isoc.org/pubs/>. Here you will find The ISP Column, Member Briefings, Articles of Interest, and links to other material.

We didn't have room for a book review in this issue, but we have several in store for future editions. If you'd like to contribute a book review for publication in IPJ, please contact me.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Anatomy: A Look Inside Network Address Translators

by Geoff Huston, APNIC

Over the past decade numerous IP-related technologies have generated some level of technical controversy. One of these is the *Network Address Translator*, or NAT. This article describes the inner workings of NATs in some detail, and then looks at the issues that have accompanied the deployment of NATs in the Internet that appear to have fueled this technical controversy. NATs are a very widespread feature of today's Internet, and this article attempts to provide some insight as to how they operate, why there is such a level of technical controversy about NATs, and perhaps some pointers to what we have learned about technology and the process of standardization of technology along the way.

NAT Motivation

The first RFC document describing NATs was by Kjeld Egevang and Paul Francis in 1994^[1]. The original motivation behind the NAT work was based on efforts in the early 1990s associated with a successor protocol to IPv4. The overall effort of a successor protocol to IPv4 was to devise a protocol that would directly address the issues of accelerating address consumption in IPv4 that appeared to be leading to the prospect of imminent address exhaustion. Although IPv4 was capable of uniquely addressing some 4.4 billion devices, it was evident by as early as 1992 that the world was heading down a path of very intensive deployment of devices that included communications capabilities, and that IPv4 was not going to be able to extend across the full range of future device deployment. The objective with NAT was to define a mechanism that allowed IP addresses to be shared across numerous devices. In addition, it was intended that NATs could be deployed in a piecemeal fashion within the Internet, without causing changes to hosts or other routers. Other forms of address-sharing technologies relied on intermittent connectivity, whereas NATs were intended to allow a collection of connected devices to share an address pool dynamically. The original RFC portrays this approach as being a measure that can “provide temporarily relief while other, more complex and far-reaching solutions are worked out.”

So, as documented, the original intent of NATs was to be a possible short-term response to address exhaustion while longer-term solutions were being devised. NATs were also intended to be unmanaged devices that are transparent to end-to-end protocol interaction, requiring no specific interaction between the end systems and the NAT device.

A decade later NATs are attaining a status of near-ubiquitous deployment across the Internet, and although IPv6 has been defined and deployment is commencing, NATs appear to be a very well-entrenched part of the network landscape. And, for the most part, NATs continue to function as unmanaged devices.

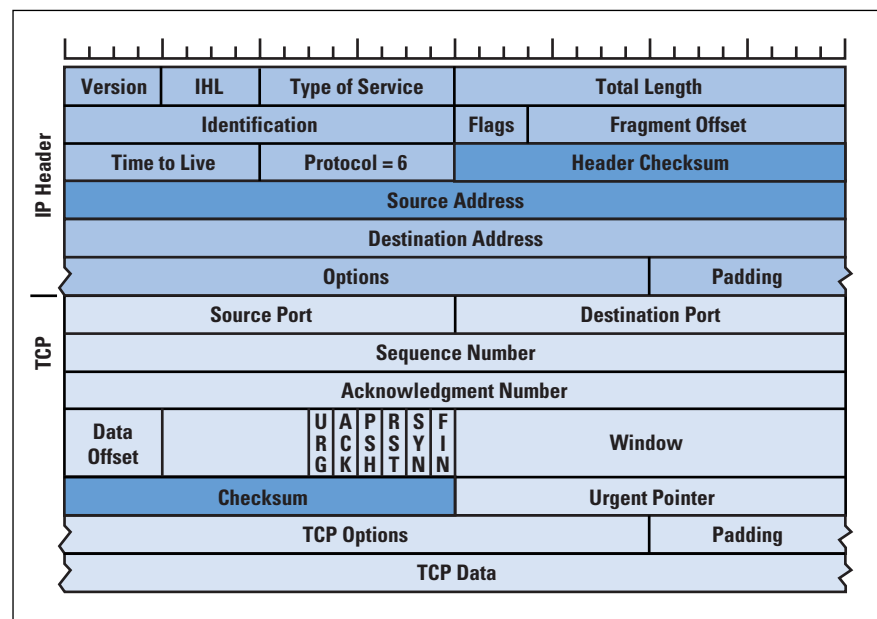
They can be transparent to some forms of protocol interaction, but, as the voice-over-IP folks are finding out, they can be very obvious to the point of being highly disruptive to other forms of protocol operation.

NAT Operation

The operation of NATs is deceptively easy to describe in general terms. They are active units placed in the data path, usually as a functional component of a border router or site gateway. NATs intercept all IP packets, and may forward the packet onward with or without alteration to the contents of the packet, or may elect to discard the packet. The essential difference here from a conventional router or a firewall is the discretionary ability of the NAT to alter the IP packet before forwarding it on. NATs are similar to firewalls, and different from routers, in that they are topologically sensitive. They have an “inside” and an “outside,” and undertake different operations on intercepted packets depending on whether the packet is going from inside to outside, or in the opposite direction.

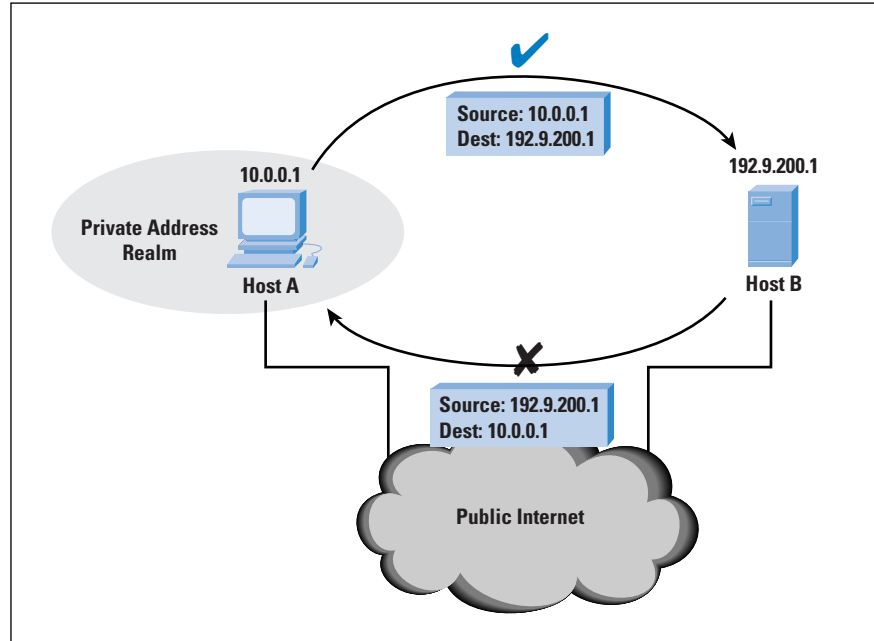
NATs are IP header translators, and, in particular, NATs are IP *address translators*. The header of an IP packet contains the source and destination IP addresses. If the packet is being passed in the direction *from* the inside *to* the outside, a NAT rewrites the source address in the packet header to a different value, and alters the IP and TCP header checksums in the packet at the same time to reflect the change of the address field. When a packet is received *from* the outside destined *to* the inside, the destination address is rewritten to a different value, and again the IP and TCP header checksums are recalculated (Figure 1). The “inside” does not use globally unique addresses to number every device within the network served by the NAT. The inside (or “local”) network may use addresses from private address blocks, implying that the uniqueness of the address holds only for the site. Let’s look at this using an example.

Figure 1: TCP/IP Header Fields Altered by NATs (Outgoing Packet)



As shown in Figure 2, how can local (private) host A initiate and maintain a TCP session with remote (public) host B? Host A first uses the *Domain Name System* (DNS) to find the public IP address for host B, and then creates an IP packet using host B's address as the destination address and host A's local address as the source, and passes the packet to the local network for delivery. If the packet was delivered to host B without any further alteration, then host B would be unable to respond. The public Internet does not (or should not at any rate!) carry private addresses, because they are not globally unique addresses.

Figure 2: Public/Private Communication



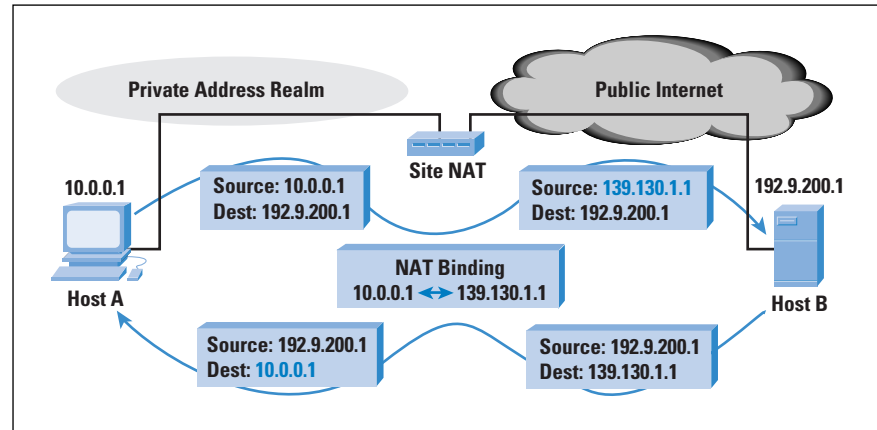
With a NAT between hosts A and B, the NAT intercepts host A's outgoing packet and rewrites the source address with a public address. NATs are configured with a pool of public addresses, and when an "inside" host first sends an outbound packet, an address is drawn from this pool and mapped as a temporary alias to the inside host A's local address. This mapped address is used as the new source address for the outgoing packet, and a local session state is set up in the NAT unit for the mapping between the private and the public addresses.

After this mapping is made, all subsequent packets within this application stream, from this internal address to the specified external address, will also have their source address mapped to the external address in the same fashion.

When an incoming packet arrives on the external interface, the destination address is checked. If it is one of the NAT pool addresses, the NAT box looks up its translation table. If it finds a corresponding table entry, the destination address is mapped to the local internal address, the packet checksums are recalculated, and the packet is forwarded. If there is no current mapping entry for the destination address, the packet is discarded.

The mode of operation of a NAT is shown in Figure 3. So, continuing our example, the local host at address A is directing packets to the external server host at address B. Because the NAT is in the path, the NAT has altered the packets so that address A is translated to address X. Host A is aware that it is communicating with host B, and from host A's perspective this is a normal session. Host B believes that it is communicating with a host at address X, and is entirely unaware of address A. From host B's perspective this is a normal session with a host at address X.

Figure 3: NAT Traversal

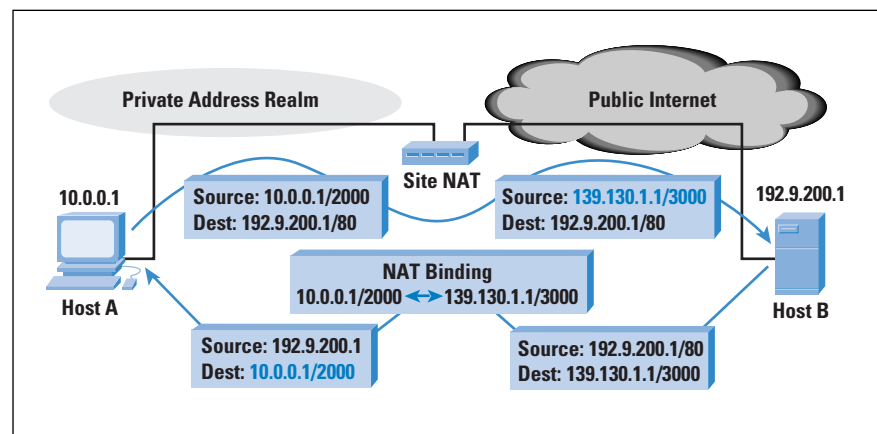


Dynamically created mapping entries (or “bindings”) are typically maintained by the NAT with a *timer*. If no packets that use the mapping are received by the NAT within a certain time window, then the binding is removed from the NAT and the public address is returned to the NAT pool.

NATs

A variant of the NAT is the *Port-Translating NAT*, or NAPT. This form of NAT is used in the context of TCP and *User Datagram Protocol* (UDP) sessions, where the NAT maps the local source address and source port number to a public source address and a public-side port number for outgoing packets. Incoming packets addressed to this public address and port pair are translated to the corresponding local address and port. Again, the binding is maintained by a NAT idle timer, and upon expiration of the timer the public address and port pair are returned to the NAT pool (Figure 4).

Figure 4: NAPT Traversal



Again the NAPT is attempting to be transparent in terms of providing a consistent view of the session to each end, using a symmetric binding of a local address and port pair to an external address and port pair.

A reasonable question to ask is: Why should NAPT's bother with port translation? Are straight address translations not enough? Surprisingly, NATs can be relatively profligate with addresses. If each TCP session from the same local host is assigned a different and unique external pool address, then the peak address demands on the external address pool could readily match or exceed the number of local hosts, in which case the NAT could be consuming more public addresses than if there were no NAT at all! NAPT's allow concurrent outgoing sessions to be distinguished by the combination of the mapped address and mapped port value. In this way each unique external pool address may be used for up to 65,535 concurrent mapped sessions.

For a while the terminology distinction between NATs and NAPT's was considered important, but this has faded over time. For the remainder of this article we use current terminology, and look at NATs and NAPT's together and refer to them collectively as "NATs."

NAT Behavior

The use of NATs involves two basic issues: One is that NATs make applications "brittle" in that NATs support a particular style of application operation, and if the application deviates in any way from this style then the application no longer works. The second is of much more concern, and that is that NATs differ from each other in quite fundamental ways. What works across one NAT may not work at all for another class of NAT. It has also been reported that NATs differ not only on a vendor-by-vendor basis, but even on a model-by-model basis within a single vendor's range of NAT units. The implication here is that such differences of behavior become a matter for discovery by applications rather than something applications can predict in advance. This section explores this behavioral aspects of NATs in further detail.

Symmetry and Sessions

NATs can manage address mapping in numerous ways, and many implementations of NATs use a form of binding termed a "symmetric" binding.

A *symmetric* binding is where the mapping of a local address to a public address is exclusively tied to the destination address used in the initial trigger outgoing packet for the lifetime of the binding. Incoming external packets with the mapped public address as their destination are translated to the local address only if the source address of the incoming packet matches the destination address of the original mapping. Multiple sessions to different public hosts may use the same mapped public address, or may use different public addresses for each session. This mapping is "endpoint" sensitive. Symmetric NATs represent a restricted model of operation, where each NAT binding represents a window through the NAT that is visible only to the destination host (Figure 5).

By comparison, a *full-cone* NAT allows any external host to use this opened window, where all incoming packets addressed to the mapped external address are translated to the mapped internal address and forwarded through the NAT. Symmetric NATs represent the most restrictive form of behavior, whereas full-cone NATs represent a far more permissive mode of operation.

In the context of NATs, this symmetric mode of operation refers to the session state 5-tuple, made up of Transport Protocol, the local IP address and port number, and the destination IP address and port number. When a session is opened from the local host to a remote service port on a remote host, then only that remote service can pass packets back through the NAT to the local host on that port. As with NATs, a full-cone NAT allows any remote service entity to direct packets back through the port window.

NATs can be further refined by having different behaviors for TCP and UDP transports. A NAT may behave in a symmetric manner for TCP sessions, and operate in a full-cone mode for UDP transactions. The variations in NAT behavior has led to an exercise in categorizing NAT behaviors and developing a discovery protocol whereby a pair of cooperating systems can discover if one or more NATs is on the network path between them, as well as attempting to establish the type of NAT.

Discovering NAT Behaviors and STUN

NAT behavior has not been the topic of any industry standardization efforts, and it should not be surprising to learn that, given that a range of possible NAT behaviors exist under certain conditions, the market contains NAT offerings that cover the full spectrum of possibilities. In the absence of common specifications or standards, implementers have been placed in the position of having to make some creative guesses as to what the “right” behavior should be under such circumstances. This is a significant problem for the application designer, given the prospect that in today’s Internet any popular application must have a means of being able to function correctly in the face of one or more NATs on the path between two hosts that are communicating using the application.

One of the more pressing problems here is that NATs commonly enforce an application model where the local “hidden” host must initiate a transaction in order to create a window in the NAT to allow the packets of the remote host back into the local network.

Some applications may wish to undertake “referral,” where the correspondent host on the external side may want to pass the externally presented address and port details of the local host to a third party in order to commence a further part of the transaction. Other application transactions may simply want to be initiated from the external side. Although this may have been thought of as a relatively obscure condition, it was brought into the forefront of attention when various forms of voice-over-IP and peer-to-peer applications gained popularity. In particular, the question of “how can the external side initiate a packet flow in the presence of a NAT?” has become increasingly important.

Given that the application needs to perform some additional gymnastics in such a case, there is the additional question that the application must answer, namely: “How does the application learn that there are NATs in the path in the first place?”

At this point the application is placed in the role of performing a forensic exercise of establishing whether or not its packets are being altered by one or more NATs when it attempts to establish an end-to-end packet transaction. If so, what types of implementation decisions have been made by the NAT in terms of the way in which packets are being systematically modified? In other words, what is the anatomy of the particular NATs that have been discovered along the path? This anatomy exercise is further complicated by the observation that NATs are silent devices, so the application cannot directly interrogate the NAT to establish its behavior. All that is left is a somewhat unsatisfying guessing game for the application. It is forced to send particular types of test packets through the NAT to some pre-defined counterpart on the other side. The application must then compare the self-view of the IP address and port number of the local host to the remote view of its IP address and port number, and then attempt to guess the nature of the systematic transforms that the NAT is applying.

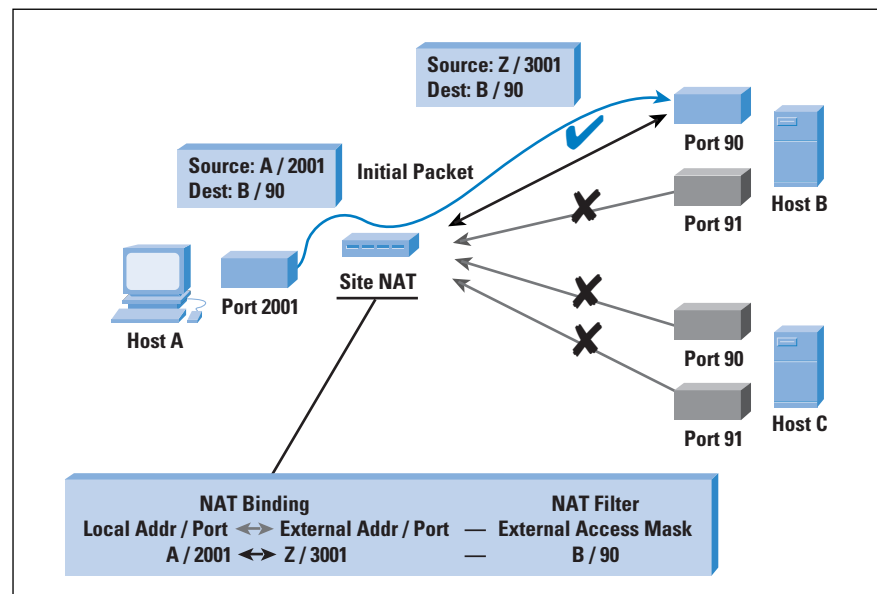
In the case of TCP it appears that the prevalent NAT behavior is that of a symmetric NAT based on address and port bindings. This implies that when the local host opens up a TCP session with a remote host, the NAT address and port bindings for the local host are coupled with the address and port of the destination host. Only packets with a source field of the destination host can pass packets back through the NAT to the TCP session of the local host. In other words, when a TCP session has been established within a NAT, only the two endpoints of the TCP session can access the NAT bindings, and attempts by others to direct packets to the external-side presented address and port meet with the NAT discard response. The fine-grained behavior of NATs with respect to TCP sessions can vary according to the amount of TCP state maintained by the NAT. At a basic level, the NAT can maintain a binding based on the local address and port and the remote address and port. The NAT also can keep the binding timer at a high value until a **FIN** exchange is observed, or until the session is reset through the **RST** flag being set, at which point the binding timer can be reduced to a very short interval. The NAT can also track the sequence number windows of the two sides and associated window sequence number scaling values and not adjust the binding timer of the session for TCP packets with sequence numbers outside the sequence number window with their **FIN** or **RST** flags set.

These NAT behaviors are based on the explicit signaling of changes in session state within the TCP packet exchange, and the consequent ability of the NAT to track the session state and adjust the associated binding timer in response to this state information. UDP is not so straightforward, because there is no explicit session state within a UDP packet exchange, and various NATs behave differently with respect to UDP-based bindings.

Various classes of NAT behavior relate to how UDP bindings are managed within a NAT. These have been classified into four types of behaviors^[11]:

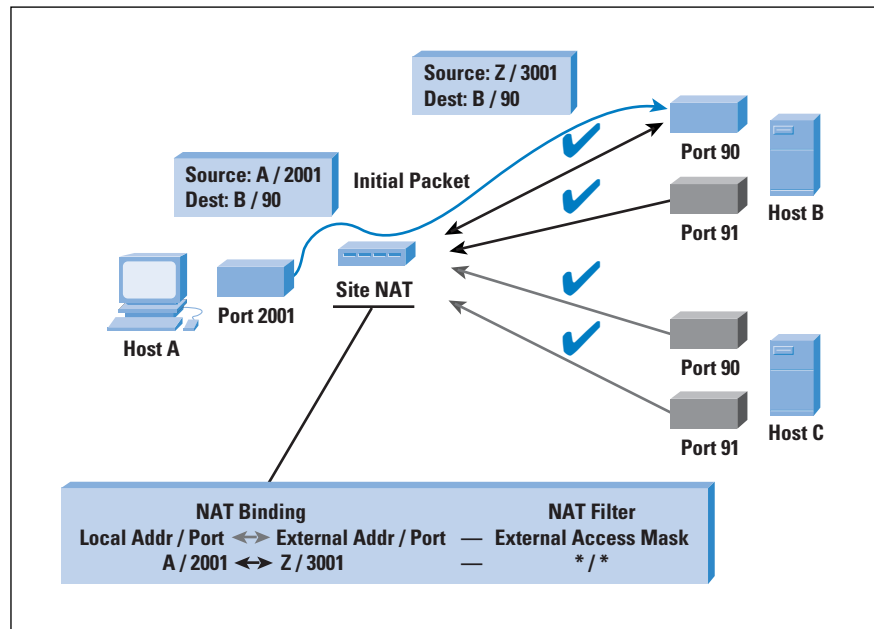
- *Symmetric*: We have already encountered the symmetric NAT, where the NAT mapping refers specifically to the connection between the local host address and port number and the destination address and port number and a binding of the local address and port to a public-side address and port. Any attempts to change any one of these fields requires a different NAT binding. This is the most restrictive form of NAT behavior under UDP, and it has been observed that this form of NAT behavior is becoming quite rare, because it prevents the operation of all forms of applications that undertake referral and handover.

Figure 5: Symmetric NAT



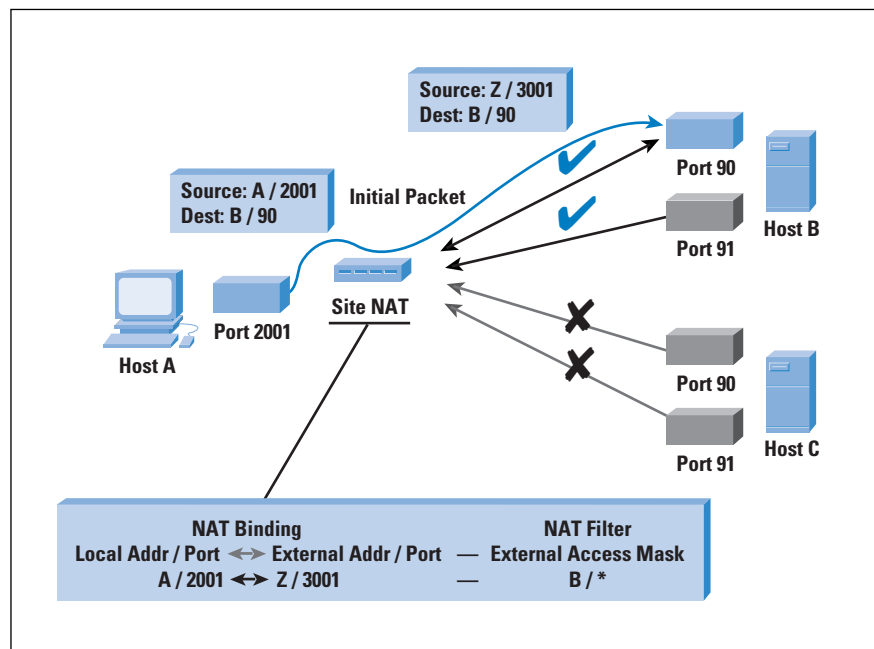
- *Full-cone*: A full-cone NAT is the least restrictive form of NAT behavior, where the binding of a local address and port to a public-side address and port, when established, can be used by any remote host on any remote port address. (Refer to Figure 6.)

Figure 6: Full Cone NAT



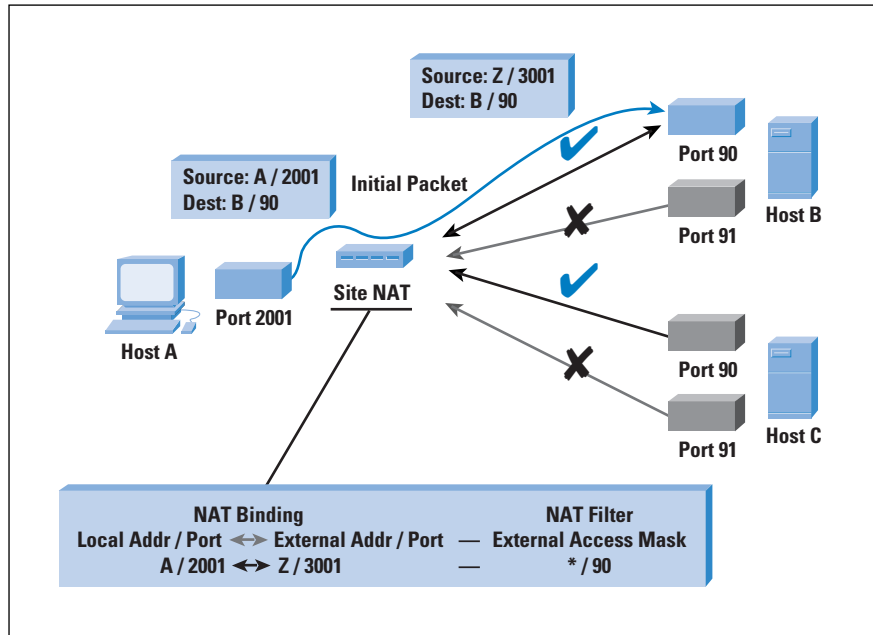
- *Restricted-cone*: A restricted-cone NAT is one where the NAT binding is accessible only by the destination host, although in this case the destination host can send packets from any port address after the binding is created. (Refer to Figure 7.)

Figure 7: Restricted-Cone NAT



- *Port-restricted-cone*: A port-restricted-cone NAT is one where the NAT binding is accessible by any remote host, although in this case the remote host must use the same source port address as the original port address that triggered the NAT binding. (Refer to Figure 8.)

Figure 8: Port-Restricted-Cone NAT



So can an application tell if one or more NATs are in the path, and, if so, what form of behavior the NAT is using? For this purpose the *Simple Traversal of UDP through NATs* (STUN) protocol has been developed^[11]. STUN is a probe system that examines the interchange between a STUN client that may lie behind a NAT and a STUN server that is positioned on the public side of the NAT. The STUN-server host must be configured with two IP addresses, and the STUN itself should respond to queries on two UDP port numbers. The protocol is a simple UDP request-response protocol that uses embedded addresses in the data payload, and compares these addresses with header values in order to determine the type of NAT that may lie in the path between client and server.

The basic operation of STUN is a request-response protocol, using a common request of the form: “Please tell me what public address and port values were used to send this query to you.”

STUN can be used to discover if a NAT is on the path between a client and server, and attempt to discover the type of NAT by a structured sequence of requests and responses. The client sends an initial request to the STUN server. If the public address and port in the returned response are the same as the local address, then the client can conclude that there is no NAT in the path between the client and the server. If the values differ, the client can conclude that there is a NAT on the path. STUN then uses subsequent requests to determine the type of NAT. One critical additional item of information returned by the STUN server in the initial response is an alternate IP address and port number that can also reach the same STUN server.

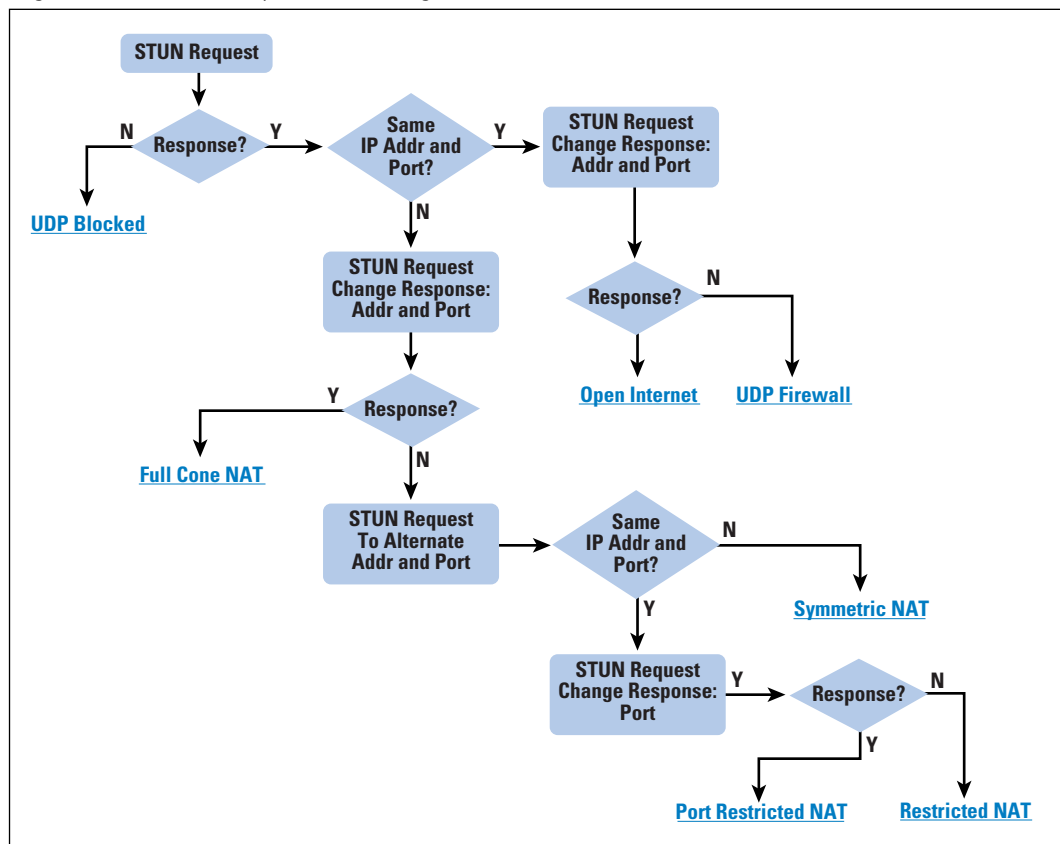
The second STUN request is directed to the same address and port as the initial request, but this time the request includes a control flag that requests the STUN server to respond using its alternate source address and port values. If the STUN client receives this alternate-sourced response, then it can conclude that it is behind a full-cone NAT. This is because the initial NAT binding of the local host address to the external presentation address can evidently be accessed by third-party external hosts.

If no response is received to the second request, then the STUN client sends the original probe request, but this time the request is addressed to the alternate destination address and port pair for the STUN client. If the returned address and port values relating to the new NAT binding are different from those of the first request, then the client can conclude that it is behind a symmetric NAT.

If the values are unaltered, then a further request can be made to determine the form of restricted-cone behavior. This fourth request includes a control flag to direct the STUN server to respond using the same IP address, but with the alternate port value. A received response indicates the presence of a port-restricted cone, and the lack of a response indicates the presence of a restricted cone.

Periodic exchanges between the STUN client and server can also discover the timer used by the NAT to maintain address bindings. Additional components of STUN are intended to provide some reasonable level of integrity in the packet exchange. A flowchart of a STUN-based NAT discovery process is shown in Figure 9.

Figure 9: NAT Discovery Process Using STUN



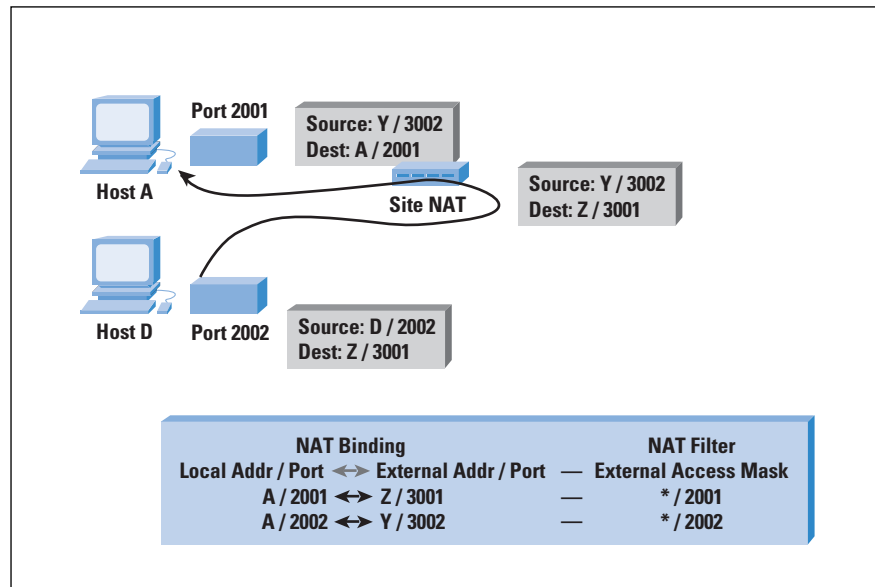
Further Behaviors: Hairpins and Determinism

It would be good if NAT behavior remained that simple. However, it does not, and some further tests on NATs reveal further differences in various NAT implementations^[16].

The first area of difference is whether the NAT supports the so-called *hairpin* operation, where a local host directs a packet to the public address and port of an already mapped local host, or even to its own mapped address and port. If successful, then the NAT supports hairpin operation, where the NAT bindings, when created, are available to either side of the NAT. (Refer to Figure 10.)

Furthermore, the NAT may generate a binding for this operation—or not—thereby presenting the hairpin packet with an external address and port, indicating that an outbound binding has been performed in conjunction with the inbound binding, or with an internal address and port, indicating that only an inbound binding is being performed.

Figure 10: Hairpin NAT Operation



The second is in the general class of NAT determinism. Nondeterministic NATs change their binding behavior when a binding conflict of some sort occurs in the NAT. This is further based on the classification of whether “primary,” “secondary,” or even “tertiary” NAT behaviors differ. To explain primary, secondary, and tertiary behaviors, it is first noted that some NATs attempt to preserve the port address in the binding, so that the local source port and the externally bound port are the same whenever possible. This is the “primary” binding of the NAT. If another local host obtains a NAT binding using the same source port number, then the behavior of the NAT for this conflicting port binding may differ from that where the port number is preserved. The first conflict of port allocations in bindings is the “secondary” binding. In some cases the primary behavior is that of a full cone, or a restricted cone, while the NAT behaves in a symmetric fashion for the secondary instance where the port number has been mapped to a new value by the NAT.

A tertiary behavior occurs when a third binding is added to the NAT, because, again, the behavior of the NAT may be different for this binding.

It is also possible that the NAT may elect to preserve the binding in any case, and remove the current binding and replace it with a new binding that refers to the most recent packet that the NAT has processed.

All these behaviors can be classified as *nondeterministic*, in that the NAT behavior becomes one that is determined by the order of out-bound traffic. The implication is that repetitions of the same STUN test at different times may produce different classifications of the type of NAT. The inference is that if an application uses STUN to determine the type of NAT in the path, and then selects a certain behavior based on this STUN-derived knowledge of the NAT type, nondeterministic NATs may behave differently between the STUN test and the application. The NAT response for a particular binding cannot be predicted in advance, and even when a binding state is established it may be disrupted or altered by subsequent traffic.

Another Approach to Classifying NATs

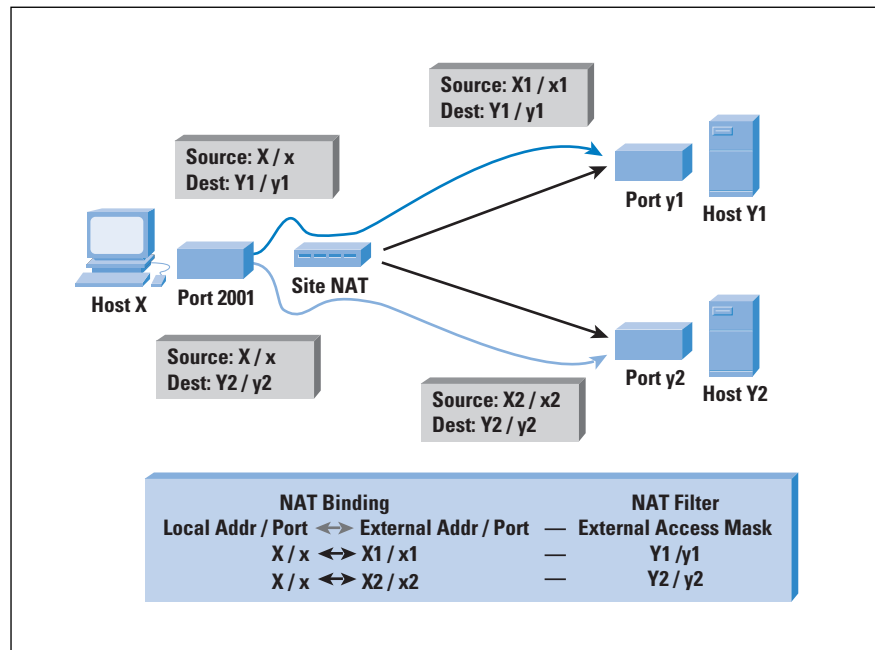
Further tests on NATs reveal that the various behaviors are yet more complex, and that different sequences of tests across a NAT will lead the test routine to come to different conclusions as to the type of NAT^[13]. The key observation here is that NATs are the conjunction of two distinct behavior sets:

- *Binding*, or context-based packet translation: Detecting those packets that can be associated with a current binding and using that binding in a manner according to the logical direction of the packet to perform packet header transforms
- *Filtering*, or packet discard: Discarding those packets that cannot be associated with current bindings and discarding them

If a STUN-like test sequence was for a local host to send a packet to one destination and obtain a response of what NAT binding was used, and then to send a packet to a second destination and compare the results, the observation of the NAT using a different binding for each request may lead the tester to conclude that the NAT is a fully symmetric NAT. If the test sequence is for the NAT to send one packet to a destination and have the destination respond using a different source address, then the observation that the response packet is successfully delivered through the NAT back to the originating local host may lead the tester to the conclusion that the same tested NAT is some form of cone NAT.

The STUN approach classifies NAT behaviors on the basis of a single binding being established by the local host when contacting an external host, and then considers what constraints are placed on third-party external hosts as they attempt to access this initial binding. An adjunct to this approach is based on the local host establishing two bindings to two distinct external hosts, and looking for any relationship between these two bindings. (See Figure 11).

Figure 11: Outbound Connections from a Common Source



The behaviors of NATs under this condition can be classified under numerous behavioral aspects.

Binding

Binding behavior can be seen as the amalgam of three somewhat distinct design decisions, namely the manner in which a binding is generated, the behavior of the NAT in managing external ports used in bindings, and the manner in which expiration timers that govern the continued existence of the binding are refreshed.

NAT Binding Behavior:

- *Endpoint independent:* The NAT reuses the port binding for subsequent sessions initiated from the same internal IP address and port to any external IP address and port. This is analogous to a full-cone NAT.
- *Endpoint address dependent:* The NAT reuses the port binding for subsequent sessions initiated from the same internal IP address and port only for sessions to the same external IP address, regardless of the external port. This is a looser form of symmetric NAT, where the binding is created on the basis of the external address, rather than the external address and port.
- *Endpoint address and port dependent:* The NAT reuses the port binding for subsequent sessions initiated from the same internal IP address and port only for sessions to the same external IP address and port. This is a more precise form of UDP symmetry where the binding is available only to a single session, where a session is the 5-tuple of protocol, source address, source port, destination address, and destination port.

Port Binding Behavior:

- *Port preservation:* In addition to the differences in the binding between the two cases, the NAT may attempt to preserve the local port number, if possible. The terminology proposed here is port preservation to describe this NAT action.
- *Port overloading:* Some NATs attempt to undertake port preservation at all times, so that when a different local host establishes a binding using a port that is already being preserved, the new binding will usurp the existing binding. This behavior is proposed to be termed port overloading.
- *Port multiplexing:* The alternative to port overloading is use of the external entity to perform the demultiplexing of the port. In this case if two local systems use the same source port to send packets to two different external hosts, the NAT preserves the source port in the two bindings. If the NAT is using a single external address, the external view is two packets with the same source address and source port, sent to two different external addresses. The reverse packets have the same destination address and port, and the NAT determines the appropriate binding based on the source address and port in the reverse packets. This requires an endpoint address and port-dependant binding behavior. If two internal hosts are directing packets to the same external endpoint using the same source port addresses, then it is necessary for one of the sessions to use a binding with an altered port number. This could be considered as nondeterministic behavior.

Binding Timer Refresh:

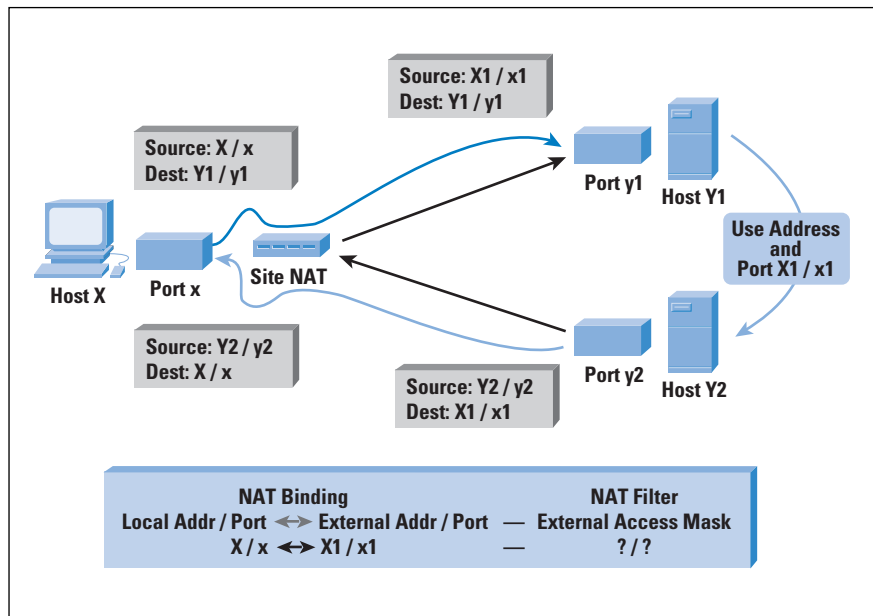
- *Bidirectional:* The NAT does not keep the binding active indefinitely, and normally removes the binding if there are no further packets that use the binding within a certain time period. However, there are variations in the classification of packets that the NAT considers as packets that reset the timer. In the case of bidirectional binding timer refresh, packets from either the local hosts or an external host that uses the NAT binding cause the NAT binding expiration time to be reset.
- *Outbound:* An outbound binding timer refresh NAT resets the expiration timer only when packets pass from the local host to the external host within the context of the binding. The implication is that a local host may have to use some form of keepalive operation to maintain a NAT binding in the face of an inbound UDP unidirectional traffic flow. Additionally, the expiration timer may be on a per-session basis, or may be on a per-binding basis if multiple sessions are associated to a single binding in the NAT.
- *Inbound:* As the name suggests, this is the opposite of the previous case, where only inbound packets cause the expiration timer of the binding to be refreshed.

- *Transport Protocol state:* Although these forms are useful in the case of UDP-based sessions, when the binding is based on a transport session (such as TCP), the NAT can base its binding timer refresh on the transport session state. For TCP this would infer a binding refresh time that is refreshed by any session packet in either direction (bidirectional), with the exception of packets with the TCP **RST** or **FIN** flags set. Although it would be an option to drop the NAT binding state when such packets are seen, this makes the NAT vulnerable to denial-of-service attacks by third-party injection of TCP **RST** packets, so there is some merit in using the binding timer for TCP sessions.

Filtering

The second phase of the test has two external hosts directing a probe to the same binding address, and classifying the behaviors based on what packets are filtered and discarded by the NAT (Figure 12).

Figure 12: Inbound Test



External Filtering:

- *Endpoint independent:* The NAT does not filter and discard packets that are addressed to the external part of the binding, irrespective of the source values in the packet. This is analogous to a full-cone NAT.
- *Endpoint address dependent:* The NAT filters and discards packets that are addressed to the external part of the binding, unless the source address of the packet matches the destination address used in the binding. This is analogous to a restricted-cone NAT.
- *Endpoint address and port dependent:* The NAT filters and discards packets that are addressed to the external part of the binding, unless the source address and port number of the packet matches the destination address used in the binding. This is analogous to a port-restricted-cone NAT or a symmetric NAT.

External Filtering Timer Refresh:

As with binding timers, these timers can be refreshed bidirectionally, inbound or outbound.

NAT Behaviors

The approach of carefully identifying the areas where NAT behaviors differ and classifying these behavioral differences in a methodical manner is one that has the potential to at least allow us to use the same sets of words when we talk about NAT behaviors, and hopefully also refer to the same set of actual behaviors when we use the same descriptions. The original approach with the STUN work used the terms *symmetric*, *full-cone*, and forms of *restricted-cone* to describe variations of NAT behaviors. Experience with this form of classification has exposed further variations in NAT behaviors, and this has led to a form of NAT classification that first uses a delineation of binding and filtering behaviors, and then classifies the various ways in which these bindings and filters are maintained within the NAT. Additional classification attributes include whether the NAT supports hairpin connections or not and whether it operates in a deterministic or nondeterministic manner.

This exercise is not another study in comparative taxonomies. A NAT has no standard way in which to advertise its presence, nor does it have any standard way in which to advise protocols or applications of the particular behaviors it applies to packets being passed through the NAT. In the absence of such explicit advertisements of the presence of a NAT, it is left to the application to make the necessary adjustments that allow it to function in the presence of NATs. The aim of behavioral classification is to associate test sequences that expose the presence of a NAT, and to determine its behavior. This allows applications to invoke a test procedure that exposes a particular choice of behaviors of a NAT implementation, and then allows the application to invoke a mode of operation that can operate across the particular NAT.

The choices available to application environments include the use of *agents* as session initiation intermediaries, where the endpoints make initial contact through agents, who then assist in passing binding information to the endpoints, allowing them to directly communicate. Other forms of application behavior need to be invoked when the NAT is endpoint address and port dependant for both binding and filtering. Different application responses are applicable when one endpoint is behind a NAT and when both endpoints are behind NATs. A typical application response in this latter case where both endpoints are behind highly restrictive NATs is for the endpoints to use agents as session intermediaries, so that the application payload is then passed through the intermediaries because an end-to-end pair of NAT bindings cannot be established.

Living in a NAT World

It would be a reasonable conclusion to draw from the previous sections that we are left in the somewhat unsatisfying position of observing that there is near-universal deployment in today's Internet of NAT devices that do not conform to any particular well-defined behavior set. NAT behavior varies across implementations, and NATs have no ability to disclose their particular behaviors to applications that are attempting to compensate for their presence in the path. It is extremely challenging for applications to reliably predict the behavior of the NATs that lie in the path, and more so in the face of multiparty applications, such as interactive game environments, where the application is attempting to understand the level to which this silent intermediary is capable of supporting a relatively promiscuous NAT binding state in terms of external entities that wish to send packets to the local host, and communicate between themselves about the local host as a single entity.

NATs, Client-Server, Peer-to-Peer, and Multiparty Applications

NATs, as a class of devices, have strong associations with a client-server model of communications. As long as all the servers have a consistent external visibility, with stable addresses in terms of an IP address and port number, and as long as clients initiate connections with servers in a fixed two-party communications model using TCP as a transport protocol, and refrain from turning on *IP Security* (IPSec), then NATs generally behave in a relatively stable and unobtrusive manner. Applications that operate conservatively in this limited mode can be unaware of the presence of NATs in their path. The relatively widespread deployment of NATs and the continued use of client-server-based applications on the Internet attests to the capability of the NAT to perform transparently and effectively within the strict confines of this particular mode of communication.

However, peer-to-peer applications are more problematic for NATs, because they have extended the model of a NAT beyond its original realm of capability. If the desire is to continue to support the NAT dynamic binding, but also allow external parties to initiate a communication to a local host, then the NAT ceases to be transparent and unobtrusive, and in this extended environment the NAT transforms itself into an application-visible network element. It is overly presumptuous to claim that NATs have led to the increasing deployment of multiparty applications on the Internet, but certainly multiparty applications have been seen to be useful in circumventing some of the more aggravating shortcomings of NATs in various peer-to-peer realms.

In this latter context, the local party is forced to advertise its willingness to participate in a peer-to-peer realm by communicating with an external agent. The local agent performs a NAT discovery test, and then selects a mode of operation that is consistent with the discovered behaviors of a NAT that may be on the path between the client and the agent. The agent then advertises itself as the local party's intermediary to other peers within the application realm. Attempts to initiate a connection with the local party are directed to the external agent, who then undertakes to perform a rendezvous function in order to establish a session.

Depending on the NATs that may exist between the two parties, the rendezvous function may need to perform a convoluted handshake process, or, in some instances, may not be able to set up a peer-to-peer session at all. This topic of establishing connectivity in the face of NATs in the path is sufficiently complex to warrant a separate examination, and the various techniques and approaches are not examined in this article other than providing some suggestions for further reading.

The salient general observation is that NATs have fueled a new generation of applications that use intermediaries and rendezvous protocols. This shift in application behavior has implied greater attention to security frameworks for applications, because intermediaries represent an additional active element in the trust model. This, in turn, has implied that the application level has to turn to other chains of derivation of trust, because the basic Internet model of some form of persistent identity as being an attribute of an IP address is no longer a workable proposition in the face of NATs. The position we are reaching here is that identity and trust need to be derived from other attributes of the end host and the application that it has invoked.

ICMP

If an *Internet Control Message Protocol* (ICMP) message is passed through NAT, there is not only the outer IP header to consider, but also the ICMP payload. Most ICMP messages contain part of the original IP packet in the body of the message, so for the NAT to behave as transparently as possible, the IP address of the IP header contained in the data part of the ICMP packet should be modified according to the NAT binding state, as well as the IP header Checksum field of this inner packet header.

NATs and IP Fragmentation

NATs that use bindings that include both address and port values do not have a clear and uniform response to fragments of an IP packet. The TCP or UDP header is resident only in the initial IP fragment, and subsequent IP packet fragments do not contain a copy of the transport layer packet header.

Some NATs attempt packet reassembly as if they were the end host, and they perform the NAT translation only when the original IP packet has been reassembled. Of course the reassembled packet may be too large to be forwarded onward, and the NAT may be forced to further fragment the packet. The interplay between this behavior and various forms of path *Maximum Transmission Unit* (MTU) discovery become a source of frustration.

Other NAT packet fragmentation behaviors do not attempt packet reassembly, but rely on a stored packet fragment translation state that directs the translation to be performed on subsequent packet fragments after the initial packet header translation has been performed on the initial IP packet fragment.

This form of behavior has weaknesses in terms of out-of-order fragments, when following fragments are received by the NAT prior to the initial IP packet fragment, and in such cases the NAT often has little choice but to silently discard the out-of-order fragment as untranslatable.

NATs and Application Level Gateways

This brings up one of the more vexing questions regarding NAT behavior, namely, should the NAT include knowledge of the payload of certain applications? Numerous applications, including FTP and the DNS resolution protocol, include IP addresses within the payload of the application. In an effort to achieve complete transparency of operation, some NATs have included *Application Level Gateway* (ALG) functionality for certain applications so that this use of IP addresses in the payload can be detected and altered according to the current NAT translation bindings.

The case of ICMP represents one of the simpler forms of gateway functionality, because it can be performed in the same manner as the basic NAT transform, on a per-packet basis while attempting to maintain retained session state. Payload transformations in the case of a TCP-based application have implications in terms of requiring subsequent alteration of TCP sequence numbers, length fields, and even the repacketization of the payload data stream, given that the data transform required by the address change may imply a change of payload length.

Some units attempt to combine the functionality of a NAT with that of an ALG, such that the NAT is an active intermediary in the transport session. This allows the NAT/ALG to perform “deep” inspection of the packets, and use both application protocol knowledge and per-application-session retained state in order to apply the NAT binding transforms to the application payload as well as to the outer IP packet header.

The most widely deployed application that can use IP addresses in the payload is FTP, where IP addresses are passed in the payload of the control channel in order to allow data sessions to be initiated on distinct transport sessions. The variability and reliability of FTP ALG support in NATs has led to the widespread use of the passive mode of FTP operation, where the data flow is passed within the control session.

A related question is that of the use of IPSec and NATs. IPSec with *Authenticated Header* protection attempts to protect what it believes is the fixed part of the IP packet header, including the source and destination addresses. The NAT changes to the IP packet invalidate the Authentication Header integrity check. Also the NAT changes the IP and UDP or TCP checksums, and this disrupts the *Encapsulating Security Payload* (ESP) function of IPSec. The implication is that IPSec needs to operate upon a TCP or UDP payload, as in the IPSec operating tunnel model, or IPSec carried as a payload within other types of tunnel operation.

It is also the case that NATs today are heavily enmeshed with the UDP and TCP transport protocols. Other transport protocols exist, including the *Streams Control Transport Protocol* (SCTP) and the *Datagram Congestion Control Protocol* (DCCP), and doubtless more transport protocol offerings will follow over time. In each case it is a matter of individual choice how NAT implementations define NAT responses to such additional transport protocols. Although it is tempting to propose that NATs should fall back to an address-only form of binding that was not address-and-port based, this does not appear to be practical guidance. Another aspect of today's NAT deployment is that the most common scenario appears to be that of a single external address and mapping each locally initiated session into a binding that uses this common external IP address and a variable external port number. This means that NATs need to be able to identify and transform port addresses from the Transport Protocol section of the IP header.

Another salient factor here is the common association of NATs and firewalls into a single unit, and the coupling of address utilization compression properties of the NAT with its associated packet-filtering actions. Deploying a NAT at the external interface of a site does lead to more restrictive site filtering outcomes and a more restrictive model of application interaction, where the model attempts to impose the constraint that applications are initiated from within the site, and that unknown or unidentifiable external traffic is considered hostile and should be subject to firewall-based inspection and filtering. From this perspective there is little desire to make more permissive NATs as an isolated exercise, and there is instead a codependence between NAT behaviors and popularly used applications. Applications that work across today's NATs appear to enjoy popular uptake, and applications that enjoy popular uptake appear to determine what forms of traffic pass across NATs.

Popular or not, there are a class of applications that simply cannot work in a "native mode" across NATs, nor can ALGs assist here. These are applications that attempt to impose some level of end-to-end protection on the IP header fields, or use the IP address of the endpoint in a context of some form of persistent identity token. When the NAT alters the IP address, an application that uses strong forms of header validation rejects such packets as corrupted. Within this class of applications and tools, one of the more commonly referenced tools is that of IPSec with Authentication Header. There is a certain sense of irony in the observation that NATs are often seen as part of an overall approach to site security, yet cannot support a "native mode" operation of some of the basic tools that applications could use to support secure end-to-end data transfer.

Views on NATs

It is certainly the case that NATs are very common in today's Internet, and it is worth understanding why NATs have enjoyed such widespread deployment while other technologies appear to be meeting some considerable resistance to widespread deployment. As the original NAT document points out:

“The huge advantage of this approach is that it can be installed incrementally, without changes to either hosts or routers. (A few unusual applications may require changes.) As such, this solution can be implemented and experimented with quickly. If nothing else, this solution can serve to provide temporarily relief while other, more complex and far-reaching solutions are worked out.”

—Egevang and Fancis,
“Network Address Translator,” RFC 1631

More generally, the positive attributes of NATs include the following considerations:

- End hosts and local routers do not change. Whether there is a NAT in place between the local network and the Internet or not, local devices can use the same software and support the same applications. NATs do not require customized versions of operating systems or router images.
- As long as you accept the limitation that sessions must be initiated from the “inside,” NATs can work in an entirely transparent fashion for a set of client-server classes of applications.
- If you accept the perspective that services and usage scenarios that are not supported by NATs are “unwelcome” or “unsafe,” then NATs can be placed into a role as a component of a site’s security architecture, providing protection from attacks launched from the outside toward the inside network.
- NAT conserves its use of public address space.
- NAT allows previously disconnected privately addressed networks to connect to the global Internet without any form of renumbering or host changes—and renumbering networks can be a very time-consuming, disruptive, and expensive operation, or, in other words, renumbering is difficult.
- NAT address space is an effective, provider-independent addressing solution with multihoming capabilities. NAT allows for rapid switching to a different upstream provider, by renumbering the NAT address pool to the new provider’s address space. In essence, NATs provide the local network manager with the flexibility of using provider-independent space without having to meet certain size and use requirements that would normally be required for an allocation of public, provider-independent address space.
- NAT allows the network administrator to exercise some control over the form of network transactions that can occur between local hosts and the public network.
- NATs require no local device or application changes. This is perhaps one of the major “features” of NATs, in that the local network requires no changes in configuration to operate behind a NAT.

- NATs do not require a coordinated deployment. There is no transition, and no “flag day” across the Internet. Each local network manager can make an independent decision whether or not to use a NAT. This allows for incremental deployment without mutual dependencies.
- These days the common theme of the public address assignment policy stresses conservative use of address space with minimum waste. The standard benchmark is to be able to show that a target of 80 percent of assigned address space is assigned to a number of connected devices. Achieving such a very high usage rate is a challenging task in many network scenarios, and NATs represent an alternative approach where the local network can be configured using private addresses without reference to the use of public addresses.
- NATs are very widely available and bundled into a large variety of gateway and firewall units. In many units NATs are not an optional extra—they are configured in as a basic item of product functionality.

The market has taken NATs and embraced them wholeheartedly. And in a market-oriented business environment, what is wrong with that?

Unfortunately NATs represent a set of design compromises, and no delving into the world of NATs would be complete without exploring some of their shortcomings. So, after enumerating what are commonly seen as their benefits, it is now necessary to enumerate some of the broken aspects of the world of NATs.

“This solution has the disadvantage of taking away the end-to-end significance of an IP address, and making up for it with increased state in the network.”

—Egevang and Francis,
“Network Address Translator,” RFC 1631

“An opposing view of NAT is that of a malicious technology, a weed which is destined to choke out continued Internet development. While recognizing there are perceived address shortages, the opponents of NAT view it as operationally inadequate at best, bordering on a sham as an Internet access solution. Reality lies somewhere in between these extreme viewpoints.”

—Tony Hain,
“Architectural Implications of NAT,” RFC 2993

- First, NATs cannot support applications where the initiator lies on the “outside.” The external device has no idea of the address of the local internal device, and, therefore, cannot direct any packets to that device in order to initiate a session. This implies that peer-to-peer services, such as voice, cannot work unaltered in a NAT environment.

- The workaround to this form of shortcoming is to force an altered deployment architecture, where service platforms used by external entities are placed “beside” the NAT, allowing command and control from the interior of the local network, and having a permanent (non-NAT) interface to the external network. Obviously this implies some further centralization of IT services within the NATed site.
- Even this approach does not work well for applications such as voice-over-IP, where the “server” now needs to operate as some form of proxy agent. The generic approach here for applications to traverse NATs in the “wrong” direction is for the inside device to forge a UDP connection to the outside agent, and for the inside device to then establish what NAT translated address has been used, and the nature of the NAT in the path, and then republish this address as the local entity’s published service rendezvous point. Sounds fragile? Unfortunately, it is. The other approach is to shift the application to use a set of endpoint identifiers that are distinct from IP addresses, and use a distributed set of “agents” and “helpers” to dynamically translate the application level identifiers into transport IP addresses as required. This tends to create added complexity in application deployment, and also embarks on a path of interdependency that is less than desirable. In summary, workarounds to reestablish a peer-to-peer networking model with NATs tend to be limited, complex, and often fragile.
- The behavior of NATs varies dramatically from one implementation to another. Consequently, it is very difficult for applications to predict or expose the precise behavior of one or more NATs that may exist on the application data path.
- Robust security in IP environments typically operates on an end-to-end model, where both ends include additional information in the packet that can detect attempts to alter the packet in various ways. In IPSec the header part of the packet is protected by the Authentication Header, where an encrypted signature of certain packet header fields is included in the IPSec packet. If the packet header is changed in transit in unexpected ways, the signature check will fail. Obviously IPSec attempts to protect the packet address fields—the very same fields that NATs alter! This leads to the observation that robust security measures and NATs do not mix very well. NATs inhibit implementation of security at the IP level.
- NATs have no inherent failover. NATs are an active in-band mechanism that cannot fail into a safe operating fallback mode. When a NAT goes offline, all traffic through the NAT stops. NATs create a single point where fates are shared in the NAT device maintaining connection state and dynamic mapping information.

- NATs sit on the data path and attempt to process every packet. Obviously bandwidth scaling requires NAT scaling.
- NATs are not backed up by industry-standardized behavior. Although certain NAT-traversal applications make assumptions about the way NATs behave, it is not the case that all NATs necessarily behave in precisely the same way. Applications that work in one context may not necessarily operate in others.
- Multiple NATs can get very confusing with “inside” and “outside” concepts when NATs are configured in arbitrary ways. NATs are best deployed in a strict deployment model of an “inside” being a stub private network and an “outside” of the public Internet. Forms of multiple interconnects, potential loops, and other forms of network transit with intervening NATs lead to very strange failure modes that are at best highly frustrating.
- With NATs there is no clear, coherent, and stable concept of network identity. From the outside these NAT-filtered interior devices are visible only as transient entities.
- Policy-based mechanisms that are based on network identity (for example, *Policy Quality of Service* [QoS]) cannot work through NATs.
- Normal forms of IP mobility are broken when any element behind the NAT attempts to roam beyond its local private domain. Solutions are possible, generally involving specific NAT-related alterations to the behavior of the Home Agent and the mobile device.
- Applications that work with identified devices, or that actually identify devices (such as the *Simple Network Management Protocol* [SNMP] and DNS) require very careful configuration when operating in a NAT environment.
- NATs may drop IP packet fragments in either direction: without complete TCP/UDP headers, the NAT may not have sufficient stored state to undertake the correct header translation.
- NATs often contain ALGs that attempt to be context-sensitive, depending on the source or destination port number. The behavior of the ALGs can be difficult to anticipate, and these behaviors have not always been documented.
- Most NAT implementations with ALGs that attempt to translate TCP application protocols do not perform their functions correctly when the substrings they must translate span across multiple TCP segments; some of them are also known to fail on flows that use TCP option headers, for example timestamps.

From this perspective, NATs are a short-term expediency that is currently turning into a longer-term set of overriding constraints placed on the further evolution of the Internet. Not only do new applications need to include considerations of NAT traversal, but we appear to be entering into a situation where if an application cannot work across NATs, then the application itself fails to gain acceptance. We seem to be locking into a world that is almost the antithesis of the Internet concept. In this NAT-based world, servers reside within the network and are operated as part of the service provider's role, whereas end devices are seen as "dumb" clients, who can establish connections to servers but cannot establish connections between each other. The widespread use of NATs appears to be reinforcing a reemergence of the model of "smart network, dumb clients," whereas others would argue that the network is getting no smarter, it is just that the number of obstacles and amount of network debris is increasing while clients are getting worse at maintaining coherent end-to-end state in the face of such changes.

However, despite their shortcomings, despite the problems NATs create for numerous applications and their users, and despite the continued grappling over a common language to understand how NATs behave, numerous NATs are deployed, and, at least in the IPv4 realm, NATs appear to be a firmly fixed part of the future of the Internet. NATs continue to proliferate in today's Internet.

Moving on with NATs

One commonly held belief is that deployment of IPv6 will eliminate the problem of NATs within the Internet. Certainly it is reasonable to observe that if achieving high address utilization densities is no longer the objective, then there will be plentiful public IPv6 address space and that particular reason to deploy NATs is significantly discounted in an IPv6 realm.

That does not say that IPv6 NATs will not be implemented, nor used. Indeed IPv6 NATs are already available, and they are being used, albeit to some small extent. NATs are, rightly or wrongly, considered to be part of a security solution for a site because of their filtering properties that prevent incoming packets from entering the site unless the NAT already has a permitting binding initiated from the inside. In addition, NATs allow a site to use an internally persistent naming and addressing scheme based on some form of deployment of IPv6 unique *site local* address, and deploy NATs at the edge to create an external view of the site that fits within a provider-based address aggregated view of the IPv6 Internet.

So it would perhaps be too enthusiastic a level of conjecture to suppose that IPv6 will drive away all forms of NAT use in IPv6. It is reasonable to predict that some use of NAT will be seen in IPv6, although many would be highly disappointed if the level of IPv6 NAT use rose to anywhere approaching that of NAT in IPv4.

However, the Internet is still largely a network that uses IPv4 and NATs, and efforts continue along the lines of reducing the amount of friction and frustration in a world in which NATs are prolific. One of the ways to progress here is to treat NAT boxes as yet another instance of Internet middleware, and attempt to apply the same sets of processes to NATs that appear in other instances of middleware. The work of the IETF in the *Middlebox Communication Working Group* uses a model that attempts to expose NATs, as well as firewalls, performance-enhancing proxies, application proxies, and relay agents, to the application, and allows the application to specify the policy that the middlebox should apply. In the case of NATs, this could allow an application to communicate to a NAT that it does not require any form of third-party access, and that a fully symmetric behavior could be applied to the binding without any loss in application functionality. Equally, an application could indicate to the NAT that it expects third parties to be able to use the NAT binding, and that the binding that the NAT will set up for the application should be managed as a port-restricted cone. There is much that could be achieved here that would allow applications to function with some level of determinism, rather than attempting to equip an application with a large and complex toolset of all the relevant techniques of NAT traversal that may be required by the application when confronted by various NAT behaviors.

In the meantime the NAT-behavior guessing game continues. The generic class of techniques that support this function is termed *Unilateral Self-Address Fixing* (UNSAF). This is a process whereby the local entity attempts to determine the address and port by which the entity is known externally, and to determine the characteristics of this association to understand in what contexts the external address may be used as a service rendezvous point for externally initiated communication. Work in this area^[10] has exposed many relevant considerations, including a set of deficiencies noted in the previous section.

So, what would a NAT implementation look like if there were standards relating to NAT behaviors and the implementation were to comply with these standards? Numerous efforts have been made to document various forms of network- and application-friendly ways in which NATs could behave, but it would appear that such an effort will require the imprimatur of a standard in order to attain a level of general acceptance from NAT implementations. However, it is possible to predict that any such effort at a “standardized” form of NAT behavior will include the following considerations. The following set of behaviors is based on that enumerated in^[13]:

- NATs must show endpoint-independent behavior for UDP-based bindings. This is to ensure that the NAT can support application rendezvous without the need for various multiparty relays and agents.
- NAT should not use port preservation nor port overloading, and should operate in a deterministic manner. Port preservation exposes the NATs to nonstandard behaviors when port preservation cannot be enforced. In addition, NATs must have deterministic behavior.

- A dynamic NAT UDP binding timer should be 5 minutes, and should avoid expiration timers of 2 minutes or less. This is to ensure that the timeout is long enough to avoid excessively frequent timer refresh packets.
- The NAT UDP timeout binding must use a timer refresh based on outbound traffic, and all sessions that use a particular binding should use a common refresh timer. This requirement is a security consideration, in that letting inbound traffic refresh the timer allows an external party to keep a port open on the NAT.
- The NAT filtering function should be address dependant. This represents a balance between security and utility.
- The timeout behavior of the NAT UDP filter must be the same as that of the NAT UDP binding timeout. This is intended to reduce the complexity of applications that are reliant on long-held NAT state.
- The NAT should support hairpin connections, using the external address and port.
- If the NAT includes ALG support, the ALGs should be configurable in terms of being able to turn off the ALG function on a per-application basis.
- NATs should support fragmentation and forwarding of packet fragments.
- NATs must support ICMP *Destination Unreachable* messages, and the ICMP timeout should be greater than 2 seconds.

Learning from NATs

At this stage we can observe a few relevant lessons about NATs:

The first is that we need standards and we rely on standards. For many years the IETF has viewed standardization of NATs and their behavior as being an action that would encourage further deployment of a technology that was apparently considered undesirable. The result has been that NATs have been deployed for reasons entirely unconnected with the IETF and standardization, but because the original specification of NAT behavior was at such a general level each NAT implementor has been forced into making local decisions as to how the NAT should behave under specific circumstances. We now enjoy a network with widespread deployment of an active device that does not have consistent implementations and, in the worst cases, exhibits nondeterministic behaviors. This has made the task of deployment of certain applications on the Internet, including voice-based applications, incredibly difficult.

Whether NATs are good or bad, they would be less of a collective headache today if they shared a common standard core behavior. NATs for IPv6 may be considered to be unnecessary today, and it can be argued they represent no real value to an IPv6 site. But a collection of IPv6 NAT implantations with no common core behavior would constitute a far worse problem to application users. Standardization of technology at least eliminates some of the worst aspects of application level guesswork out of technology deployment.

Secondly, a little bit of security is often far worse than no security. NATs are very poor security devices, and in terms of their behavior with UDP, NATs afford only minor levels of protection. The task of securing a site from various forms of attack and disruption remains one of a careful exercise of assessment of acceptable risk coupled with detailed consideration of site-management functions. NATs are not a quick way out of this effort.

In considering NATs it seems that we are back to the very basics of networking. The basic requirements of any network are “who,” “where,” and “how,” or “identity,” “location,” and “forwarding.” In the case of IP, all these elements were included in the semantics of an IP address, and when addresses get translated dynamically we lose track of IP-level identity across the network. Maybe, just maybe, as we look at the longer-term developments of IP technology, one potential refinement may be the separation of endpoint identity to that of location, and as a potential outcome, NATs could readily manipulate location-based addresses while applications could look to a different token set as a means of establishing exactly who is the other party to the communications.

Of course, if we ever venture down such a path, I trust that such a move toward the use of explicit identities does not generate a complementary deployment of *Network Identity Translators*, or NITs, as an adjunct to the current set of NATs. Too many NITs and NATs will definitely send us all NUTs!

Further Reading

There is no shortage of material on NATs from a wide variety of sources. The following is a list of IETF-related documents, encompassing both published *Request for Comments* (RFCs) and works in progress, that have been circulated as *Internet Drafts*.

RFCs:

- [1] Egevang, K., and P. Francis, “The IP Network Address Translator (NAT),” RFC 1631, May 1994.
- [2] Srisuresh, P., and D. Gan, “Load Sharing Using IP Network Address Translation (LSNAT),” RFC 2391, August 1998.

- [3] Srisuresh, P., and M. Holdrege, “IP Network Address Translator (NAT) Terminology and Considerations,” RFC 2663, August 1999.
- [4] Tsirtsis, G., and P. Srisuresh, “Network Address Translation—Protocol Translation (NAT-PT),” RFC 2776, February 2000.
- [5] Hain, T., “Architectural Implications of NAT,” RFC 2993, November 2000.
- [6] Srisuresh, P., and K. Egevang, “Traditional IP Network Address Translator (Traditional NAT),” RFC 3022, January 2001.
- [7] Holdrege, M., and P. Srisuresh, “Protocol Complications with the IP Network Address Translator,” RFC 3027, January 2001.
- [8] D. Senie, “Network Address Translator (NAT)-Friendly Application Design Guidelines,” RFC 3235, January 2002.
- [9] Srisuresh, P., J. Kuthan, J. Rosenberg, A. Molitor, and A. Rayhan, “Middlebox Communication Architecture and Framework,” RFC 3303, August 2002.
- [10] Daigle, L., and IAB, “IAB Considerations for Unilateral Self-Address Fixing (UNSAF) Across Network Address Translation,” RFC 3424, November 2002.
- [11] Rosenberg, J., Weinberger, J., Huitema, C., and R. Mahy, “STUN—Simple Traversal of User Datagram Protocol (UDP) Through Network Address Translators (NATs),” RFC 3489, March 2003.
- [12] Aboba, B., and W. Dixon, “IPsec—Network Address Translation (NAT) Compatibility Requirements,” RFC 3715, March 2004.

Internet Drafts:

Internet Drafts enjoy a fleeting existence, and the following documents may not be available when you read this article. In such cases it is often the case that a decent Internet search will locate the document, or its successor.

- [13] Audet, F., and C. Jennings, “NAT/Firewall Behavioral Requirements,” work in progress, **draft-audet-nat-behave**, July 2004.
- [14] Ford, B., P. Srisuresh, and D. Kegel, “Peer-to-Peer(P2P) Communication across Network Address Translators (NATs),” work in progress, **draft-ford-midcom-p2p**, June 2004.

- [15] Rosenberg, J., “Interactive Connectivity Establishment (ICE): A Methodology for Network Address Translator (NAT) Traversal for the Session Initiation Protocol (SIP),” work in progress, **draft-ietf-mmusic-ice**, July 2004.
- [16] Jennings, C., “NAT Classification Results Using STUN,” work in progress, **draft-jennings-midcom-stun-results**, July 2004.
- [17] J. Rosenberg, J. Weinberger, R. Mahy, and C. Huitema, “Traversal Using Relay NAT (TURN),” work in progress, **draft-rosenberg-midcom-turn-01**, July 2004.

Other Resources:

NAT Check: Ford, B. and D. Andersen, Nat Check Website:
<http://midcom-p2p.sourceforge.net>

STUN Client and Server:
<http://sourceforge.net/projects/stun>

Phifer, Lisa, “The Trouble with NAT,” *The Internet Protocol Journal*, Volume 3, No. 4, December 2000.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector, and has served his time with Telstra, where he was the Chief Scientist in their Internet area. Geoff is currently the Internet Research Scientist at the *Asia Pacific Network Information Centre* (APNIC). He is also the Executive Director of the Internet Architecture Board, and is a member of the Board of the Public Interest Registry. He is author of *The ISP Survival Guide*, ISBN 0-471-31499-4, *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, and co-author of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. All three books are published by John Wiley & Sons. E-mail: **gih@apnic.net**

Letters to the Editor

Content Networks

Dear Editor,

Christophe Deleuze's article on Content Networks (*The Internet Protocol Journal*, Volume 7, Number 2, June 2004) made me realize that there are very different ways to look at this issue. I would like to use the term *Content Addressable Network* for a network that is used to retrieve information not by specifying its location but the identity of the content itself. The term points to similar concepts in electronics (*Content Addressable Memory*) and storage (*Content Addressable Storage*). One could argue that a Content Addressable Network is in fact a distributed Content Addressable Storage.

In a very real sense the Internet already is content addressable. Several of my non-IT friends use the "Search" field in the *Google* toolbar even for regular URLs, foregoing the Address field in their browsers. In doing so, they simply ignore the distinction between *content* and *location*. It usually gets them where they want to go.

Let's define content as a static binary object, for example, a document, picture, song, or movie. How can we identify content if not by location? We can create a hash of the object as a handle or placeholder. (A hash is the result of a calculation that takes the whole object as input. A good hashing algorithm ensures that if you change a bit in the object, at least one bit in its hash changes too.) If we know the placeholder, we can retrieve a copy of the original object, even if we don't know the location of any of the copies out there on the net. I could mail you the hash of a paper, song, or movie and you would be able to retrieve a copy, although not necessarily from the same place as where I got it. (You might have to pay to get it though!)

Suppose that the Google bot, while traversing the Internet to build its index, calculates the hash for each object it encounters. It can then build an index of all hash codes, relating them to the URLs where they were found. (This requires no change in Google: the hash is just one more word it found in the document). We can then google a hash code to find all occurrences of the object. (You can simulate this today by selecting a line of text from a document and launching a search for that sequence of words. Google will often find multiple copies. Just one line of text is an extremely poor hash, so you may get a few false hits, but in my experience not many.)

Simply by adding these hashes, we have turned the Internet into a Content Addressable Network. If our purpose is to make ourselves independent of any single copy on any particular server, this is all we need. For other applications, the objective is to optimize the network paths to the servers that hold a copy of our object (for example, a movie). We need a metric that tells us which of the listed locations is "closest" to our point of entry. This is complicated by the fact that the Internet is a weird space. The shortest route between Amsterdam and Brussels might well go via London or Paris.

Fortunately, there is a database that keeps track of all the available routes and their cost. It is the *Border Gateway Protocol* (BGP) routing table. BGP divides the Internet in chunks called *Autonomous Systems* or ASs and tracks the cost of the routes to each AS. If the Google bot would record the AS along with each URL, our client system could query our local BGP router (or a proxy holding a copy of its database) to find the AS and thus the copy that is closest in terms of network costs. Note that these costs also reflect policy rules such as peering arrangements between ISPs.

If our objective is to dynamically optimize the load on the servers, we cannot avoid querying (a local subset of) these servers for a bid. Distributing the load over servers in different time zones may sometimes be more important than keeping the transports local. Our client should select a server that is not too busy but no further away than necessary.

The Content Networks as discussed by Christophe Deleuze were created as a commercial offering that would require no cooperation from the clients—in every sense an operator’s approach. It is restricted to the case where all copies of the object are published by a single entity. The way ahead is to create protocols for requesting network cost for a list of sites, and service costs from a list of servers, independent of the nature of the object and the servers that hold copies of it.

It may seem more efficient to let the publisher add the hash code to the objects. HTML files would be labeled with a `<MD5=` tag, obviating the need for bots and users (for “content bookmarks”) to do the calculation. This would allow publishers to change content without changing the hash, to correct typos or remove scenes deemed unsuitable for local viewers. But it would no doubt result in fake objects, purporting to be copies of popular objects but peddling dubious commercial proposals. Creating fake objects is more difficult if the hash code is calculated by an independent and unrecognizable bot, although I’m sure the problem is not completely solved with that.

—Ernst Lopes Cardozo, *Aranea Consult BV, The Netherlands*
`e.lopes.cardozo@aranea.nl`

I noticed that the IPJ page footer only says “The Internet Protocol Journal” but neither the Volume/Issue number, nor the issue date. That makes it a bit hard to correctly reference a given article when you only have a copy of that article and not the whole issue. I propose that you add something like (from the August issue of CACM):

Communications of the ACM August 2004/Vol. 47, No. 8

(I only checked the archived PDF files but I suppose the hardcopy has the same problem.)

—Örjan Petersson
orjan.petersson@logcode.com

We could certainly add the Volume/Issue identifier to the footer, but since this would have to be done retroactively for all 26 issues to date it is probably better to use our soon-to-be-deployed ASCII index. This will allow you to find any article with a simple search. A short sample of the index is shown below.

The Internet Protocol Journal Volume 1, 1998

Article	Author(s)	Page

* Volume 1, No. 1, June 1998:		
What Is a VPN? - Part I	Ferguson/Huston	2
SSL: Foundation for Web Security	William Stallings	20
Book Review: Groupware	Dave Crocker	31
Book Review: High-Speed Networks	Neophytos Iacovou	33
* Volume 1, No. 2, September 1998:		
What Is a VPN? - Part II	Ferguson/Huston	2
Reliable Multicast Protocols and Applications	C. Kenneth Miller	19
Layer 2 and Layer 3 Switch Evolution	Thayumanavan Sridhar	38
Book Review: Gigabit Ethernet	Ed Tittel	44
* Volume 1, No. 3, December 1998:		
Security Comes to SNMP: SNMPv3	William Stallings	2
CATV Internet Technology	Mark Laubach	13
Digital TV	George Abe	27
I Remember IANA	Vint Cerf	38
Book Review: Internet Messaging	Dave Crocker	40
Book Review: Web Security	Richard Perlman	42
Book Review: Internet Cryptography	Frederick M. Avolio	44

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

IPv6 Address “Glue” added to the Root DNS Zone

The *Internet Corporation For Assigned Names and Numbers* (ICANN) recently announced that for the first time, an IPv6 nameserver address has been added to the Internet’s root DNS zone. This next generation version of the Internet Protocol provides trillions more addresses than the IPv4 system that is in use by most networks today. By taking this significant step forward in the transition to IPv6, ICANN is supporting the innovations through which the Internet evolves to meet the growing needs of a global economy.

On 20 July 2004 at 18:33 UTC the IPv6 AAAA records for the Japan (**.jp**) and Korea (**.kr**) *country code Top Level Domain* (ccTLD) nameservers became visible in the root zone file with serial number 2004072000. It is expected that the IPv6 records for France (**.fr**) will be added shortly. Other requests are pending and will be added in accordance with documented procedure, which was developed through ICANN’s unique multi-stakeholder consensus-based approach. See: <http://www.iana.org/procedures/delegation-data.html>

Recognizing the importance of IPv6 to the Internet community, ICANN has coordinated with its *Root Server System Advisory Committee*, *Top Level Domain* managers, *Security and Stability Advisory Committee*, and other interested parties in careful analysis of this issue. After a period of thorough examination, the decision was made to move forward with deployment of the IPv6 address records in the manner prescribed by the community.

ICANN is the global public-benefit non-profit organization responsible for coordinating the Internet’s naming and numbering systems. For more information please visit: <http://www.icann.org>

Formation of Asia Pacific ENUM Engineering Team

China Network Information Center (CNNIC), *Japan Registry Service* (JPRS), *Korea Network Information Center* (KRNIC), *Singapore Network Information Center* (SGNIC) and *Taiwan Network Information Center* (TWNIC) recently announced the formation of the *Asia Pacific ENUM Engineering Team* (APEET), an informal technical project team formed to coordinate and synergize ENUM activities in the Asia Pacific region.

The proposal to form APEET was discussed during an ENUM BoF (Birds-of-a-Feather) session at the *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) in February 2004. Founding member organizations of APEET shared a common vision that as a collective group, they will be able to achieve greater community awareness and better interoperability of ENUM-based trials.

“ENUM allows IP devices to be assigned a telephone number which is globally interoperable,” said James Seng, Chairman of APEET. “It is a key enabling technology for seamless IP Telephony that will greatly benefit the end-users.”

Before the formation of APEET, each member organization has been conducting its own ENUM trials, most of which are isolated trials conducted within each member organization's country/region. With the formation of APEET, member organizations will be able to implement technical solutions that facilitate ENUM trials across Asia Pacific.

"We are extremely excited about the formation of this much needed organization," said Hiro Hotta, Director JPRS. "We are ready to bring ENUM trials to the next level."

One of APEET's key project is to implement a live ENUM trial at APRICOT 2005, Kyoto, Japan. The live trial will allow hundreds of APRICOT participants to experience IP Telephony using wireless SIP Phones and calling each another with standard 10-key telephone interface via ENUM. The live trial, believed to be the first of its kind, will serve to demonstrate and educate the technical community on the power, capabilities and feasibility of ENUM together with SIP.

"This looks like one of the most exciting events of 2005 with a demonstration of technologies to rock Asia Pacific," said Richard Shockey, co-Chair of the ENUM Working Group of the IETF.

The formation of APEET has been well received by the Industry. The *Asia Pacific Network Information Centre* (APNIC) has extended its goodwill to host DNS records of **apenum.org**, the selected "golden root" of APEET technical trials. APEET is also fortunate to have individual experts member such as Richard Shockey.

APEET welcomes all Asia Pacific ccTLD administrators (or its designated representatives) to join and contribute towards the success of ENUM adoption in Asia Pacific. For more information, please visit <http://www.apenum.org>

Phill Gross Receives Postel Award

Phill Gross is this year's recipient of the prestigious *Jonathan B. Postel Service Award*. A co-founder of the *Internet Engineering Task Force* (IETF), Gross has been instrumental in defining and shaping the way in which the IETF standards process functions. He was awarded the Postel Service Award in recognition of his early leadership of the IETF and for firmly establishing the principles that are essential for its success. The Postel Award was presented on August 5th, during the 60th meeting of the IETF in San Diego, California.

"The Internet Society is pleased to recognize Phill's significant contribution to the area of Internet standardization by awarding him this year's Postel Award," said Internet Society President and CEO Lynn St. Amour. "The continued success of the IETF's consensus-based processes shows the importance of Phill's pioneering work in developing the IETF's foundations."

According to Steve Crocker, noted Internet authority and chair of this year's Postel award committee, "Many of the IETF's current structures, including Working Groups, Technical Areas, Proceedings and Internet Drafts came about thanks to Phill's dedication and passion for the Internet standards area. And we're delighted to be presenting the award to Phill in San Diego, the location of the first ever IETF meeting back in 1986."

Gross, who is currently Director of Academics and Technology for the Northern Virginia ECPI College of Technology, has worked with the Internet community for over 20 years. His career has taken him from working with government-funded research projects through to networking engineering responsibilities for large corporations and startups, including leading the development of MCI Corporation's first national network.

In 1986 Gross helped found the IETF. He became the first official chair in 1987—a position he held for seven years. During his chairmanship, the IETF evolved from a government-sponsored research group to an industry-wide Internet standards body. As well as contributing to developing the IETF standards process itself, Gross played an active role as co-chair of the IETF Routing and Addressing Working Group. This group led to solutions for growth-related Internet problems and was instrumental in specifying the initial direction for the next generation *Internet Protocol* (IPv6) in RFC 1719. He also served as a member of the *Internet Architecture Board* (IAB) from 1987 to 1996.

Expressing his appreciation for the award, Gross said "It was very gratifying to be there at the beginning and to work with such an incredible group of people. And, working with Jon over the years gives me a special appreciation for the honor that comes with this award."

The Jonathan B. Postel Service Award was established by the Internet Society to honor those who have made outstanding contributions in service to the data communications community. The award is focused on sustained and substantial technical contributions, service to the community, and leadership. With respect to leadership, the nominating committee places particular emphasis on candidates who have supported and enabled others in addition to their own specific actions. The award is named after Dr. Jonathan B. Postel, who embodied all of these qualities during his extraordinary stewardship over the course of a thirty-year career in networking. He served as the editor of the RFC series of notes from its inception in 1969, until 1998. He also served as the ARPANET "Numbers Czar" and the *Internet Assigned Numbers Authority* (IANA) over the same period of time. He was a founding member of the Internet Architecture Board and the first individual member of the Internet Society, where he also served as a trustee. Previous recipients of the Postel Award include Jon himself (posthumously and accepted by his mother), Scott Bradner, Daniel Karrenberg, Stephen Wolff and Peter Kirstein. For more information, please visit:

<http://www.isoc.org>

Where did my copy of IPJ go?

Each time we mail out a new issue of IPJ, a certain number of copies are returned to us as undeliverable by the postal authorities around the globe. These so-called “Nixies” can take as much as a year to arrive back in San Jose, California, and almost all of them are returned without any updated delivery information. Obviously we cannot do much other than delete these records from our database. However, if you tell us when you move, we can make sure your address is up-to-date so that you will receive the next issue of IPJ. You can update your own record using the subscription tool at <http://www.cisco.com/ipj> or just send your updates via e-mail to: ipj@cisco.com

Where did you go?
Do let us know!

<p>退回 RETOUR 邮 2608-CN15 (舊 C 33/CP10) (ancien C 33/CP10)</p> <p><input checked="" type="checkbox"/> 寄件人 Inconnu <input type="checkbox"/> 拒收 Refusé</p> <p><input type="checkbox"/> 迁居 Déménagé <input type="checkbox"/> 无人认领 Non réclamé</p> <p><input type="checkbox"/> 地址不详 Adresse insuffisante</p> <p><input type="checkbox"/> Inconnu <input type="checkbox"/> Refusé</p> <p><input checked="" type="checkbox"/> Déménagé <input type="checkbox"/> Non réclamé</p> <p><input type="checkbox"/> Adresse insuffisante</p> <p><input type="checkbox"/> Adresse postale changée</p> <p><input type="checkbox"/> Nouvelle adresse N'y address</p>	<p>UNDELIVERED ONAFGELEVERD NON-DISTRIBUE</p> <p><input checked="" type="checkbox"/> UNKNOWN UNBEKANT INCORNU A CETTE ADRESSE</p> <p><input type="checkbox"/> NO SUCH NUMBER NO SO N'NUMERO NIE NIENIG NUMERO N'EXISTE PAS</p> <p><input type="checkbox"/> NO SUCH STREET NO SO N'STRAAT NIE NIENIG STRAAT N'EXISTE PAS</p> <p><input type="checkbox"/> ADDRESS UNKNOWN ADRESSE ONBEPALDENDE ADRESSE INCOMPLETE</p> <p><input type="checkbox"/> ADDRESS ILLEGIBLE ADRESSE ONLEESBAAR ADRESSE ILLISIBLE</p> <p><input type="checkbox"/> UNCLAIMED ONVERKLAARDE EN BOVENFRANCE</p> <p><input type="checkbox"/> REFUSED ONTREFUSE</p> <p><input type="checkbox"/> BOX CLOSED BOVENSLUIT BOITE FERMEE</p> <p><input type="checkbox"/> GONE AWAY - NO ADDRESS LEFT VERSTRAK - GEEN ADRESSE LAAT NIE PARTI SANS L'ADRESSE</p> <p>COMPLETED BY / VOLTOOI DEUR / COMPLETE PAR</p> <p>NAME NAAM NOM</p> <p>RETURN CHARGE PAYABLE TERUGSENDINGSCHIEF RETALJAAR PRIJS DE RETOUR PAYABLE</p> <p>R</p> <p>DATE DATUM DATE</p> <p>WT 2463 701387</p>	<p>RETURN TO SENDER</p> <p><input type="checkbox"/> No such Street/Number</p> <p><input type="checkbox"/> NonIdentically Addressed</p> <p><input checked="" type="checkbox"/> Unknown at Address</p> <p><input type="checkbox"/> Left Address</p> <p><input type="checkbox"/> Refused</p> <p><input type="checkbox"/> Box / Bag Cancelled</p> <p><input type="checkbox"/> Unclaimed</p> <p>8838221 (JUL06)</p> <p>N'habite pas à l'adresse indiquée. Retour à l'expéditeur QL 07</p> <p>Adresse und Irrfahrten-/Postfach: Anschrift stimmen nicht überein</p> <p>Adresse de l'envoi et de la boîte aux lettres postale ne concordent pas</p> <p>Indirizzo e intestazione della buca postale/cassa postale non corrispondono</p> <p>F DIE POST LA POSTE LA POSTA</p>
<p>RETOUR</p> <p>Nicht (mehr) unter dieser Adresse</p> <p>RETURN TO SENDER</p> <p>NOT (ANY) LONGER AT THIS ADDRESS</p>	<p>Royal Mail</p> <p>We were unable to deliver this item because</p> <p><input checked="" type="checkbox"/> addressee has gone away</p> <p><input type="checkbox"/> no answer <input type="checkbox"/> addressee unknown</p> <p><input type="checkbox"/> address incomplete <input type="checkbox"/> refused</p> <p><input type="checkbox"/> address inaccessible <input type="checkbox"/> not called for</p> <p>no such address in</p> <p>date initials</p> <p>badge number P3966/97302943</p>	<p>Zurück/Retour CN 15</p> <p>Empfänger/Firma unter der angegebenen Anschrift nicht zu errichten</p> <p><input checked="" type="checkbox"/> Inconnu/Adresse insuffisante</p> <p><input type="checkbox"/> Déménagé</p> <p>Empfänger verzogen, Einwilligung zur Weitergabe der neuen Anschrift liegt nicht vor</p> <p>Annahme verweigert <input type="checkbox"/> Refusé</p> <p>Nicht abgeholt <input type="checkbox"/> Non réclamé</p> <p>Nicht zulässig <input type="checkbox"/> Non admis</p> <p>Rücksendung am/Retour le: <i>ATC</i></p>

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Technology Strategy
MCI, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.
Copyright © 2004 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSR STD U.S. Postage PAID Cisco Systems, Inc.
--

The Internet Protocol Journal

December 2004

Volume 7, Number 4

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Network Processors	2
Denial of Service Attacks	13
Letter to the Editor	36
Book Review	37
Call for Papers	39

FROM THE EDITOR

The electronics industry is full of examples of devices which contain one or two “special-purpose” chips. Your computer probably has a modem that is implemented with a single chip and a few analog components. It probably also contains a dedicated graphics processor responsible for driving your display. In networking, vendors have long since realized that in order to design highly efficient routers or switches, a custom-designed *network processor* is a good solution. We asked Doug Comer to give us an overview of network processors.

Attacks against individual computers on a network have become all too common. Usually these attacks take the form of a virus or worm which arrives via e-mail to the victim’s machine. The industry has been relatively quick in responding to such attacks by means of antivirus software, as well as sophisticated filtering of content “on the way in.” A more serious form of attack is the *Distributed Denial-of-Service* (DDoS) attack which may render an entire network unusable. Charalampos Patrikakis, Michalis Masikos, and Olga Zouraraki give an overview of the many variants of denial-of-service attacks and what can be done to prevent them.

Although we make every effort to provide you with an error-free journal, mistakes do happen occasionally. Sometimes it takes careful analysis by a reader to spot the mistake, and we are grateful for the correction provided in the “Letter to the Editor” on page 36. Other times, technology just gets in our way, such as when all the non-printing end-of-line and TAB characters became very much “printing”—see page 35 of the printed version of Volume 7, No. 3. At least it didn’t show up in the PDF or HTML versions.

Take a moment to visit our Website: <http://www.cisco.com/ipj> and update your mailing address if necessary. You will also find all back issues and index files at the same address.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Network Processors: Programmable Technology for Building Network Systems

by Douglas Comer, Cisco Systems (on leave from Purdue University)

Chip vendors have defined a new technology that can be used to implement packet-processing systems such as routers, switches, and firewalls. The technology offers the advantages of being software-programmable and sufficiently high-speed to accommodate interfaces running at 10 Gbps.

This article provides an overview of the technology, describes the motivations, and presents a brief survey of hardware architectures. It also discusses the relationship between programming and the underlying hardware.

A wide variety of packet-processing systems are used in the Internet, including DSL modems, Ethernet switches, IP routers, *Network Address Translation* (NAT) boxes, *Intrusion Detection Systems* (IDS), Soft-switches used for *Voice over IP* (VoIP), and security firewalls. Such systems are engineered to provide maximal functionality and performance (for example, operate at wire speed) while meeting constraints on size, cost, and time to market.

Engineers who design network systems face the additional challenges of keeping designs scalable, general, and flexible. In particular, because industry trends change rapidly, typical engineering efforts must accommodate changes in requirements during product construction and changes in the specification for a next-generation product.

Generations of Network Systems

During the past 20 years, engineering of network systems has changed dramatically. Architectures can be divided broadly into three generations:

- *First generation* (circa 1980s): Software running on a standard processor (for example, an IP router built by adding software to a standard minicomputer),
- *Second generation* (mid 1990s): Classification and a few other functions offloaded from the CPU with special-purpose hardware, and a higher-speed switching fabric replacing a shared bus.
- *Third generation* (late 1990s): Completely decentralized design with *Application-Specific Integrated Circuit* (ASIC) hardware plus a dedicated processor on each network interface offloading the CPU and handling the fast data path.

The change from a centralized to a completely distributed architecture has been fundamental because it introduces additional complexity. For example, in a third-generation IP router, where each network interface has a copy of the routing table, changing routes is difficult because all copies must be coordinated to ensure correctness and the router should not stop processing packets while changes are propagated.

Motivation for Network Processors

Although the demand for speed pushed engineers to use ASIC hardware in third-generation designs, the results were disappointing. First, building an ASIC costs approximately US\$1 million. Second, it takes 18 to 22 months to generate a working ASIC chip. Third, although engineers can use software simulators to test ASIC designs before chips are manufactured, networking tasks are so complex that simulators cannot handle the thousands of packet sequences needed to verify the functionality. Fourth, and most important, ASICs are inflexible.

The inflexibility of ASICs impacts network systems design in two ways. First, changes during construction can cause substantial delay because a small change in requirements can require massive changes in the chip layout. Second, adapting an ASIC for use in another product or the next version of the current project can introduce high cost and long delays. Typically, a silicon respin takes an additional 18 to 20 months.

Network-Processor Technology

In the late 1990s as demand for rapid changes in network systems increased, chip manufacturers began to explore a new approach: programmable processors designed specifically for packet-processing tasks. The goal was clear: combine the advantage of software programmability, the hallmark of the first-generation network systems, with high speed, the hallmark of third-generation network systems.

Chip vendors named the new technology *network processors*, and predicted that in the future, most network systems would be constructed using network processors. Of course, before the prediction could come true, vendors faced a tough challenge: programming introduces an extra level of indirection, meaning that functionality implemented directly in hardware always performs faster than the same functionality implemented with software. Thus, to make a network processor fast enough, packet-processing tasks need to be identified and special-purpose hardware units constructed to handle the most intensive tasks.

Interestingly, vendors also face an economic challenge: although an ASIC costs a million dollars to produce, subsequent copies of the chip can be manufactured at very low cost. Thus, the initial development cost can be amortized over many copies. In contrast, purchasing conventional processors does not entail any initial development cost, but vendors typically charge at least an order of magnitude more per unit than for copies of an ASIC. So, vendors must consider a pricing strategy that entices systems builders to use network processors in systems that have many network interfaces with multiple processors per interface.

A Plethora of Architectures

As vendors began to create network processors, fundamental questions arose. What are the most important protocol-processing tasks to optimize? What hardware units should a network processor provide to increase performance? What I/O interfaces are needed? What sizes of instruction store and data store are needed? What memory technologies should be used (for example, *Static Random-Access Memory* [SRAM], *Dynamic Random-Access Memory* [DRAM], or others)? How should functional units on the network-processor chip be organized and interconnected (for example, what on-chip bus infrastructure should be used)?

Interestingly, although they realized that it was essential to identify the basic protocol-processing tasks before hardware could be built to handle those tasks efficiently, chip vendors had little help from the research community. Much effort had been expended considering how to implement specific protocols such as IP or TCP on conventional processors. However, researchers had not considered building blocks that worked across all types of network systems and all layers of the protocol stack. Consequently, in addition to designing network-processor chips, vendors needed to decide which protocol functions to embed in hardware, which to make programmable, and which (if any) to leave for special-purpose interface chips or coprocessors. Finally, chip vendors needed to choose software support including programming language(s), compilers, assemblers, linkers, loaders, libraries, and reference implementations.

Faced with a myriad of questions and possibilities about how to design network processors and the recognition that potential revenue was high if a design became successful, chip vendors reacted in the expected way: each vendor generated a design and presented it to the engineering community. By January 2003, more than 30 chip vendors sold products under the label “network processor.”

Unfortunately, the euphoria did not last, and many designs did not receive wide acceptance. Thus, companies began to withdraw from the network-processor market, and by January 2004, fewer than 30 companies sold network processors.

Basic Architectural Approaches

Hardware engineers use three basic techniques to achieve high-speed processing: a single processor with a fast clock rate, parallel processors, and hardware pipelining. Figure 1 illustrates packet flow through a single processor, which is known as an *embedded processor architecture* or a *run-to-completion model*. In the figure, three functions must be performed on each packet.

Figure 1: Embedded Processor Architecture in Which a Single Processor Handles All Packets

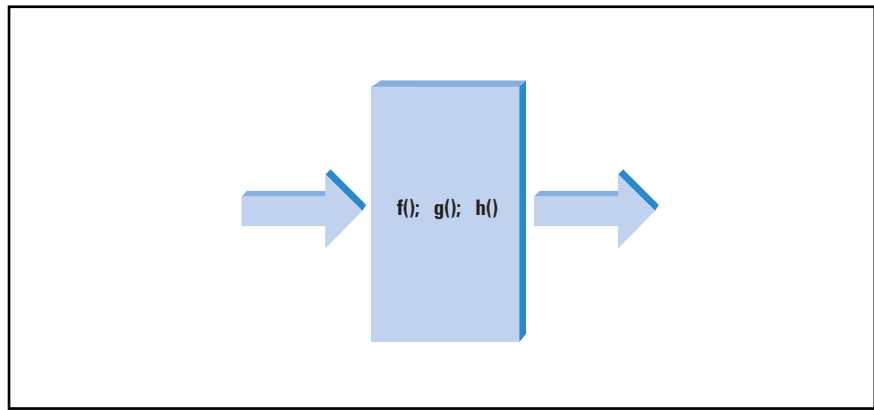


Figure 2 illustrates packet flow through an architecture that uses a parallel approach. A coordination mechanism on the ingress side chooses which packets are sent to which processor. Coordination hardware can use a simplistic round-robin approach in which a processor receives every N th packet, or a sophisticated approach in which a processor receives a packet whenever the processor becomes idle.

Figure 2: Parallel Architecture in Which the Incoming Packet Flow Is Divided Among Multiple Processors

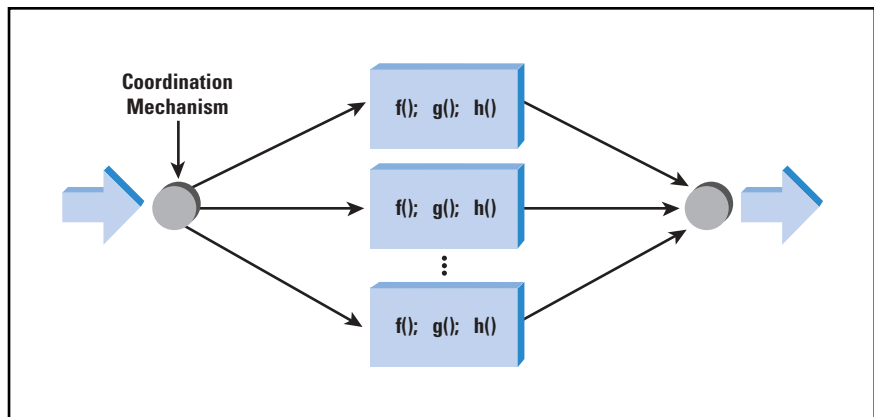
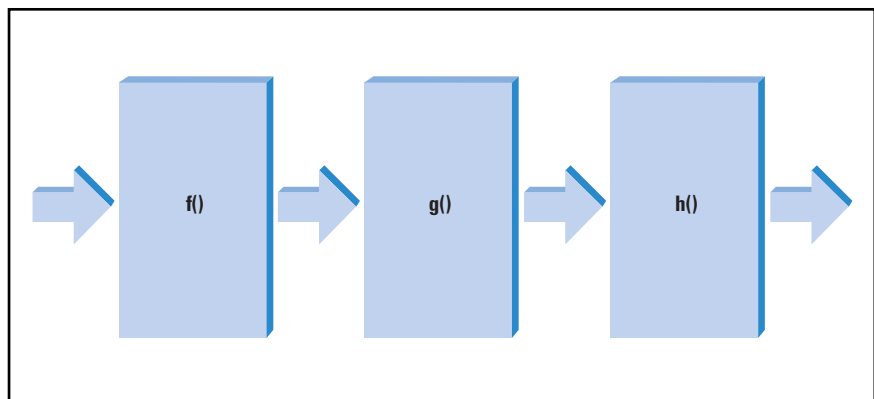


Figure 3 illustrates packet flow through a pipeline architecture. Each packet flows through the entire pipeline, and a given stage of the pipeline performs part of the required processing.

Figure 3: Pipeline Architecture in Which Each Incoming Packet Flows Through Multiple Stages of a Pipeline



As we will see, pipelining and parallelism can be combined to produce hybrid designs. For example, it is possible to have a pipeline in which each individual stage is implemented by parallel processors or a parallel architecture in which each parallel unit is implemented with a pipeline.

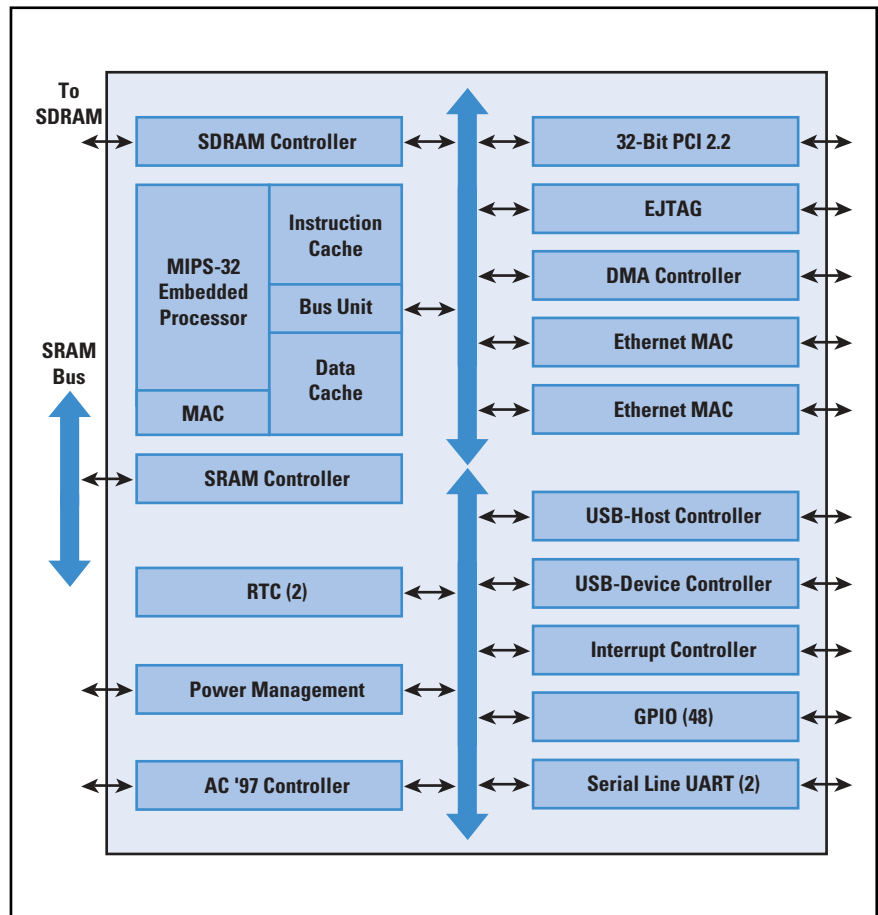
Examples of Commercial Architectures

To appreciate the broad range of network-processor architectures, we will examine a few commercial examples. Commercial network processors first emerged in the late 1990s, and were used in products as early as 2000. The examples contained in this article are chosen to illustrate concepts and show broad categories, not to endorse particular vendors or products. Thus, the examples are not necessarily the best, nor the most current.

Augmented RISC (Alchemy)

The first example, from Alchemy Semiconductor (now owned by Advanced Micro Devices), illustrates an embedded processor augmented with special instructions and I/O interfaces.

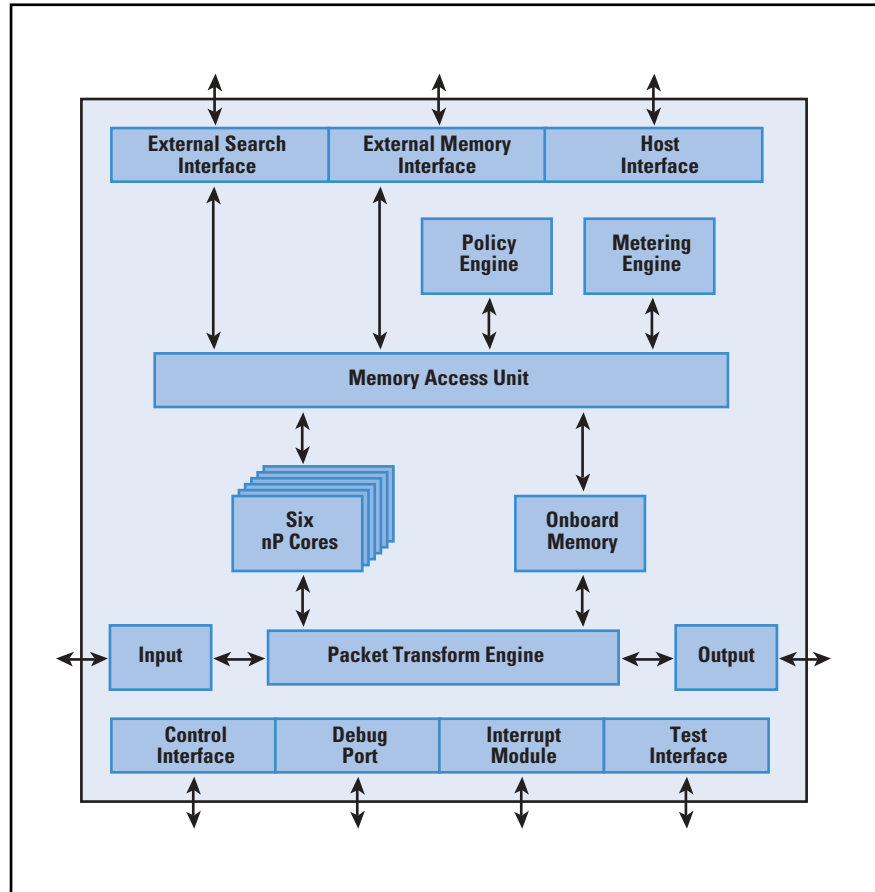
Figure 4: An Example Embedded Processor Architecture: The Processor Has Extra Instructions to Speed Packet Processing



Parallel Processors Plus Coprocessors (AMCC)

A network processor from AMCC uses an architecture with parallel processors plus coprocessors that handle packet-processing tasks. When a packet arrives, one of the parallel processors, called *cores*, handles the packet. The coprocessors are shared—any of the parallel processors can invoke a coprocessor, when needed.

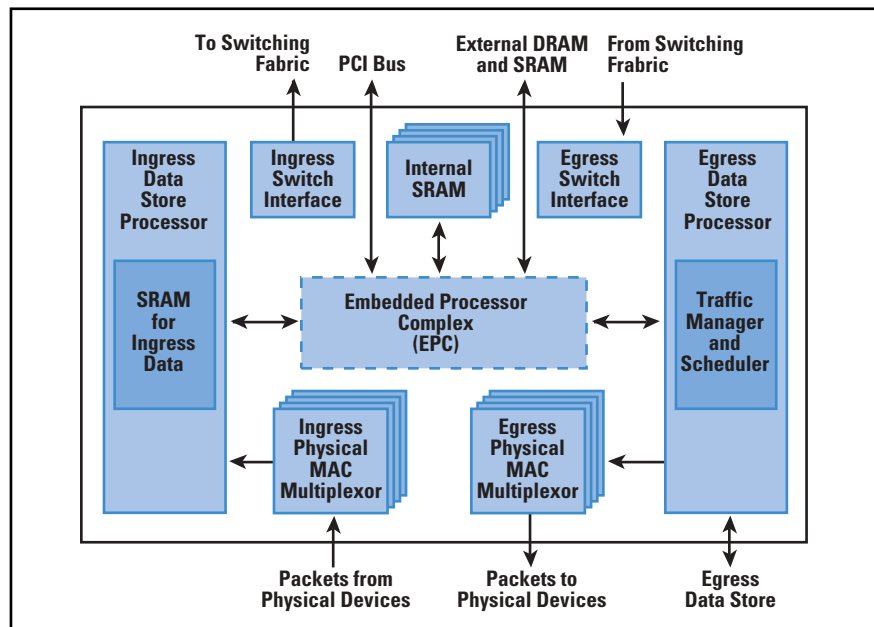
Figure 5: An Example Parallel Architecture that Uses Special-Purpose Coprocessors to Speed Execution



Extensive and Diverse Processors (Hifn)

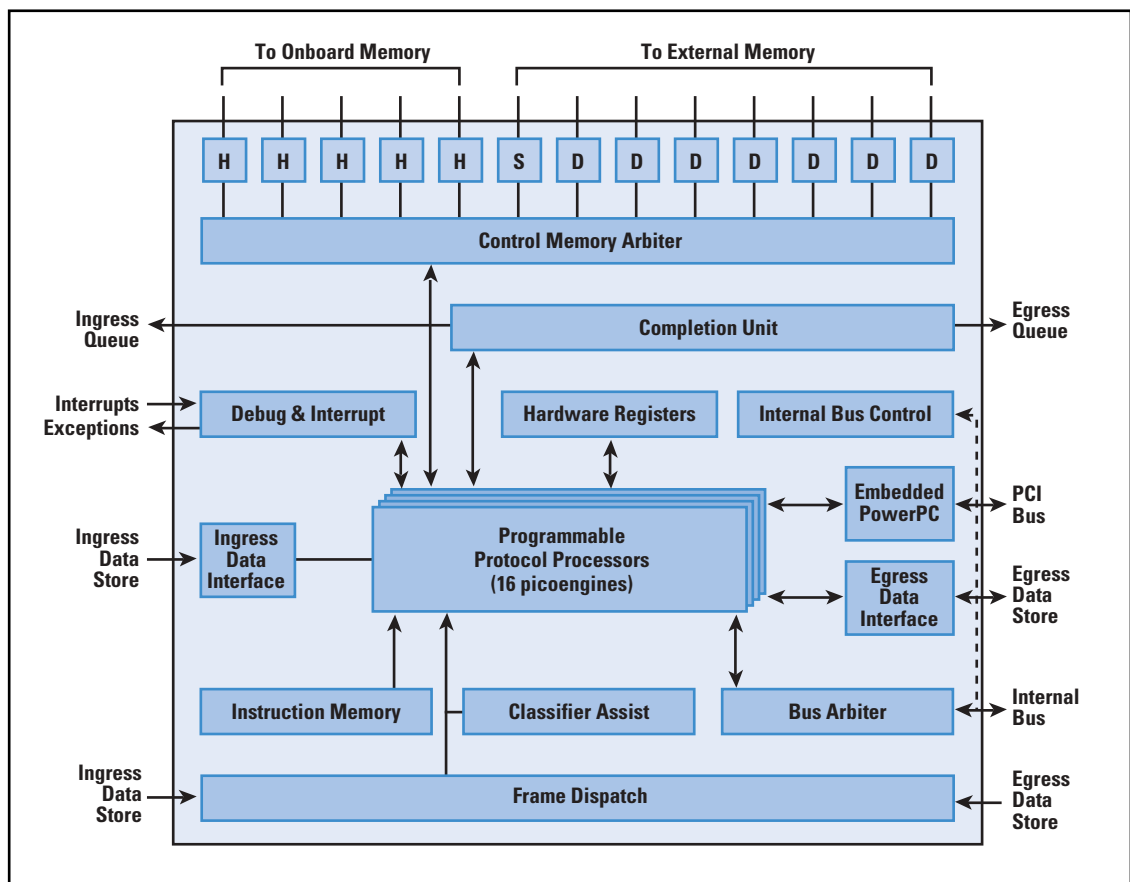
A network processor (named *Rainier*) originally developed by IBM and now owned by Hifn Corporation uses a parallel architecture, and includes a variety of special-purpose and general-purpose processors. For example, the chip provides parallel ingress and egress hardware to handle multiple high-speed network interfaces. It also has intelligent queue-management hardware that enqueues incoming packets in an ingress data store, a switching fabric interface built onto the chip, and an intelligent egress data store. Figure 6 illustrates the overall architecture of the Hifn chip.

Figure 6: An Example Parallel Architecture that Includes Hardware Support for Ingress and Egress Processing as well as Intelligent Queuing



The *Embedded Processor Complex* (EPC) on the Hifn chip contains 16 programmable packet processors, called *picoengines*, as well as various other coprocessors. In addition, the EPC contains an embedded PowerPC to handle control and management tasks. Figure 7 shows a few of the many processors in the EPC.

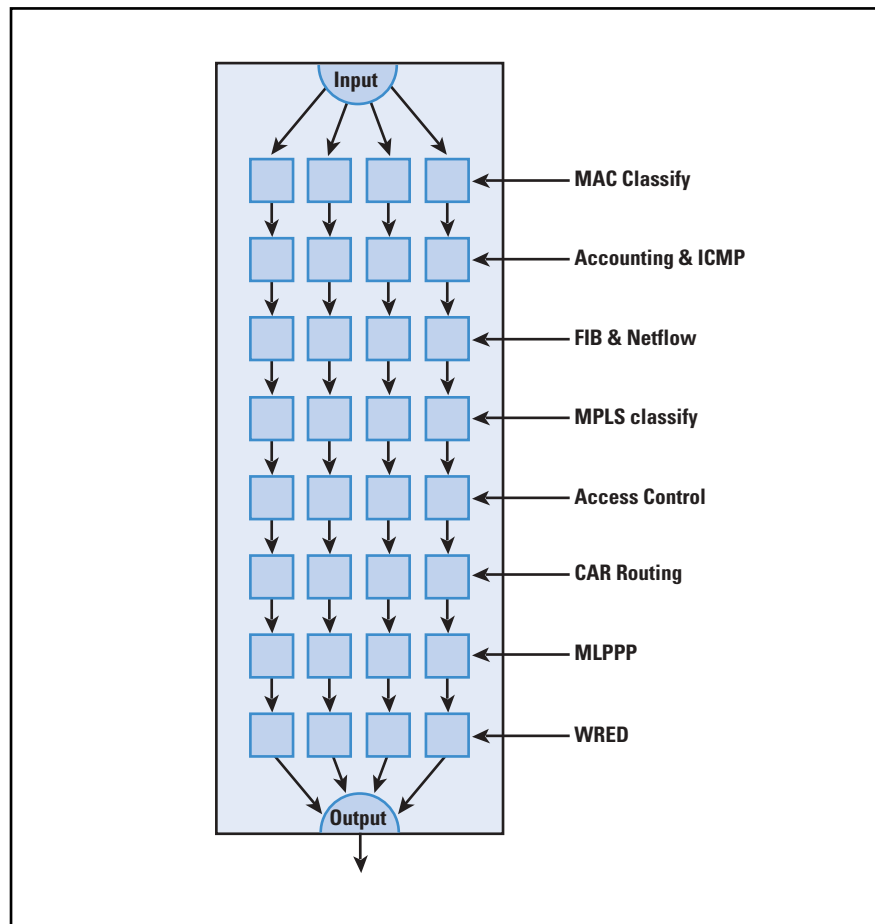
Figure 7: Structure of the Embedded Processor Complex on the Example Network Processor in Figure 6



Parallel Pipelines of Homogeneous Processors (Cisco)

Although it is not a chip vendor, Cisco Systems uses network processors in its products, and has developed network processors for internal use. One of the more interesting designs employs parallel pipelines of homogeneous processors. Figure 8 illustrates the architecture of the Cisco chip. When a packet enters, the hardware selects one of the pipelines, and the packet travels through the entire pipeline.

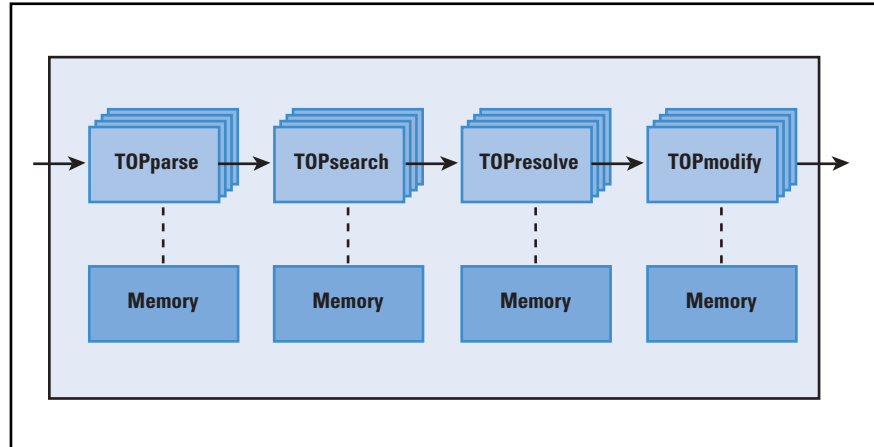
Figure 8: An Example Architecture that Uses Parallel Pipelines of Homogeneous Processors



Pipeline of Parallel Heterogeneous Processors (EZchip)

EZchip Corporation sells a network processor that combines pipelining and parallelism by using a four-stage pipeline in which each stage is implemented by parallel processors. However, instead of using the same processor type at each stage, the EZchip architecture employs heterogeneous processors, with the processor type at each stage optimized for a certain task (for example, the processor that runs forwarding code is optimized for table lookup). Figure 9 illustrates the architecture.

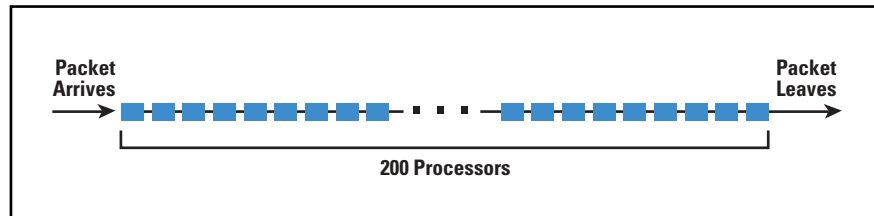
Figure 9: An Example Architecture that Uses a Pipeline of Parallel Stages with Heterogeneous Processors



Extremely Long Pipeline (Xelerated)

Xelerated Corporation sells an interesting network processor that uses a pipelining approach. Unlike other network processors, the Xelerated chip uses an extremely long pipeline of 200 stages. Figure 10 illustrates the overall architecture. To achieve high speed, each stage is limited to executing four instructions per packet.

Figure 10: An Example of an Extremely Long Pipeline with 200 Stages



In fact, the Xelerated architecture is more complex than the figure shows because the pipeline contains special hardware units after every 10 stages that allow external communication (for example, access to external memory or a call to a coprocessor).

More Details and Example Network-Processor Source Code

The previous survey is not meant to be complete. Two notable network processors have been omitted. Agere Systems and Intel each manufacture a network processor. Agere's design consists of a short pipeline that has two basic stages. Agere's architecture is both interesting and unusual because the two stages are composed of unconventional processors. For example, the processor used for classification performs high-speed pattern matching, but does not have conventional instructions for iteration or conditional testing. For details about the Agere network processor see^[1], which includes the source code for an example *Differentiated Services* (DiffServ) network system.

Intel's chip uses a parallel approach in which a set of *microengines* are programmed to handle packets. The Intel hardware allows a programmer to pass packets between microengines, meaning a programmer can decide to arrange microengines in a software pipeline. For details about the Intel network processor see^[2], which includes the source code for an example NAT implementation.

Programming Network Processors

Although the general idea of building programmable devices seems appealing, most network-processor designs make programming difficult. In particular, to achieve high speed, many designs use low-level hardware constructs and require a programmer to accommodate the hardware by writing low-level code. Many network processors are much closer to a microcontroller than a conventional processor, and are programmed in *microassembly* language. Programmers must be conscious of details such as register banks.

Programming is especially difficult in cases where the network-processor hardware uses explicit parallelism and requires a programmer to plan program execution in such a way that processors do not contend for resources simultaneously or otherwise stall. For example, on one vendor's chip, a packet processor can execute several hundred instructions while waiting for a single memory access to complete. Thus, to achieve high performance, a programmer must start a memory operation, go on with other calculations while the memory operation proceeds, and then check that the operation has completed.

In addition to considering processing, some network processors provide a set of memory technologies, and require a programmer to allocate each data item to a specific memory. A programmer must understand memory latency, the expected lifetime of a data object, and the expected frequency of access as well as properties of the hardware such as memory banks and interleaving.

A few pleasant exceptions exist. For example, Agere Systems provides special-purpose, high-level programming languages to program its network processors. Thus, it is easy to write classification code or traffic-management scripts for an Agere processor. More important, an Agere chip offers implicit parallelism: a programmer writes code as if a single processor is executing the program; the hardware automatically runs multiple copies on parallel hardware units and handles all details of coordination and synchronization.

Another pleasant exception comes from IP Fabrics, which has focused on building tools to simplify programming. Like Agere, IP Fabrics has developed a high-level language that allows a programmer to specify packet classification and the subsequent actions to be taken. The language from IP Fabrics is even more compact than the language from Agere.

Summary

To provide maximal flexibility, ease of change, and rapid development for network systems, chip vendors have defined a new technology known as network processors. The goal is to create chips for packet processing that combine the flexibility of programmable processors with the high speed of ASICs.

Because there is no consensus on which packet-processing functions are needed or which hardware architecture(s) are best, vendors have created many architectural experiments. The basic approaches comprise an embedded processor, parallelism, and hardware pipelining. Commercial chips often combine more than one approach (for example, a pipeline of parallel stages or parallel pipelines).

Programming network processors can be difficult because many network processors provide low-level hardware that requires a programmer to use a microassembly language and handle processor, memory, and parallelism details. A few exceptions exist where a vendor provides a high-level language.

References

- [1] Comer, D., *Network Systems Design Using Network Processors, Agere Version*, Prentice Hall, 2005.
- [2] Comer, D., *Network Systems Design Using Network Processors, Intel 2xxx Version*, Prentice Hall, 2005.

This article is based on material in *Network Systems Design Using Network Processors, Agere Version*, and *Network Systems Design Using Network Processors, Intel 2xxx Version* by Doug Comer. Both books are published by Prentice Hall in 2005. Used with permission.

DOUGLAS E. COMER is a Visiting Faculty at Cisco Systems, a Distinguished Professor of Computer Science at Purdue University, a Fellow of the ACM, and editor-in-chief of the journal *Software—Practice and Experience*. As a member of the IAB, he participated in the formation of the Internet, and is considered a leading authority on TCP/IP and Internetworking. He is the author of 16 technical books that have been translated into 14 languages, and are used around the world in industry and academia. Comer has been working with network processors for several years, and has reference platforms from three leading vendors in his lab at Purdue. E-mail: comer@cs.purdue.edu

Distributed Denial of Service Attacks

By Charalampos Patrikakis, Michalis Masikos, and Olga Zouraraki
National Technical University of Athens

The Internet consists of hundreds of millions of computers distributed all around the world. Millions of people use the Internet daily, taking full advantage of the available services at both personal and professional levels. The interconnectivity among computers on which the World Wide Web relies, however, renders its nodes an easy target for malicious users who attempt to exhaust their resources and launch *Denial-of-Service* (DoS) attacks against them.

A DoS attack is a malicious attempt by a single person or a group of people to cause the victim, site, or node to deny service to its customers. When this attempt derives from a single host of the network, it constitutes a DoS attack. On the other hand, it is also possible that a lot of malicious hosts coordinate to flood the victim with an abundance of attack packets, so that the attack takes place simultaneously from multiple points. This type of attack is called a *Distributed DoS*, or DDoS attack.

DDoS Attack Description

DoS attacks attempt to exhaust the victim's resources. These resources can be network bandwidth, computing power, or operating system data structures. To launch a DDoS attack, malicious users first build a network of computers that they will use to produce the volume of traffic needed to deny services to computer users. To create this attack network, attackers discover vulnerable sites or hosts on the network. Vulnerable hosts are usually those that are either running no antivirus software or out-of-date antivirus software, or those that have not been properly patched. Vulnerable hosts are then exploited by attackers who use their vulnerability to gain access to these hosts. The next step for the intruder is to install new programs (known as *attack tools*) on the compromised hosts of the attack network. The hosts that are running these attack tools are known as *zombies*, and they can carry out any attack under the control of the attacker. Many zombies together form what we call an *army*.

But how can attackers discover the hosts that will make up the attack network, and how can they install the attack tools on these hosts? Though this preparation stage of the attack is very crucial, discovering vulnerable hosts and installing attack tools on them has become a very easy process. There is no need for the intruder to spend time in creating the attack tools because there are already prepared programs that automatically find vulnerable systems, break into these systems, and then install the necessary programs for the attack. After that, the systems that have been infected by the malicious code look for other vulnerable computers and install on them the same malicious code. Because of that widespread scanning to identify victim systems, it is possible that large attack networks can be built very quickly.

The result of this automated process is the creation of a DDoS attack network that consists of handler (master) and agent (slave, daemon) machines. It can be inferred from this process that another DDos attack takes place while the attack network is being built, because the process itself creates a significant amount of traffic.

Recruiting the Vulnerable Machines

Attackers can use different kinds of techniques (referred to as *scanning techniques*) in order to find vulnerable machines^{[1][2][3]}. The most important follow:

- *Random scanning*: In this technique, the machine that is infected by the malicious code (such a machine can be either the attacker's machine or the machine of a member of their army, such as a zombie) probes IP addresses randomly from the IP address space and checks their vulnerability. When it finds a vulnerable machine, it breaks into it and tries to infect it, installing on it the same malicious code that is installed on itself. This technique creates significant traffic, because the random scanning causes a large number of compromised hosts to probe and check the same addresses. An advantage (to attackers) of this scanning method is that the malicious code can be spread very quickly because the scans seem to come from everywhere. However, the fast rate at which the malicious code is dispersed cannot last forever. After a small period of time, the spreading rate reduces because the number of the new IP addresses that can be discovered is smaller as time passes. This becomes obvious if we consider the analysis of David Moore and Colleen Shannon^[4] on the spread of the Code-Red (CRv2) Worm, which uses random scanning to spread itself.
- *Hit-list scanning*: Long before attackers start scanning, they collect a list of a large number of potentially vulnerable machines. In their effort to create their army, they begin scanning down the list in order to find vulnerable machines. When they find one, they install on it the malicious code and divide the list in half. Then they give one half to the newly compromised machine, keep the other half, and continue scanning the remaining list. The newly infected host begins scanning down its list, trying to find a vulnerable machine. When it finds one, it implements the same procedure as described previously, and in this way the hit-list scanning takes place simultaneously from an enduringly increasing number of compromised machines. This mechanism ensures that the malicious code is installed on all vulnerable machines contained in the hit list in a short period of time. In addition, the hit list possessed by a new compromised host is constantly reducing because of the partitioning of the list discussed previously.

As has been mentioned, the construction of the list is carried out long before the attackers start scanning. For that reason, the attackers can create the list at a very slow rate and for a long period of time. If the attackers conduct a slow scan, it is possible that this activity would not be noticed because a scanning process in a network usually occurs at extremely high frequencies, so a slow scan could occur without anyone realizing that it is a malicious scan.

It should also be mentioned that there are public servers such as the Netcraft Survey^[2] that can create such hit lists without scanning.

- *Topological scanning*: Topological scanning uses information contained on the victim machine in order to find new targets. In this technique, an already-compromised host looks for URLs in the disk of a machine that it wants to infect. Then it renders these URLs targets and checks their vulnerability. The fact that these URLs are valid Web servers means that the compromised host scans possible targets directly from the beginning of the scanning phase. For that reason, the accuracy of this technique is extremely good, and its performance seems to be similar to that of hit-list scanning. Hence, topological scanning can create a large army of attackers extremely quickly and in that way can accelerate the propagation of the malicious code.
- *Local subnet scanning*: This type of scanning acts behind a firewall in an area that is considered to be infected by the malicious scanning program. The compromised host looks for targets in its own local network, using the information that is hidden in “local” addresses. More specifically, a single copy of the scanning program is running behind a firewall and tries to break into all vulnerable machines that would otherwise be protected by the firewall. This mechanism can be used in conjunction with other scanning mechanisms: for example, a compromised host can start its scans with local subnet scanning, looking for vulnerable machines in its local network. As soon as it has probed all local machines, it can continue the probing process by switching to another scanning mechanism in order to scan off-local network machines. In that way, an army with numerous zombies can be constructed at an extremely high speed.
- *Permutation scanning*: In this type of scanning, all machines share a common pseudorandom permutation list of IP addresses. Such a permutation list can be constructed using any block cipher of 32 bits with a preselected key^[3]. If a compromised host has been infected during either the hit-list scanning or local subnet scanning, it starts scanning just after its point in the permutation list and scans through this list in order to find new targets. Otherwise, if it has been infected during permutation scanning, it starts scanning at a random point. Whenever it encounters an already-infected machine, it chooses a new random start point in the permutation list and proceeds from there. A compromised host can recognize an already-infected machine among noninfected ones, because such machines respond differently than other machines. The process of scanning stops when the compromised host encounters sequentially a predefined number of already-infected machines without finding new targets during that period of time. Then, a new permutation key is produced and a new scanning phase begins. This mechanism serves two major purposes: first, it prevents unnecessary reinfections of the same target because when a compromised host recognizes an already-compromised machine, it changes the way it scans according to the process described previously.

Second, this mechanism maintains the advantages (to attackers) of random scanning, because the scanning of new targets takes place in a random way. Hence, permutation scanning can be characterized as a coordinated scanning with an extremely good performance, because the randomization mechanism allows high scanning speeds.

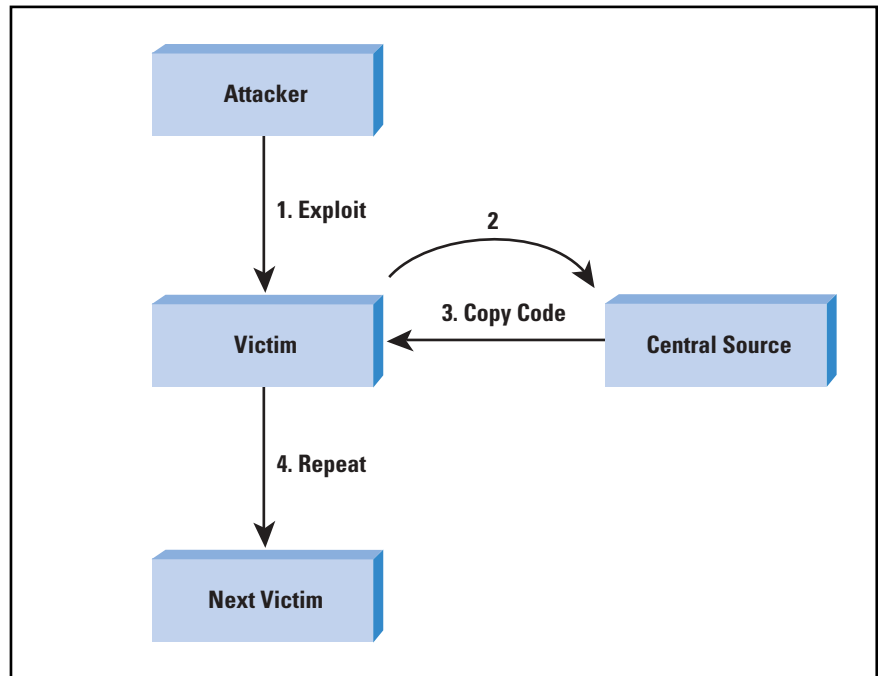
An improved version of permutation scanning is *partitioned permutation scanning*. This type of scanning is a combination of permutation and hit-list scanning. In this scenario, the compromised machine has a permutation list, which is cut in half when it finds a new target. Then it keeps one section of the list and gives the other section to the newly compromised machine. When the permutation list that an infected machine possesses reduces below a predefined level, the scanning scheme turns from partitioned permutation scanning into simple permutation scanning.

Propagating the Malicious Code

We can identify three groups of mechanisms for propagating malicious code and building attack networks^[4]:

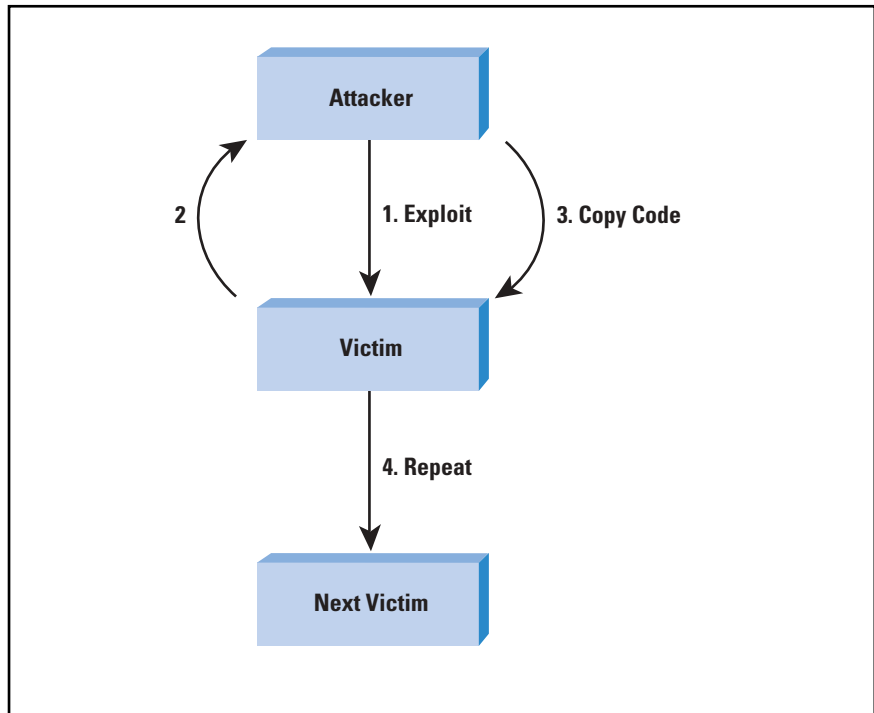
- *Central source propagation*: In this mechanism, after the discovery of the vulnerable system that will become one of the zombies, instructions are given to a central source so that a copy of the attack toolkit is transferred from a central location to the newly compromised system. After the toolkit is transferred, an automatic installation of the attack tools takes place on this system, controlled by a scripting mechanism. That initiates a new attack cycle, where the newly infected system looks for other vulnerable computers on which it can install the attack toolkit using the same process as the attacker. Like other file-transfer mechanisms, this mechanism commonly uses HTTP, FTP, and *remote-procedure call* (RPC) protocols. A graphical representation of this mechanism is shown in Figure 1.

Figure 1: Central Source Propagation



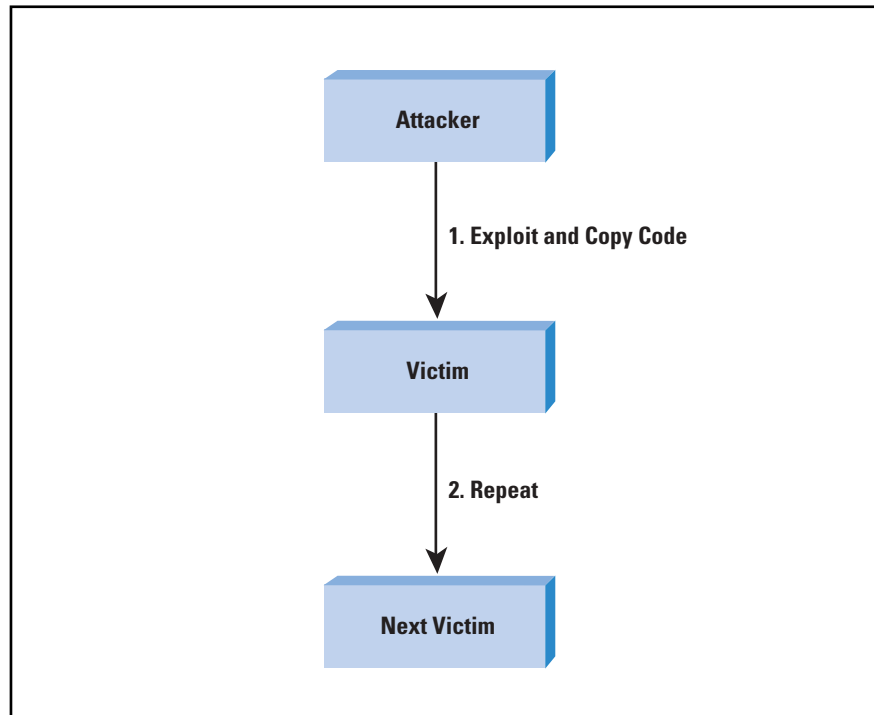
- *Back-chaining propagation*: In this mechanism, the attack toolkit is transferred to the newly compromised system from the attacker. More specifically, the attack tools that are installed on the attacker include special methods for accepting a connection from the compromised system and sending a file to it that contains the attack tools. This back-channel file copy can be supported by simple port listeners that copy file contents or by full intruder-installed Web servers, both of which use the *Trivial File Transfer Protocol (TFTP)*. Figure 2 presents the this mechanism:

Figure 2: Back-Chaining Propagation



- *Autonomous propagation*: In this mechanism, the attacking host transfers the attack toolkit to the newly compromised system at the exact moment that it breaks into that system. This mechanism differs from the previously mentioned mechanisms in that the attack tools are planted into the compromised host by the attackers themselves and not by an external file source. Figure 3 shows the autonomous propagation.

Figure 3: Autonomous Propagation



After the construction of the attack network, the intruders use handler machines to specify the attack type and the victim's address and wait for the appropriate moment in order to mount the attack. Then, either they remotely command the launch of the chosen attack to agents or the daemons “wake up” simultaneously, as they had been programmed to do. The agent machines in turn begin to send a stream of packets to the victim, thereby flooding the victim's system with useless load and exhausting its resources. In this way, the attackers render the victim machine unavailable to legitimate clients and obtain unlimited access to it, so that they can inflict arbitrary damage. The volume of traffic may be so high that the networks that connect the attacking machines to the victim may also suffer from lower performance. Hence the provision of services over these networks is no longer possible, and in this way their clients are denied those services. Thus, the network that has been burdened by the attack load can be considered as one more victim of the DDos attack.

The whole procedure for carrying out a DDoS attack is mostly automated thanks to various attack tools. According to^[5], the existence of the first controllable DDOS tool was reported by the *CERT Coordination Center* (CERT/CC) in early 1998 and it was called “Fapi.” It is a tool that does not provide easy controls for setting up the DDoS network and does not handle networks with more than 10 hosts very well. In mid-1999 Trinoo arrived. Later that year the existence of *Tribe Flood Network* (TFN) and its upgraded version TFN2K (or TFN2000) was reported. Stacheldraht (German for “barbed wire”) evolved out of the latter two tools (Trinoo and TFN). This tool is remarkable because it has full-control features and a Blowfish-encrypted control channel for the attacker. Moreover, in early 2000 it mutated into StacheldrahtV4, and later into Stacheldraht v1.666.

However, the development of attack tools did not stop, and many tools were later introduced, such as Mstream, Omega, Trinity, Derivatives, myServer, and Plague^[6]. Dave Dittrich and his partners have provided the most comprehensive analyses of the Trinoo, Tribe Flood Network, Stacheldraht, shaft, and mstream DDoS attack tools^[7]. Through this work, a lot of malicious code was captured, important observations were made about DDoS attack tools, and solutions were proposed toward detection and defense.

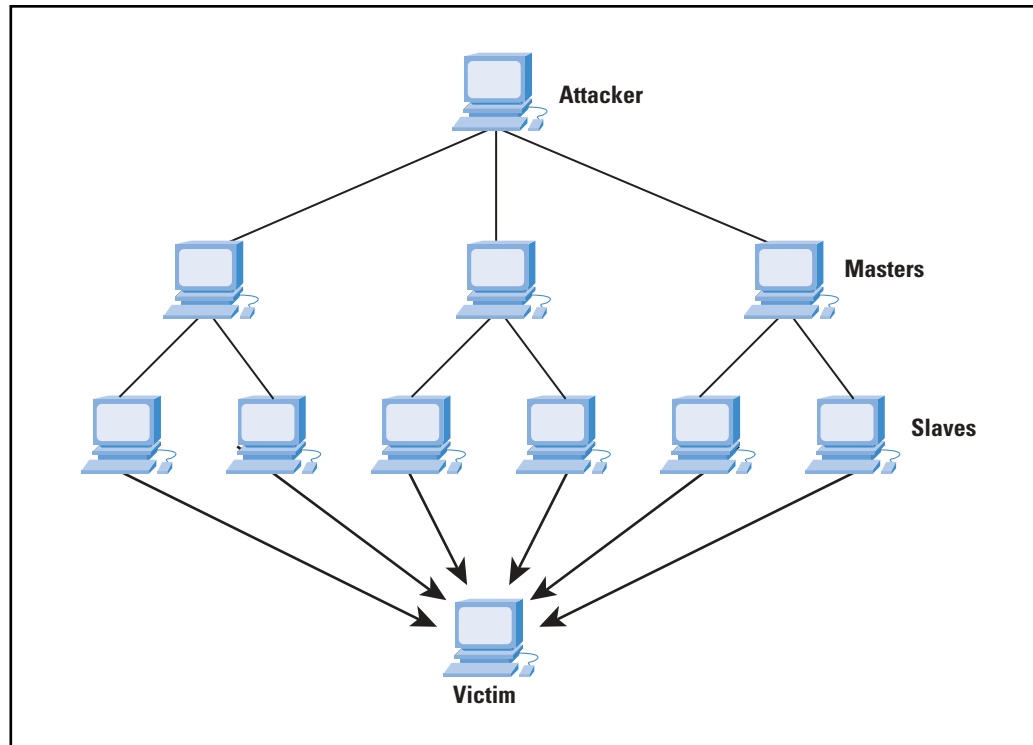
DDoS Attack Taxonomy

As has been already said, a DDoS attack takes place when many compromised machines infected by the malicious code act simultaneously and are coordinated under the control of a single attacker in order to break into the victim's system, exhaust its resources, and force it to deny service to its customers. There are mainly two kinds of DDoS attacks^[10]: typical DDoS attacks and *distributed reflector DoS* (DRDoS) attacks. The following paragraphs describe these two kinds analytically.

Typical DDoS Attacks

In a typical DDoS attack, the army of the attacker consists of *master zombies* and *slave zombies*. The hosts of both categories are compromised machines that have arisen during the scanning process and are infected by malicious code. The attacker coordinates and orders master zombies and they, in turn, coordinate and trigger slave zombies. More specifically, the attacker sends an attack command to master zombies and activates all attack processes on those machines, which are in hibernation, waiting for the appropriate command to wake up and start attacking. Then, master zombies, through those processes, send attack commands to slave zombies, ordering them to mount a DDoS attack against the victim. In that way, the agent machines (slave zombies) begin to send a large volume of packets to the victim, flooding its system with useless load and exhausting its resources. Figure 4 shows this kind of DDoS attack.

Figure 4: A DDoS Attack



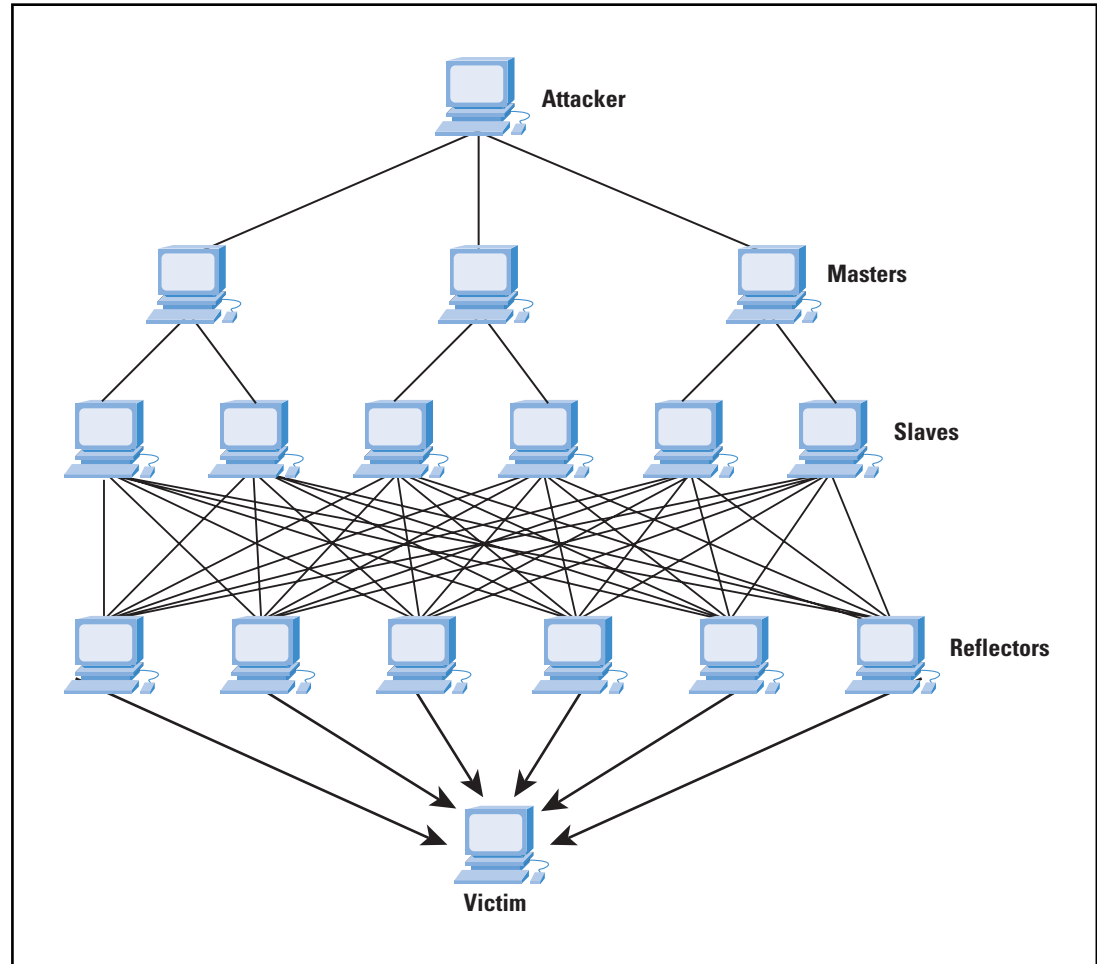
In cases of DDoS attacks, spoofed source IP addresses are used in the packets of the attack traffic. An attacker prefers to use such counterfeit source IP addresses for two major reasons: first, the attackers want to hide the identity of the zombies so that the victim cannot trace the attack back to them. The second reason concerns the performance of the attack. The attackers want to discourage any attempt of the victim to filter out the malicious traffic.

DRDoS Attacks

Unlike typical DDoS attacks, in DRDoS attacks the army of the attacker consists of master zombies, slave zombies, and reflectors^[11]. The scenario of this type of attack is the same as that of typical DDoS attacks up to a specific stage. The attackers have control over master zombies, which, in turn, have control over slave zombies. The difference in this type of attack is that slave zombies are led by master zombies to send a stream of packets with the victim's IP address as the source IP address to other uninfected machines (known as *reflectors*), exhorting these machines to connect with the victim. Then the reflectors send the victim a greater volume of traffic, as a reply to its exhortation for the opening of a new connection, because they believe that the victim was the host that asked for it. Therefore, in DRDoS attacks, the attack is mounted by noncompromised machines, which mount the attack without being aware of the action.

Comparing the two scenarios of DDoS attacks, we should note that a DRDoS attack is more detrimental than a typical DDoS attack. This is because a DRDoS attack has more machines to share the attack, and hence the attack is more distributed. A second reason is that a DRDoS attack creates a greater volume of traffic because of its more distributed nature. Figure 5 graphically depicts a DRDoS attack.

Figure 5: A DRDoS Attack



Well-Known DDoS Attacks

This article would be incomplete without reference to some of the most well-known DDoS attacks. Some of the most famous documented DDoS attacks^{[12][13]} are summarized in the following:

- *Apache2*: This attack is mounted against an Apache Web server where the client asks for a service by sending a request with many HTTP headers. However, when an Apache Web server receives many such requests, it cannot confront the load and it crashes.

- *ARP Poison: Address Resolution Protocol (ARP) Poison* attacks require the attacker to have access to the victim's LAN. The attacker deludes the hosts of a specific LAN by providing them with wrong MAC addresses for hosts with already-known IP addresses. This can be achieved by the attacker through the following process: The network is monitored for "arp who-has" requests. As soon as such a request is received, the malevolent attacker tries to respond as quickly as possible to the questioning host in order to mislead it for the requested address.
- *Back:* This attack is launched against an apache Web server, which is flooded with requests containing a large number of front-slash (/) characters in the URL description. As the server tries to process all these requests, it becomes unable to process other legitimate requests and hence it denies service to its customers.
- *CrashIIS:* The victim of a CrashIIS attack is commonly a Microsoft Windows NT IIS Web server. The attacker sends the victim a malformed GET request, which can crash the Web server.
- *DoSNuke:* In this kind of attack, the Microsoft Windows NT victim is inundated with "out-of-band" data (MSG_OOB). The packets being sent by the attacking machines are flagged "urg" because of the MSG_OOB flag. As a result, the target is weighed down, and the victim's machine could display a "blue screen of death."
- *Land:* In Land attacks, the attacker sends the victim a TCP SYN packet that contains the same IP address as the source and destination addresses. Such a packet completely locks the victim's system.
- *Mailbomb:* In a Mailbomb attack, the victim's mail queue is flooded by an abundance of messages, causing system failure.
- *SYN Flood:* A SYN flood attack occurs during the three-way handshake that marks the onset of a TCP connection. In the three-way handshake, a client requests a new connection by sending a TCP SYN packet to a server. After that, the server sends a SYN/ACK packet back to the client and places the connection request in a queue. Finally, the client acknowledges the SYN/ACK packet. If an attack occurs, however, the attacker sends an abundance of TCP SYN packets to the victim, obliging it both to open a lot of TCP connections and to respond to them. Then the attacker does not execute the third step of the three-way handshake that follows, rendering the victim unable to accept any new incoming connections, because its queue is full of half-open TCP connections.
- *Ping of Death:* In Ping of Death attacks, the attacker creates a packet that contains more than 65,536 bytes, which is the limit that the IP protocol defines. This packet can cause different kinds of damage to the machine that receives it, such as crashing and rebooting.

- *Process Table*: This attack exploits the feature of some network services to generate a new process each time a new TCP/IP connection is set up. The attacker tries to make as many uncompleted connections to the victim as possible in order to force the victim's system to generate an abundance of processes. Hence, because the number of processes that are running on the system cannot be boundlessly large, the attack renders the victim unable to serve any other request.
- *Smurf Attack*: In a "smurf" attack, the victim is flooded with *Internet Control Message Protocol* (ICMP) "echo-reply" packets. The attacker sends numerous ICMP "echo-request" packets to the broadcast address of many subnets. These packets contain the victim's address as the source IP address. Every machine that belongs to any of these subnets responds by sending ICMP "echo-reply" packets to the victim. Smurf attacks are very dangerous, because they are strongly distributed attacks.
- *SSH Process Table*: Like the Process Table attack, this attack makes hundreds of connections to the victim with the *Secure Shell* (SSH) Protocol without completing the login process. In this way, the daemon contacted by the SSH on the victim's system is obliged to start so many SSH processes that it is exhausted.
- *Syslogd*: The Syslogd attack crashes the *syslogd* program on a Solaris 2.5 server by sending it a message with an invalid source IP address.
- *TCP Reset*: In TCP Reset attacks, the network is monitored for "tcp-connection" requests to the victim. As soon as such a request is found, the malevolent attacker sends a spoofed TCP RESET packet to the victim and obliges it to terminate the TCP connection.
- *Teardrop*: While a packet is traveling from the source machine to the destination machine, it may be broken up into smaller fragments, through the process of fragmentation. A Teardrop attack creates a stream of IP fragments with their offset field overloaded. The destination host that tries to reassemble these malformed fragments eventually crashes or reboots.
- *UDP Storm*: In a *User Datagram Protocol* (UDP) connection, a character generation ("chargen") service generates a series of characters each time it receives a UDP packet, while an echo service echoes any character it receives. Exploiting these two services, the attacker sends a packet with the source spoofed to be that of the victim to another machine. Then, the echo service of the former machine echoes the data of that packet back to the victim's machine and the victim's machine, in turn, responds in the same way. Hence, a constant stream of useless load is created that burdens the network.

The first DoS attack occurred against Panix, the New York City area's oldest and largest *Internet Service Provider* (ISP), on September 6, 1996, at about 5:30 p.m.^[14]. The attack was against different computers on the provider's network, including mail, news, and Web servers, user "login" machines, and name servers. The Panix attack was a SYN Flood attack deriving from random IP addresses and directed toward server *Simple Mail Transfer Protocol* (SMTP) ports. More specifically, Panix's computers were flooded by, on average, 150 SYN packets per second (50 per host), so Panix could not respond to legitimate requests^[15]. Because the attackers used spoofed source IP addresses in their packets, the addresses could not be traced and malicious traffic could not be filtered. For that reason the attack was not immediately confronted. The solution was to use a special structure, instead of full *Transmission Control Block* (TCB), to hold half-open connections until the last ACK packet was received. In that way, the listen queue was large enough to keep all the SYN requests before the half-open connection timed out. The timeout, on the other hand, was adjusted to 94 seconds^[16]. However, although Panix overcame this attack, the new threat (DoS attacks) made administrators worry.

Problems Caused and Countermeasures

The results of these attacks are disastrous. DDoS attacks have two characteristics: they are both distributed attacks and denial-of-service attacks. Distributed means that they are large-scale attacks having a great impact on the victims. Denial of service means that their goal is to deny the victim's access to a particular resource (service). This is not too difficult because the Internet was not designed with security in mind.

First, available *bandwidth* is one of the "goods" that attackers try to consume. Flooding the network with useless packets, for example, prevents legitimate ICMP echo packets from traveling over the network. Secondly, attackers try to consume *CPU power*. By generating several thousands of useless processes on the victim's system, attackers manage to fully occupy memory and process tables. In this way the victim's computer cannot execute any process and the system breaks down. Using this method, the attacker manages to prevent clients from accessing the victim's services and disrupts the current connections. Finally, attackers try to occupy victims' *services* so that no one else can access them. For example, by leaving TCP connections half open, attackers manage to consume the victim's data structures, and when they do so, no one else can establish a TCP connection with that victim.

The impact of these attacks is catastrophic, especially when victims are not individuals but companies. DDoS attacks prevent victims either from using the Internet, or from being reached by other people. Consequently, when the victim is an ISP, the results of such an attack are far more severe. ISPs' clients will not be served. E-business is also top on the "hit list." Being off line for a few hours could result in the loss of large sums of money for an ISP. Finally, the fact that companies use the Internet more and more for advertising or for providing goods and services increases the severity of such incidents.

Defense Mechanisms

From the beginning, all legitimate users have tried to respond against these threats. University communities and software corporations have proposed several methods against the DDoS threat. Despite the efforts, the solution remains a dream. The attackers manage to discover other weaknesses of the protocols and—what is worse—they exploit the defense mechanisms in order to develop attacks. They discover methods to overcome these mechanisms or they exploit them to generate false alarms and to cause catastrophic consequences.

Many experts have tried to classify the DDoS defense mechanisms in order to clarify them. This classification gives users an overall view of the situation and helps defense-mechanism developers cooperate against the threat. The basic discrimination is between *preventive* and *reactive* defense mechanisms.

Preventive Mechanisms

The preventive mechanisms try to eliminate the possibility of DDoS attacks altogether or to enable potential victims to endure the attack without denying services to legitimate clients. With regard to attack prevention, countermeasures can be taken on victims or on zombies. This means modification of the system configuration to eliminate the possibility of accepting a DDoS attack or participating unwillingly in a DDoS attack. Hosts should guard against illegitimate traffic from or toward the machine. By keeping protocols and software up-to-date, we can reduce the weaknesses of a computer. A regular scanning of the machine is also necessary in order to detect any “anomalous” behavior. Examples of system security mechanisms include monitoring access to the computer and applications, and installing security patches, firewall systems, virus scanners, and intrusion detection systems automatically. The modern trend is toward security companies that guard a client’s network and inform the client in case of attack detection to take defending measures. Several sensors monitor the network traffic and send information to a server in order to determine the “health” of the network. Securing the computer reduces the possibility of being not only a victim, but also a zombie. Not being a zombie is very important because it wipes out the attacker’s army. All these measures can never be 100-percent effective, but they certainly decrease the frequency and strength of DDoS attacks.

Many other measures can be taken in order to reduce the attacker’s army or restrict its “power.” Studying the attack methods can lead to recognizing loopholes in protocols. For example, administrators could adjust their network gateways in order to filter input and output traffic. The source IP address of output traffic should belong to the subnet-work, whereas the source IP address of input traffic should not. In this way, we can reduce traffic with spoofed IP addresses on the network^[28].

Furthermore, over the last few years, several techniques have been proposed to test systems for possible drawbacks, before their shipment to the market. More precisely, by replacing the components of a system with malicious ones we can discover whether the system can survive an attack situation^[38]. If the system breaks down, a drawback has been detected and developers must correct it.

On the other hand, DoS prevention mechanisms enable the victim to endure attack attempts without denying service to legitimate clients. Until now, two methods have been proposed for this scenario. The first one refers to policies that increase the privileges of users according to their behavior. When users' identities are verified, then no threat exists. Any illegitimate action from those users can lead to their legal prosecution. The second method is usually too expensive; it involves increasing the effective resources to such a degree that DDoS effects are limited. Most of the time application of such a measure is impossible.

Reactive Mechanisms

The reactive mechanisms (also referred to as *Early Warning Systems*) try to detect the attack and respond to it immediately. Hence, they restrict the impact of the attack on the victim. Again, there is the danger of characterizing a legitimate connection as an attack. For that reason it is necessary for researchers to be very careful.

The main detection strategies are *signature detection*, *anomaly detection*, and *hybrid systems*. Signature-based methods search for patterns (signatures) in observed network traffic that match known attack signatures from a database. The advantage of these methods is that they can easily and reliably detect known attacks, but they cannot recognize new attacks. Moreover, the signature database must always be kept up-to-date in order to retain the reliability of the system.

Anomaly-based methods compare the parameters of the observed network traffic with normal traffic. Hence it is possible for new attacks to be detected. However, in order to prevent a false alarm, the model of "normal traffic" must always be kept updated and the threshold of categorizing an anomaly must be properly adjusted.

Finally, hybrid systems combine both these methods. These systems update their signature database with attacks detected by anomaly detection. Again the danger is great because an attacker can fool the system by characterizing normal traffic as an attack. In that case an *Intrusion Detection System* (IDS) becomes an attack tool. Thus IDS designers must be very careful because their research can boomerang.

After detecting the attack, the reactive mechanisms respond to it. The relief of the impact of the attack is the primary concern. Some mechanisms react by limiting the accepted traffic rate. This means that legitimate traffic is also blocked. In this case the solution comes from traceback techniques that try to identify the attacker. If attackers are identified, despite their efforts to spoof their address, then it is easy to filter their traffic. Filtering is efficient only if attackers' detection is correct. In any other case filtering can become an attacker's tool.

The University of Washington provides an example of attack detection. Dave Dittrich and his team of 40 people discovered that more than 30 of their systems were zombies exploited by a single attacker^[39]. By monitoring network traffic, Dittrich's team located directory and file names uncommon to the Windows operating systems the attacker ran on the network, as well as the port through which all these files were running communications.

Difficulties in Defending

Development of detection and defending tools is very complicated. Designers must think in advance of every possible situation because every weakness can be exploited. Difficulties involve:

- DDoS attacks flood victims with packets. This means that victims cannot contact anyone else in order to ask for help. So it is possible for a network neighbor to be attacked, but nobody would know it and nobody can help. Consequently, any action to react can be taken only if the attack is detected early. But can an attack be detected early? Usually traffic flow increases suddenly and without any warning^{[34][35][36]}. For this reason defense mechanisms must react quickly.
- Any attempt of filtering the incoming flow means that legitimate traffic will also be rejected. And if legitimate traffic is rejected, how will applications that wait for information react? On the other hand, if zombies number in the thousands or millions, their traffic will flood the network and consume all the bandwidth. In that case filtering is useless because nothing can travel over the network.
- Attack packets usually have spoofed IP addresses. Hence it is more difficult to trace back to their source. Furthermore, it is possible that intermediate routers and ISPs may not cooperate in this attempt. Sometimes attackers, by spoofing source IP addresses, create counterfeit armies. Packets might derive from thousands of IP addresses, but zombies number only a few tens, for example.
- Defense mechanisms are applied in systems with differences in software and architecture. Also systems are managed by users with different levels of knowledge. Developers must design a platform independent of all these parameters.^[37]

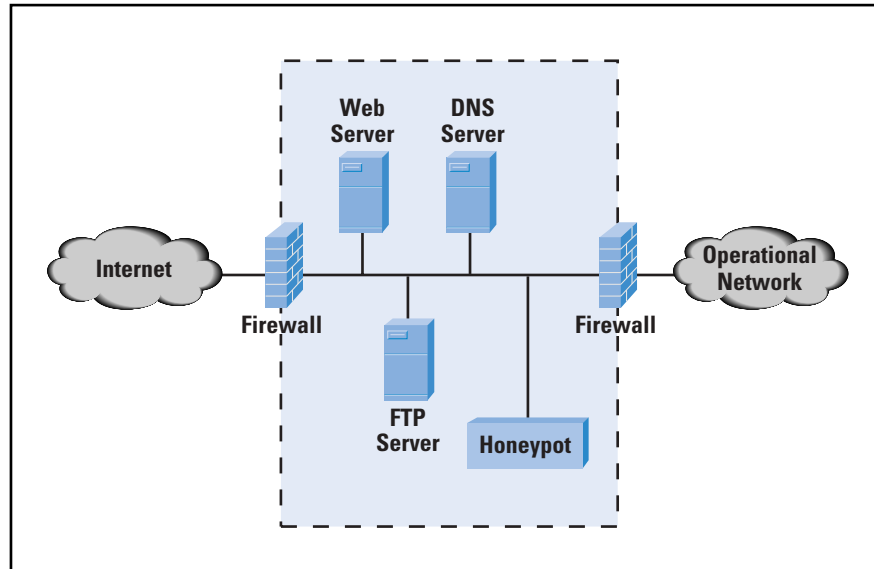
Modern Tendencies in Defending Against DDoS Attacks

Until now, developers have not managed to develop a 100-percent-effective defense mechanism. All mechanisms that have been presented either can confront only specific DDoS attacks or are being finally compromised by the attackers. Therefore, developers are currently working on DDoS diversion systems. *Honeypots* are the best representative of this category (See Figure 6).

Honeypots

There are two basic types of honeypots: *low-interaction honeypots* and *high-interaction honeypots*. The first ones refer to emulating services and operating systems. It is easy and safe to implement them. Attackers are not allowed to interact with the basic operating system, but only with specific services. For that reason, this type of honeypot cannot provide detailed informations for attackers' actions and they can easily be detected. However, they can detect communication attempts toward unused IP addresses. In that case an alarm is triggered, warning that someone is trying to compromise or attack the network. But what happens if the attack is not directed against the emulated service?

Figure 6: Honeypot



The answer comes from high-interaction honeypots. In [41], HoneyNet is proposed. HoneyNet is not a software solution that can be installed on a computer but a whole architecture, a network that is created to be attacked. Within this network, every activity is recorded and attackers are being trapped. Encrypted SSH sessions, e-mails, file uploads, and every possible attacker's action is captured. Moreover, a Honeywall gateway allows incoming traffic, but controls outgoing traffic using intrusion prevention technologies. This allows the attacker to interact with HoneyNet systems, but prevents the attacker from harming other non-HoneyNet systems. By studying the captured traffic, researchers can discover new methods and tools and they can fully understand attackers' tactics. However, HoneyNet systems are more complex to install and deploy and the risk is increased as attackers interact with real operating systems and not with emulations. But what would happen if someone did compromise such a system? The consequences could be disastrous.

Route Filter Techniques

Different suggestions for defending against DDoS attacks derive from the *Border Gateway Protocol* (BGP) community. When routing protocols were designed, developers did not focus on security, but effective routing mechanisms and routing loop avoidance. Early on, attackers started directing their attention towards routers. By gaining access to a router, they could direct the traffic over bottlenecks, view critical data, and modify them. Cryptographic authentication mitigates these threats. Because of neighbor authentication, the routing update comes from a trusted source and there is no possibility that someone can give routers invalid routing information in order to compromise a network. On the other hand, routing filters are necessary for preventing critical routes and subnetworks from being advertised and suspicious routes from being incorporated in routing tables. In that way, attackers do not know the route toward critical servers and suspicious routes are not used.

Two other route filter techniques, *blackhole routing* and *sinkhole routing*, can be used when the network is under attack. These techniques try to temporarily mitigate the impact of the attack. The first one directs routing traffic to a null interface, where it is finally dropped. At first glance, it would be perfect to “blackhole” malicious traffic. But is it always possible to isolate malicious from legitimate traffic? If victims know the exact IP address being attacked, then they can ignore traffic originating from these sources. This way, the attack impact is restricted because the victims do not consume CPU time or memory as a consequence of the attack. Only network bandwidth is consumed. However, if the attackers’ IP addresses cannot be distinguished and all traffic is blackholed, then legitimate traffic is dropped as well. In that case, this filter technique fails.

Sinkhole routing involves routing suspicious traffic to a valid IP address where it can be analyzed. There, traffic that is found to be malicious is rejected (routed to a null interface); otherwise it is routed to the next hop. A sniffer on the sinkhole router can capture traffic and analyze it. This technique is not as severe as the previous one. The effectiveness of each mechanism depends on the strength of the attack. Specifically, sinkholing cannot react to a severe attack as effectively as blackholing. However, it is a more sophisticated technique, because it is more selective in rejecting traffic.

Filtering malicious traffic seems to be an effective countermeasure against DDoS. The closer to the attacker the filtering is applied, the more effective it is. This is natural, because when traffic is filtered by victims, they “survive,” but the ISP’s network is already flooded. Consequently, the best solution would be to filter traffic on the source; in other words, filter zombies’ traffic.

Until now, three filtering possibilities have been reported concerning criteria for filters. The first one is filtering on the *source address*. This one would be the best filtering method, if we knew each time who the attacker is. However, this is not always possible because attackers usually use spoofed IP addresses. Moreover, DDoS attacks usually derive from thousands of zombies and this makes it too difficult to discover all the IP addresses that carry out the attack. And even if all these IP addresses are discovered, the implementation of a filter that rejects thousands of IP addresses is practically impossible to deploy.

The second filtering possibility is filtering on the *service*. This tactic presupposes that we know the attack mechanism. In this case, we can filter traffic toward a specific UDP port or a TCP connection or ICMP messages. But what if the attack is directed toward a very common port or service? Then we must either reject every packet (even if it is legitimate) or suffer the attack.

Finally, there is the possibility of filtering on the *destination address*. DDoS attacks are usually addressed to a restricted number of victims, so it seems to be easy to reject all traffic toward them. But this means that legitimate traffic is also rejected. In case of a large-scale attack, this should not be a problem because the victims will soon break down and the ISP will not be able to serve anyone. So filtering prevents victims from breaking down by simply keeping them isolated.

Fred Baker and Paul Ferguson developed a technique called *Ingress Filtering* for mitigating DoS attacks (and, later, DDoS attacks too). After the Panix attack and a few other attacks, Paul Ferguson wrote RFC 2267^[42], which became *Best Current Practices* (BCP) 38 in RFC 2827^[43]. This RFC presents a method for using ingress traffic filtering against DoS attacks that use forged IP addresses and try to be propagated from “behind” an ISP’s aggregation point. This method prevents the attack from forged source addresses, but nothing can be done against an attack from a valid source address. However, in that case, if the attack is detected, it is easy to trace the attacker. Finally, although this solution allows the network to protect itself from other attacks too (for example, spoofed management access to networking equipment), it can also create some problems, for example, with multihoming.

For that reason, RFC 2827 was recently (March 2004) updated by Fred Baker in BCP 84/ RFC 3704^[44]. This RFC describes and evaluates the current ingress filtering mechanisms, examines some implementation matters related to ingress filtering, and presents some solutions to ingress filtering with multihoming. According to this RFC, ingress filtering should be implemented at multiple levels in order to prohibit the use of spoofed addresses and to make attackers more traceable, even if asymmetric/multihomed networks are presented. However, although Ferguson’s work was published a long time ago, service providers in some cases ignore his suggestions.

Hybrid Methods and Guidelines

Currently researchers try to combine the advantages from all the methods stated previously in order to minimize their disadvantages. As a result, several mechanisms that implement two or more of these techniques are proposed for mitigation of the impact of DDoS attacks. The best solution to the DDoS problem seems to be the following: victims must detect that they are under attack as early as possible. Then they must trace back the IP addresses that caused the attack and warn zombies administrators about their actions. In that way, the attack can be confronted effectively.

However, as we saw previously, this is currently impossible. The lack of a 100-percent-effective defending tool imposes the necessity of private alerts. Users must care for their own security. Some basic suggestions follow:

- Prevent installation of distributed attack tools on our systems. This will help to restrict the zombies army. Several tasks also need to be performed. First, keep protocols and operating systems up-to-date. We can prevent system exploitation by eliminating the number of weaknesses of our system.
- Use firewalls in gateways to filter incoming and outgoing traffic. Incoming packets with source IP addresses belonging to the subnetwork and outgoing packets with source IP addresses not belonging to the subnetwork are not logical.
- Deploy IDS systems to detect patterns of attacks.
- Deploy antivirus programs to scan malicious code in our system.

Further Thoughts

The Internet is not stable—it reforms itself rapidly. This means that DDoS countermeasures quickly become obsolete. New services are offered through the Internet, and new attacks are deployed to prevent clients from accessing these services. However, the basic issue is whether DDoS attacks represent a network problem or an individual problem—or both. If attacks are mainly a network problem, a solution could derive from alterations in Internet protocols. Specifically, routers could filter malicious traffic, attackers could not spoof IP addresses, and there would be no drawback in routing protocols. If attacks are mostly the result of individual system weaknesses, the solution could derive from an effective IDS system, from an antivirus, or from an invulnerable firewall. Attackers then could not compromise systems in order to create a “zombies” army. Obviously, it appears that both network and individual hosts constitute the problem. Consequently, countermeasures should be taken from both sides. Because attackers cooperate in order to build the perfect attack methods, legitimate users and security developers should also cooperate against the threat. The solution will arise from combining both network and individual countermeasures.

References

- [1] Kevin Tsui, "Tutorial-Virus (Malicious Agents)," University of Calgary, October 2001.
- [2] Nicholas Weaver, "Warhol Worms: The Potential for Very Fast Internet Plagues,"
<http://www.iwar.org.uk/comsec/resources/worms/warhol-worm.htm>
- [3] Nicholas Weaver, U.C. Berkeley BRASS group, "Potential Strategies for High Speed Active Worms: A Worst Case Analysis," February 2002
- [4] David Moore and Colleen Shannon, "The Spread of the Code Red Worm (crv2)," July 2001,
http://www.caida.org/analysis/security/codered/coderedv2_analysis.xml#animations
- [5] "A Chronology of CERT Coordination Center Involvement with Distributed Denial-of-Service Tools,"
<http://www.cdt.org/security/dos/000229senatehouse/chron.html>
- [6] "Analyzing Distributed Denial Of Service Tools: The Shaft Case," Sven Dietrich, NASA Goddard Space Flight Center; Neil Long, Oxford University; David Dittrich, University of Washington,
http://www.usenix.org/events/lisa2000/full_papers/dietrich/dietrich_html/
- [7] <http://staff.washington.edu/dittrich>
- [8] Kevin J. Houle, CERT/CC; George M. Weaver, CERT/CC, in collaboration with: Neil Long, Rob Thomas, "Trends in Denial of Service Attack Technology," V1.0, October 2001.
- [9] <http://staff.washington.edu/dittrich/misc/stacheldraht.analysis>
- [10] T. Peng, C. Leckie, and K. Ramamohanarao, "Detecting Distributed Denial of Service Attacks Using Source IP Address Monitoring," The University of Melbourne, Australia, 2003.
- [11] Steve Gibson, "Distributed Reflection Denial of Service Description and Analysis of a Potent, Increasingly Prevalent, and Worrisome Internet Attack," February 2002.
- [12] <http://www.ll.mit.edu/IST/ideval/docs/1999/attackDB.html>
- [13] Yanet Manzano, "Tracing the Development of Denial of Service Attacks: A Corporate Analogy," 2003,
<http://www.acm.org/crossroads/xrds10-1/tracingDOS.html>
- [14] <http://www.panix.com/press/synattack.html>
- [15] <http://cypherpunks.venona.com/date/1996/09/msg01055.html>
- [16] <http://cypherpunks.venona.com/date/1996/09/msg01061.html>

- [17] Larry Rogers, "What Is a Distributed Denial of Service (DDoS) Attack and What Can I Do About It?" February 2004,
<http://www.cert.org/homeusers/ddos.html>
- [18] Alefiya Hussain, John Heidemann, and Christos Papadopoulos, "A Framework for Classifying Denial of Service Attacks," 25 February 2003.
- [19] <http://www.cs.berkeley.edu/~nweaver/warhol.old.html>
- [20] CIS 659 "Introduction to Network Security – Fall 2003,"
<http://www.cis.udel.edu/~sunshine/F03/CIS659/class15.pdf>
- [21] Miguel Vargas Martin, School of Computer Science, Carleton University, "Overview of Worms and Defence Strategies," October 2003.
- [22] "Computer Security," Testimony of Richard D. Pethia, Director, CERT Centers Software Engineering Institute, Carnegie Mellon University, March 2000,
http://www.cert.org/congressional_testimony/Pethia_testimony_Mar9.html#Distributed
- [23] Jelena Mirkovic, Janice Martin, and Peter Reiher, UCLA, "A Taxonomy of DDoS Attacks and DDoS Defense Mechanisms."
- [24] Distributed Denial of Service Tools,
http://www.cert.org/incident_notes/IN-99-07.html
- [25] Barbara Fraser, Lawrence Rogers, and Linda Pesante, "Was the Melissa Virus So Different?" *The Internet Protocol Journal*, Volume 2, No. 2, June 1999.
- [26] <http://news.bbc.co.uk/1/hi/sci/tech/635444.stm>
- [27] <http://www.nta-monitor.com/newrisks/feb2000/yahoo.htm>
- [28] <http://www.cert.org/advisories/CA-1996-21.html>
- [29] S. Axelsson, "Intrusion Detection Systems: A Survey and Taxonomy," Technical Report 99-15, Department of Computer Engineering, Chalmers University, March 2000.
- [30] J. Shapiro and N. Hardy, "EROS: A principle-driven Operating System from the Ground Up," *IEEE Software*, pp. 26–33, January/February 2002.
- [31] A. Garg and A. L. Narasimha Reddy, "Mitigating Denial of Service Attacks Using QoS Regulation," Texas A & M University Tech report, TAMU-ECE-2001-06.
- [32] Y. L. Zheng and J. Leiwo, "A method to implement a Denial of Service Protection Base," *Information Security and Privacy*, Volume 1270 of *Lecture Notes in Computer Science* (LNCS), pp. 90–101, 1997.

- [33] CERT on Home Network Security:
http://www.cert.org/tech_tips/home_networks.html
- [34] CERT on SMURF Attacks:
<http://www.cert.org/advisories/CA-1998-01.html>
- [35] CERT on TCP SYN Flooding Attacks:
<http://www.cert.org/advisories/CA-1996-21.html>
- [36] CERT TRIN00 Report:
http://www.cert.org/incident_notes/IN-99-07.html#trinoo
- [37] <http://falcon.jmu.edu/~flynnngn/whatnext.htm>
- [38] Charalampos Patrikakis, Thomas Kalamaris, Vaios Kakavas, "Performing Integrated System Tests Using Malicious Component Insertion," *Electronic Notes in Theoretical Computer Science*, Volume 82 No. 6 (2003).
- [39] <http://www.paypal.com/html/computerworld-011402.html>
- [40] Ho Chung, "An Evaluation on Defensive Measures against Denial-of-Service Attacks," Fall 2002.
- [41] Nathalie Weiler, "Honeypots for Distributed Denial of Service Attacks,"
www.tik.ee.ethz.ch/~weiler/papers/wetice02.pdf
- [42] P. Ferguson and D. Senie, "Network Ingress Filtering: Defeating Denial of Service Attacks which employ IP Source Address Spoofing," RFC 2267, January 1998.
- [43] P. Ferguson and D. Senie, "Network Ingress Filtering: Defeating Denial of Service Attacks which employ IP Source Address Spoofing," RFC 2827, May 2000.
- [44] F. Baker and P. Savola, "Ingress Filtering for Multihomed Networks," RFC 3704, March 2004.
- [45] Taxonomies of Distributed Denial of Service Networks, Attacks, Tools, and Countermeasures:
www.ee.princeton.edu/~rblee/DDoS%20Survey%20Paper_v7final.doc
- [46] Lance Spitzner, "Honeypots Definitions and Value of Honeypots," May 2003, <http://www.tracking-hackers.com>
- [47] How to Get Rid of Denial of Service Attacks:
<http://www.bgpexpert.com/antidos.php>
- [48] Proposed Solutions to DDoS Information, March 2001:
http://www.cs.virginia.edu/~survive/ddos/ddos_solutions.html
- [49] Dennis Fisher, "Thwarting the Zombies," March 2003:
<http://www.eweek.com/article2/0,3959,985389,00.asp>

- [50] Merike Kaeo, "Route to Security," March 2004,
http://infosecuritymag.techtarget.com/ss/0,295796,sid6_iss346_art668,00.html
- [51] "Report to the President's Commission on Critical Infrastructure Protection," James Ellis, David Fisher, Thomas Longstaff, Linda Pesante, and Richard Pethia, January 1997,
http://www.cert.org/pres_comm/cert.rpcci.body.html
- [52] "Cisco Quality of Service and DDOS, Engineering Issues for Adaptive Defense Network," MITRE, 7/25/2001.
- [53] "Denial of Service Attacks," CERT Coordination Center, June 4, 2001,
http://www.cert.org/tech_tips/denial_of_service.html
- [54] Tom Chen, "Trends in Viruses and Worms," *The Internet Protocol Journal*, Volume 6, No. 3, September 2003.

CHARALAMPOS Z. PATRIKAKIS holds a Dipl.-Ing. and a Ph.D. degree from the Electrical Engineering and Computer Science Department of the National Technical University of Athens (NTUA). He is currently a senior research associate of the Telecommunications Laboratory of NTUA. He has participated in several European Union projects (ESPRIT, RACE, ACTS, IST). His main interests are in the area of IP service design and implementation, multicasting in IP networks, IP transport protocols, and media streaming over IP networks. He is a member of IEEE, a member of the Greek Computer Society, a certified trainer by the National Accreditation Centre of Vocational Training Structures and Accompanying Support Services, and a member of the Technical Chamber of Greece. He can be reached at: bpatr@telecom.ntua.gr

MICHALIS MASIROS holds a Dipl.-Ing. degree from the Electrical Engineering and Computer Science Department of the National Technical University of Athens (NTUA). He is currently a research associate of the Telecommunications Laboratory of NTUA. His interests are in the fields of network security, network simulation, and analysis. He can be reached at: mmasik@telecom.ntua.gr

OLGA ZOURARAKI holds a Dipl.-Ing. degree from the Electrical Engineering and Computer Science Department of the National Technical University of Athens (NTUA). She is currently a research associate of the Telecommunications Laboratory of NTUA. Her interests are in the fields of network security, Internet application design, and implementation. She can be reached at: ozour@telecom.ntua.gr

Letter to the Editor

Ole,

I was reading your latest issue of IPJ (Volume 7, No. 3, September 2004) and I could be wrong but I think you mis-typed an explanation about the STUN protocol. On page 12, 3rd paragraph, last sentence, it reads: “A received response indicates the presence of a port-restricted cone, and the lack of a response indicates the presence of a restricted cone.”

According to the definitions you gave about “restricted cone” and “port-restricted cone” on pages 10 and 11. Shouldn’t this sentence instead read: “A received response indicates the presence of a restricted cone, and the lack of a response indicates the presence of a port-restricted cone.”

—Ryan Liles
ryanliles@hotmail.com

The author responds:

Ryan is correct, there is an error here in the text.

The flow control of the sequence of STUN tests is detailed in Figure 9 of the article. The test referred to here is to determine if the NAT is a restricted cone NAT, or a port-restricted cone NAT.

The restricted cone NAT, in Figure 7, is one where the NAT binding is accessible using any source port number on the external host when responding to a UDP packet from the internal sending host.

The port-restricted cone NAT, in Figure 8, is one where the NAT binding is accessible using the same port number as originally used by the internal host, and this binding is accessible from any external IP address.

The test referenced in this section, as per Figure 9, is one where the local host requests the external agent to respond using the same port number, but an altered source address. The text should read “This fourth request includes a control flag to direct the STUN server to respond using the alternate IP address, but with the same port value,” in which case the interpretation of the response—that a response indicates the presence of a port-restricted cone NAT and the lack of response indicates the presence of a restricted cone NAT—would be correct.

Ryan is also correct in that if the test is performed the other way, requesting the agent to use the same IP address, but with the alternate port value, then the opposite interpretation would hold, namely that a response indicates the presence of a restricted cone NAT, and the lack of a response would indicate the presence of a port-restricted cone NAT, as Ryan points out.

Thanks to Ryan for following through this rather complex explanation of the STUN algorithm and spotting this error.

Regards,

—Geoff Huston, APNIC
gih@apnic.net

Book Review

The IP Multimedia Subsystem

The IP Multimedia Subsystem—Merging the Internet and the Cellular Worlds, by Gonzalo Camarillo and Miguel A. Garcia-Martin, John Wiley & Sons, 2004. ISBN 0470 87156 3.

The Internet and the cellular telephony system are the two most influential communication systems of the last half century. That the telecommunications industry would attempt to merge them into a single system was inevitable. The potential benefits are compelling—a single packet-based communication system with the capability to carry voice, video and data while providing ubiquitous wireless access and global mobility. The resulting system architecture is called the *Internet Multimedia Subsystem* (IMS) and is described comprehensively in this volume by Gonzalo Camarillo and Miguel A. Garcia Martin.

A “merging” of the two systems is only superficially what has happened. In practice, the IMS is an “embrace and extend” exercise which adapts the IP protocol suite to the existing architecture of the cellular telephony system. The cellular industry has taken a broad collection of IP protocols and mapped them onto their existing architecture, effecting a “protocol transplant” into an environment somewhat different from the Internet. Among the protocols imported are IPv6, SIP, DHCP, DNS, SDP, RTP, IPSec, and DIAMETER. Many are adopted unaltered; some are profiled by introducing new configuration data and rules; others are extended in various ways. The authors navigate their way through the various parts of the system with clarity and confidence. They can speak with authority on the subject—both were major contributors to the design through their key roles in the IETF and 3GPP (*Third Generation Partnership Project*—the standardization body for third generation cellular systems).

The book is clearly written and logically organized. The first part explains the reasoning behind adopting Internet-style packet networking for cellular mobile systems and describes the evolution of the standardization efforts. Although interesting, much of this material can be skimmed by those only interested in the meaty technical material which follows. The authors then explain the general principles behind the IMS architecture, including how various requirements of the cellular telephony industry drove the choices, and particularly the perceived need to extend and adapt the protocols rather than use them as deployed on the Internet. The majority of the book is devoted to explaining in considerable technical depth how the protocols have been modified and how they are intended to work when IMS is successfully deployed. While not for the faint of heart, the writing is extremely clear and logical and hence should be understandable by anyone with a moderate background in the principles of protocol and system design. One aspect of the organization is particularly helpful to readers unfamiliar with some of the protocols in their native Internet instantiation. The authors divide the material into blocks where they first describe the native Internet flavor of the protocol, and then introduce the IMS-specific extensions and modifications.

Much of the volume is devoted to the *Session Initiation Protocol* (SIP) as the core signaling plane for IMS. All aspects of session establishment and management are covered. In addition, the ancillary parts of the control system are covered, including *Authentication, Authorization, and Accounting* (AAA), Security, Session Policies, and Quality of Service. For completeness, the data plane is also covered briefly through a discussion of the 3GPP audio, video, and text encoders, plus material on the media transport protocols.

The book concludes with a substantial section on how services are build on top of the core IMS protocols. Two of the most important, *Presence* and *Instant Messaging*, get comprehensive treatment, with a briefer discussion of the push-to-talk application.

As an old time “IP-head,” it is hard to come away from this deep exploration of IMS without a bit of trepidation. The hallmark of IP and the Internet are simplicity and generality. IMS arguably succeeds at the latter, but at the expense of almost numbing complexity. This was perhaps inevitable given that the goal was to adapt Internet packet technology to the cellular system, which is itself quite complex. IMS will be quite a challenge to deploy. It remains to be seen if transplanting IP into a cellular telephony architectural model will result in economically sustainable services for the service providers or if a more native peer-to-peer Internet approach will simply bypass all the fancy IMS elements and just use basic packet transport. Such a market experiment is currently playing out in the broadband access arena with the broadband pipe suppliers offering telephony-oriented services themselves via customized standards like PacketCable, while third parties like Vonage and Skype simply piggyback on basic IP packet transport.

The next few years will be interesting. Whatever the outcome, anyone needing to be technically conversant with the architecture and protocols of IMS will find *The IP Multimedia Subsystem* indispensable.

—David Oran
oran@cisco.com

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. Contact us at **ipj@cisco.com** for more information.

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Technology Strategy
MCI, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.
Copyright © 2004 Cisco Systems Inc.
All rights reserved. Printed in the USA.*



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PSN STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA
--

The Internet Protocol Journal

March 2005

Volume 8, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Misbehaving Name Servers	2
Wireless Networks.....	6
Internet Governance	15
Book Review.....	29
Fragments	32
Call for Papers	35

FROM THE EDITOR

Internet Protocol Version 6 (IPv6) continues to be the focus of much work within the IETF as well as throughout the world in numerous deployment projects. The success of IPv6 depends not only on the protocol itself but also on its interaction with existing services such as the *Domain Name System* (DNS). In our first article, David Malone looks at some issues with DNS servers and IPv6. If you are interested in following the progress of IPv6 deployment, you might want to visit The IPv6 Forum's Website at: <http://www.ipv6forum.org>

A couple of years ago I signed up for GSM cellphone service and later added GPRS data service to my account. With my Bluetooth-enabled phone and laptop, I can access the Internet from almost anywhere in the world. The service is neither particularly fast nor inexpensive, but for occasional use it works very well, and has "saved the day" for me numerous times. However, GPRS is not the only wide-area wireless data network technology. Kostas Pentikousis gives an overview of the many alternatives.

The term "Internet Governance" is not well-defined, but it is being used more frequently when speaking about such organizations as the *Internet Corporation for Assigned Names and Numbers* (ICANN). The formation of the *World Summit on the Information Society* (WSIS) and its *Working Group on Internet Governance* (WGIG) has certainly brought the term into sharper focus. Although governance is certainly not a technical protocol issue, we still believe that it is important for our readers to follow both the debate about and the actual evolution of Internet Governance issues. However, we fully appreciate that this is an area where opinions differ—and that is why the article by Geoff Huston on this topic is labeled "Opinion."

We remind you to visit our Website, <http://www.cisco.com/ipj>, where you can find back issues of this journal, search the index files, or make changes to your subscription information. Your feedback is also very much appreciated, so drop us a line at ipj@cisco.com

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Misbehaving Name Servers and What They're Missing

by David Malone, Hamilton Institute, NUI Maynooth, Ireland

IPv6-capable hosts abound, and the number is growing. Evidence^[1] shows that more than 2 million Windows XP machines are probing for 6to4^[2] connectivity. When combined with deployments of Linux and BSD that have been shipping with IPv6 support enabled by default for some time, that is a sizable platform on which to build IPv6 applications. Most Web browsers (Internet Explorer, Mozilla, Opera) now support IPv6 if the underlying platform does, so that is a significant number of applications ready to start making IPv6 queries.

In fact, many of these applications are already looking for IPv6 addresses in the *Domain Name System* (DNS), even if IPv6 connectivity is not actually available. This usually does not result in a problem—the name server says there are no IPv6 records and the application falls back to IPv4. In a small number of cases, name servers running outdated or errant software are misbehaving when faced with a request for an IPv6 address.

The Problem

So, what problem are these name servers having with the request for IPv6 addresses? Well, the DNS stores different types of information, such as host names and addresses. Different types of data are stored using different record types. For example, IPv4 addresses are stored using a type “A” record and host names are stored using a type “PTR” record. Some new record types have been introduced for IPv6. The most important one is “AAAA,” which is for storing IPv6 addresses. (Another type called “A6” was also introduced, but it is now consigned to experimental status because it proved too complicated in certain situations.)

When you issue a request to the DNS, you indicate the domain and type of record that you are interested in. If the server has records of that type for that domain, it replies, including those records. If the server has no records of that type, it should respond saying “there are no records of this type.” If the domain does not exist, then the server should return a “no such domain” error.

However, the problems arise when the DNS server does something different, and some name servers behave badly when faced with a query for a type they do not explicitly know about. For the sake of simplicity, we will highlight three wrong reactions to an unknown query that have been observed. A more complete technical analysis of the problem can be found in^[3].

The first reaction that people notice is that some name servers do not reply when faced with a query for an unknown type. In this case, the person who made the request waits a while before the request is reissued. Eventually the application falls back to IPv4. “Eventually” means anything from 10 seconds to 100 seconds, depending on the operating system and application—enough to irk the casual Web user.

The second reaction is more subtle. Here the name server returns a “no such domain” response. At first glance this may seem harmless enough—the query for an IPv4 address is issued quickly. However, DNS specifications say that the “no such domain” response may be cached. This means that the “A” query is never issued, and the system acts as if the domain does not exist.

The third reaction is that the server issues some other sort of incorrect response. Usually this is less serious than the two previous reactions, because other responses at worst result in a particular name server being considered “bad” and being avoided for future queries. This means that some better-behaved name server can answer the query.

The Extent of the Problem

Although sites with these problems are sometimes discussed on mailing lists, the extent of a problem is not always proportional to the coverage it receives. Historically, numerous online advertising companies have had load-balancing DNS servers that exhibit these symptoms. Because the content of an ad server is embedded in the Web pages of many organizations, this means a single errant DNS server can give the end user the impression that this problem is more widespread than it is.

To give some idea of the scale of the problem, Table 1 shows the results of querying the name servers for the names mentioned in a month’s worth of Web proxy logs. The number of servers responding in each of the three ways mentioned (no reply, no such domain, or other error) is shown, along with a total. Also shown is the number of name servers that actually returned IPv6 addresses.

These results show that actually only a small number of name servers have this problem. Unfortunately, it also looks as if the number of name servers distributing IPv6 addresses is actually comparable. However, it does look like the proportion of problem name servers is decreasing over time.

Table 1: Responses to Name Queries

Nameservers that:	January 2004	April 2004	August 2004
<i>Responded to type A</i>	<i>16838</i>	<i>20631</i>	<i>17934</i>
<i>Did not reply to type A</i>	<i>64 (0.38%)</i>	<i>49 (0.24%)</i>	<i>36 (0.20%)</i>
<i>Returned no such domain</i>	<i>11 (0.07%)</i>	<i>19 (0.09%)</i>	<i>11 (0.06%)</i>
<i>Returned other error</i>	<i>22 (0.13%)</i>	<i>39 (0.19%)</i>	<i>11 (0.06%)</i>
<i>Had any issue with AAAA</i>	<i>97 (0.58%)</i>	<i>107 (0.52%)</i>	<i>58 (0.32%)</i>
<i>Returned AAAA records</i>	<i>105 (0.62%)</i>	<i>123 (0.60%)</i>	<i>18 (0.66%)</i>

Looking at Web logs to determine the size of the problem gives us a feeling for the number of name servers that need attention. Another interesting parameter to consider is the proportion of requests that might be subject to this problem. The answer would tell us how many queries might be mishandled if your name server cannot deal with new query types.

Looking at the queries for addresses at one authoritative name server shows that 65 percent of queries are for A records, 21 percent are for AAAA records, and 14 percent are for A6 records. Although this server is IPv6-capable and might attract more queries for AAAA records, even the root servers run by RIPE show that 10 percent of address queries are for IPv6 addresses.

The Solution

Some of the name servers that exhibit this problem are simply running old versions of DNS server software. If this is the case, then the fix is simple: *upgrade!*

A significant number of the remaining problem servers are running unusual name server software, and the only way to fix the problem is to have that software fixed. Where the name server software is maintained in house, there should be enough DNS expertise to resolve the issue when it is identified. Where DNS systems have been bought in, it can be difficult to get the relevant information to the developers who can make the necessary changes. Thus increasing awareness of the issue among DNS vendors and troubleshooters is important.

In some cases^[5,6], discussions on Internet mailing lists has alerted those responsible for the server to the problem and the issue has been resolved. In other cases, feedback provided by users and customers has marked IPv6 conformance as an issue for future upgrades of a site's DNS infrastructure. Unfortunately, on some occasions, feedback has been ignored and the problem has persisted. This is maybe not so surprising because it is a subtle problem. The fact that it is IPv6-related means it is sometimes dismissed because the organization thinks "we have not begun IPv6 deployment yet, so it cannot affect us."

Where problems have persisted, people have resorted to various practical solutions (hacks?) to avoid the issue. Some people, who do not need IPv6 at this time, have just suppressed the AAAA queries. Others, when they discover a name server that times out, add it to a blacklist. This avoids any delays, but may make a site unavailable. Mozilla includes a more forgiving style of blacklisting, in the form of a "ipv4OnlyDomains" setting, that can be set to a list of domains known to have problems^[7].

The long-term solution seems straightforward. As we have seen, the number of name servers exhibiting this problem is relatively small, though some do serve some often-queried domains. If we can ensure that no more servers with these problems get deployed, then as the existing servers are updated or retired the problem will be resolved.

To this end, it is worth testing new DNS deployments to make sure that they correctly respond to unusual query types^[8]. This will smooth the path not just for IPv6, but also for other new technology such the *Domain Name System Security Extension* (DNSSEC)^[9].

References

- [1] “Observations of 6to4 Traffic on a 6to4 Router,” Pekka Savola, preprint, October 2004.
- [2] “Connecting IPv6 Routing Domains Over the IPv4 Internet,” Carpenter, Moore, and Fink, *The Internet Protocol Journal*, Volume 3, No. 1, March 2000.
- [3] “Common Misbehavior against DNS Queries for IPv6 Addresses,” Y. Morishita and T. Jinmei, **draft-ietf-dnsop-misbehavior-against-aaaa-01.txt**, April 2004.
- [4] K-Root information page, RIPE, **<http://k.root-servers.org/>**
- [5] “**news.bbc.co.uk** NXDOMAIN problem fixed,” itojun, Simon Lockhart, et al., 6bone mailing list, April 2002.
- [6] “**ftp.perl.org** strangenes” thread, Mark Andrews, Ask Bjoern Hansen, et al., freebsd-stable mailing list, March 2004.
- [7] “IPv6: Some IPv4 addresses won’t resolve w/IPv6 OS,” Mozilla bug 68796,
https://bugzilla.mozilla.org/show_bug.cgi?id=68796
- [8] “AAAA lookup checker,” David Malone,
http://www.cnri.dit.ie/cgi-bin/check_aaaa.pl
- [9] “DNSSEC,” Miek Gieben, *The Internet Protocol Journal*, Volume 7, No. 2, June 2004.

DAVID MALONE received B.A. (mod), M.Sc., and Ph.D. degrees from Trinity College Dublin. He has been involved with system administration since 1994 and has been slowly growing IPv6 networks since 1999, when he also became a FreeBSD committer. With Niall Murphy, he is the coauthor of *IPv6 Network Administration*, ISBN 0-596-00934-8 published by O’Reilly and Associates, 2005. He is currently on secondment to the Hamilton Institute of NUI Maynooth. E-mail: **dwmalone@maths.tcd.ie**

Wireless Data Networks

by Kostas Pentikousis, VTT

Most IPJ readers are familiar with *Wireless Local-Area Networks* (WLANs; see, for example, IPJ Volume 5, No. 1). Some may even be familiar with recent developments in *Wireless Metropolitan-Area Networks* (WMANs), such as WiMAX. Although nonproprietary WMAN technologies are still in the standardization phase, the IEEE 802.11 family of protocols has reached maturity and rendered inexpensive (and often free) WLAN access increasingly popular. Both WLANs and WMANs provide high-speed connectivity (in the order of tens of Mbps), but user mobility is restricted. In fact, it is probably more appropriate to talk about “portability” rather than “mobility”^[1] when referring to WLANs and WMANs.

Wireless wide-area networks (WWANs), on the other hand, allow full user mobility but at data rates typically in the order of tens of kbps. This will change to some extent when *third-generation* (3G) cellular networks are fully deployed. Still, 3G deployment is slower than originally anticipated, a development often attributed to the combination of high spectrum license costs, the recent economic downturn, and high equipment costs. As a result, both population and geographical coverage tend to be uneven. For example, in Finland, a forerunner in wireless communications, population coverage is well below the 35-percent level, and geographical coverage is even smaller.

This article introduces several wireless network technologies, perhaps not so widely known, which deserve attention when considering how to provide mobile connectivity to field personnel, introduce *machine-to-machine* (M2M) communication, or deploy applications that require always-on connectivity. The approach taken in this article is a bit different from the one typically followed in the literature: We focus more on higher-level issues, the information that is essential for application developers, instead of modulation, channel coding, and other low-level details. Unlike WLANs and WMANs, none of the networks surveyed provide data rates in the order of tens of Mbps. Nevertheless, successful applications can be built even with stringent bandwidth limitations. For example, online gambling and several gaming applications can be served by really “thin” networks (and possibly “thick” clients).

Cellular Networks

The *Global System for Mobile Communications* (GSM) specifies a cellular, wide-area, circuit-switched, digital mobile phone network architecture^[2]. Circuit-switched networks such as GSM and IS-95, commonly referred to as *Code Division Multiple Access* (CDMA) in the United States, can provide wireless data connectivity, cover a large area, and handle mobile host handovers efficiently^[3]. Users can transfer data over, say, GSM, by establishing a “dialup” connection^[4]. Mobile hosts can roam, even at high speeds, and remain connected throughout.

Communication is full-duplex at a radio data rate of 9.6 kbps or 14.4 kbps in GSM Phase 2+^[5]. User throughput is always smaller than the nominal radio data rate.

While the user is connected using a wireless circuit-switched network, phone calls cannot be initiated or received whether data is being transferred or not. This is not much different from wire-line dialups over basic telephone service. The difference is that a dialup over a *Public Switched Telephone Network* (PSTN) takes up a resource, namely the wire-line local loop, which is dedicated to a single user, whereas a dialup over a cellular network such as GSM consumes a resource, the radio channel, which is shared among many users. Because of the *burstiness* that data traffic usually exhibits, circuit switching may lead to inefficient use of the network capacity. Establishing a GSM dialup connection usually takes several seconds, meaning that if the user has a small amount of data to send, a small e-mail message, for example, the overall experience is poor. Moreover, after the connection is established, the channel remains idle between traffic bursts and the allocated bandwidth is wasted. Packet switching is more efficient for bursty data transmission over a shared medium^[6].

Another variable that favors packet-switching over circuit-switching, especially over slow wireless networks, is *billing*. Users of circuit-switched networks are usually charged based on the duration of a connection regardless of the amount of traffic transmitted or received. On the other hand, users of packet-switched networks can be charged based solely on the amount of data transferred—not how long they remain attached to the network. In short, introducing packet switching to wireless networks can lead to better use of network resources and attract more users as data transfers become more economical.

Two-way, packet-switched WWANs permit users to roam freely indoors and outdoors, even at relatively high speeds^[7]. Most WWANs employ a cellular architecture to take advantage of frequency reuse and increase capacity while covering a larger area. Furthermore, because the coverage area of a single cell is generally large (cell diameters are typically in the order of dozens of kilometers), mobile hosts do not have to go through frequent and lengthy handovers. Hosts remain connected throughout after they attach to the network, permitting users to receive and transmit data on demand without having to dial up. The following sections survey some of the most widely deployed packet-switched wireless data networks.

Mobitex

Mobitex is the first digital data-only WWAN developed by Ericsson and Swedish Telecom. Not based on IP, Mobitex was introduced in Sweden in 1986 for emergency communications^[8]. It uses a cellular architecture with cell diameters of up to 30 km. Each service area can operate 10–30 channels^[9] and each base station is usually allocated 1 to 4 channels. Each channel is composed of a frequency pair: different frequencies are used for the uplink and the downlink.

Communication between the base station and a single mobile host is, nevertheless, effectively half-duplex. Although base stations can transmit and receive simultaneously, mobile nodes are unable to do so^[10]. The Mobitex *Maximum Transmission Unit* (MTU) is 545 bytes, with up to 512 bytes of user data. Although the system has undergone several revisions, the raw transfer rate remains only 8 kbps. Effective user throughputs range from 4 kbps (for 125-byte packets) to 4.6 kbps (for 512-byte packets)^[11], and round-trip times can be up to 10 seconds.

Mobitex deals with network lapses using a store-and-forward procedure: Packets destined for a mobile node outside the network coverage area are stored while awaiting delivery. When the mobile node reconnects, the stored packets are delivered. Mobitex uses a hierarchical routing architecture that prevents local traffic from being injected into the backbone network. In other words, packets destined for a node in the range of the same base station are switched locally^[8]. Besides supporting unicast addressing, Mobitex allows hosts to send one packet to several recipients^[10]. According to the *Mobitex Association* (www.mobitex.org), the technology features “true push functionality,” whereby data can be pushed to both a single mobile node and a predefined group of nodes, a feature that can be very useful when trying to send an urgent message to field personnel. And, because the mobile host does not have to keep querying for pending data, network traffic can be kept to a minimum. All these features can also significantly boost battery life.

According to the Yankee Group, despite the limited data rates, a variety of applications have been developed based on Mobitex, including: burglar and fire alarm systems; paging, interactive messaging, e-mail, form-based applications, and access to databases; telemetry; credit card authorizations; field service; and fleet management. Virtually all of them require small and bursty transfers. Mobitex does not lend itself to large file transfers, e-mail with large attachments, or video transmission. In fact, file transfers of more than 20 KB used to be discouraged^[8]. On the other hand, by using a slotted ALOHA^[12] variation for channel access, Mobitex can provide message delivery delay guarantees and support hundreds of users within the same cell. Parsa^[13] calculated that Mobitex can accommodate 2,000 users per channel, assuming two uplink and two downlink messages per hour. Other networks simply cannot provide tight delay bounds for such a large number of users. For example, the *Mobile Data Magazine* (No. 1, 2002) reported that a Korean operator launched real-time stock trading and horse gambling mobile applications with great commercial success, by guaranteeing delay bounds notwithstanding the low data rates.

DataTAC

DataTAC (also known as ARDIS in the United States) was developed by Motorola in the mid-1980s. DataTAC is also a non-IP based, wide-area, data-only message-oriented network. A single base station can cover an area exceeding 20 km in diameter^[14]. Like Mobitex, communication between the base station and a single DataTAC mobile node is half-duplex, and mobile hosts have to compete to get access to transmit and receive data.

Unlike Mobitex, DataTAC was designed to provide optimal in-building coverage, and it uses a cellular architecture that does not take advantage of frequency reuse. Instead, a single frequency is used, increasing the probability that a packet transmission is successful (because the same transmission can be picked up by more than one base station), but at the expense of network capacity^[8]. Bodsky notes that the U.S. DataTAC operator formerly recommended refraining from transferring files larger than 10 KB.

Although neither Mobitex nor DataTAC provides native IP support, middleware can take care of protocol translation and allow unmodified, off-the-shelf applications to communicate. The maximum Data-TAC message size is 2048 bytes^[15], but the maximum over-the-air packet size depends on the link layer. For rural areas the maximum radio data rate is 4.8 kbps, and the maximum over-the-air packet size is 256 bytes. In metropolitan areas, the radio data rate is 19.2 kbps and the maximum packet size is 512 bytes^[16]; end-user throughput does not exceed 10 kbps on average. Traditionally, DataTAC was used for dispatching and law enforcement applications. The *Worldwide Wireless Data Network Operators Group* (www.datatac.com) reports that DataTAC networks are also used for two-way messaging, wireless e-mail, telemetry, access to corporate databases, and package tracking by courier carriers.

CDPD

Cellular Digital Packet Data (CDPD) was designed by IBM and McCaw Cellular Communications in the early 1990s to take advantage of channels that do not carry voice traffic in the *Advanced Mobile Phone Service* (AMPS), the first-generation analog cellular network^[17]. Data channels are allocated dynamically, sharing the network capacity with AMPS voice traffic, which is quite different from Mobitex and DataTAC. This, for example, might mean that data can be transmitted and received only when phone calls do not consume all available capacity. One could argue that CDPD considers data traffic less important than voice. However, the standard allows network operators to specifically assign channels to data traffic only. In theory, deployment can be more economical than it is for other WWANs because CDPD takes advantage of existing AMPS infrastructure and does not require licensing new spectrum. Original projections anticipated that as CDPD gained popularity—and AMPS became obsolete—more CDPD dedicated channels would be allocated. With time, CDPD would have taken over the existing AMPS bandwidth, effectively becoming a data-only WWAN.

CDPD is based on a *Carrier Sense Multiple Access* (CSMA) variant called *Digital Sense Multiple Access*^[14] and transparently provides IP services, constituting a great advantage. CDPD allows for an MTU of 2048 bytes. However, one has to account for the *TCP/User Datagram Protocol* (UDP) and IP headers that are used to encapsulate the application payload before sending it over the CDPD network and also for the fact that CDPD user data is transmitted in much smaller blocks. Although the CDPD raw data rate is 19.2 kbps, the effective throughput is in the order of 10 kbps and response times have been reported to be in the order of 4 seconds^[18].

GPRS

The *General Packet Radio Service* (GPRS) is overlaid on a GSM network in a fashion similar to the way CDPD is embedded in AMPS: Voice and data traffic share the same bandwidth and network infrastructure^[14]. In other words, GPRS is an add-on to GSM networks, and it requires certain hardware and software upgrades and introduces packet switching to a circuit-switched architecture. GSM voice traffic is oblivious to the presence of GPRS data traffic. Similar to CDPD, GPRS is designed to appear as a regular IP subnetwork both to hosts attached over the air interface and to hosts outside the GPRS network.

The GPRS standard was finalized by the *European Telecommunications Standards Institute* (ETSI) in late 1997 as part of GSM Phase 2+^[5]. It is regarded as a transitional technology toward 3G networks^[19], and is commonly referred to as 2.5G. One of its main advantages is that the same device can be used to transmit and receive data, and initiate and accept phone calls. GPRS defines three classes with respect to simultaneous usage of voice and data. Class A mobile hosts can transmit and receive voice and data at the same time. Class B hosts can transmit and receive either voice or data but not both simultaneously. Finally, class C hosts have the user manually select if the host should be attached to the GSM (voice) or GPRS (data) network. When compared to Mobitex, DataTAC, and CDPD, GPRS class A devices can have simultaneous access to a packet-switched and circuit-switched network. Of course, GSM-only devices do not have this capability either, as mentioned earlier.

GSM uses a combination of *Frequency Division Multiple Access* (FDMA) and *Time Division Multiple Access* (TDMA) for channel allocation, as explained in detail in^[5]. In short, each frequency channel carries eight TDMA channels. Each of these channels is essentially a time slot in a TDMA frame. Thus, any GSM frequency channel can carry up to eight circuit-switched connections with each slot reserved for a single connection (read *voice call*). In GPRS, each slot is treated as a shared resource and any mobile host can use it to transmit or receive data. In addition, a mobile host can be allocated more than one of the eight available slots in the same TDMA frame. In other words, GPRS can multiplex different traffic sources in one channel and allocate several channels to the same traffic source.

GPRS defines four different channel coding schemes^[20], namely CS1, CS2, CS3, and CS4, with radio data rates 8.8 kbps, 13.3 kbps, 15.6 kbps, and 21.4 kbps, respectively. CS1 is the most “conservative” (includes more error correction bits) and is used for signaling packets and when poor channel conditions prevail. CS4 is the most “optimistic” (includes minimal error correction bits), and, assuming excellent channel conditions, allows operators to advertise a maximum radio data rate of 171.2 kbps per 200-kHz frequency channel (or TDMA frame).

In practice, CS4 is rarely used because it can lead to frequent retransmissions of lost packets and overall network underperformance. CS3 is commonly used, providing 124.8 kbps per frequency channel. Because a mobile host can be allocated multiple slots, user throughputs can range between 40 and 60 kbps. Mobile hosts typically use an MTU of 1500 bytes.

Communication between the base station and any given mobile host is full-duplex but can be *asymmetric*; that is, the downlink and uplink capacities need not be the same. The *GSM Association* has defined 12 multislot classes for GPRS. Each class is associated with a maximum number of uplink and downlink slots that can be allocated to a single mobile host. The slot allocation is usually written as $M + N$, where M is the maximum number of downlink slots and N is the maximum number of uplink slots. For example, class 1 is “1 + 1” (one downlink slot plus one uplink slot); class 2 is “2 + 1”; . . . ; and class 12 is “4 + 4” (four downlink and four uplink slots). In addition, each multislot class has an active slot constraint: A mobile host cannot use more than K active slots simultaneously. Given the number of slots and the channel coding scheme, one can calculate the peak rate. For example, for a class 12 device the sum of the physical downlink and uplink rates cannot exceed 124.8 kbps, if CS3 is used. However, the active slot constraint limits this rate even further. In the case of a class 12 mobile node, $K = 5$, that is, only “4 + 1”, “3 + 2”, “2 + 3”, or “1 + 4” slots can be used simultaneously. See www.gsmworld.com

EDGE and Beyond

Enhanced Data for GSM Evolution (EDGE), also known as Enhanced GPRS, builds on the changes introduced by GPRS to GSM. EDGE essentially increases the radio data rates by using a more efficient modulation scheme^[21], namely *8-Phase Shift Keying* (8-PSK) instead of the *Gaussian Minimum Shift Keying* (GMSK) used by both GSM and GPRS. EDGE defines nine modulation coding schemes named MCS1 to MCS9. MCS1 to MCS4 use GMSK with radio data rates similar to the four GPRS coding schemes. The real throughput improvements come from MCS6 (29.6 kbps per slot) through MCS9 (59.2 kbps per slot). The data rate usually associated with EDGE is a (shared) 384 kbps. This corresponds to using MCS7 for all 8 TDMA slots. Higher data rates are theoretically possible (up to 473 kbps using MCS9) but are not commonly deployed.

EDGE improves not only on the high end of data rates but also on the low end^[22]. First, the greater diversity of coding schemes permits an EDGE network to choose the most appropriate one depending on channel conditions. Changing coding schemes is dynamic. Second, EDGE supports *packet resegmentation*: Packets that failed to be transmitted successfully can be resegmented and retransmitted using a more “conservative” coding scheme.

Table 1 summarizes the main high-level features for the WWANs surveyed.

Table 1: WWAN Characteristics

	Transmit/ Receive	Radio Data Rate	User Throughput	MTU
<i>Mobitex</i>	<i>Half duplex</i>	<i>8.0 kbps</i>	<i><4.6 kbps</i>	<i>512 B</i>
<i>DataTAC</i>	<i>Half duplex</i>	<i>19.2 kbps</i>	<i><10 kbps</i>	<i>2048 B*</i>
<i>CDPD</i>	<i>Full duplex</i>	<i>19.2 kbps</i>	<i><10 kbps</i>	<i>2048 B</i>
<i>GPRS</i>	<i>Full duplex</i>	<i><171 kbps</i>	<i>40–60 kbps</i>	<i>1500 B</i>
<i>EDGE</i>	<i>Full duplex</i>	<i><473 kbps</i>	<i>50–60 kbps</i>	<i>1500 B</i>

* Typically 512 B

Discussion and Trends

Among the WWANs presented, Mobitex and GPRS can be singled out as the most widely deployed; they also have enjoyed significant gains in the number of users and traffic volume in recent years. The popularity of enterprise wireless e-mail (due in part to the success of the Research in Motion BlackBerry devices) allowed Mobitex and DataTAC operators to revive their business models briefly. Worldwide, however, GSM dwarfs all other technologies: There are more than 1 billion GSM subscribers compared to the 1 million Mobitex users. DataTAC enjoys an even smaller user base. Even if a small percentage of GSM subscribers use GPRS and EDGE, the potential market for wireless applications is tremendous. On the other hand, subscribers who do not take advantage of GPRS or EDGE do use the inexpensive, (two-way) *Short Message Service* (SMS), which is built in GSM. Two-way messaging was available for many years but was certainly popularized by less-affluent and younger GSM users in the late 1990s. SMS is now commonplace, and in many countries it is more popular than e-mail. Dedicated data-only networks such as Mobitex have to look elsewhere for their niche.

For some, Mobitex, let alone DataTAC and CDPD, is virtually moribund. In the United States, for example, Cingular sold its Mobitex network and is investing heavily on GPRS and EDGE. DataTAC and CDPD are phased out by service providers in the United States in favor of newer technologies. Low-speed packet radio is considered lackluster and is not popular with younger crowds. After all, narrowband WWANs had their chance and failed to attract large numbers of subscribers. Recent pricing trends, too, reveal a heavy operator push in favor of GPRS and EDGE. In Finland, for example, 100 MB over GPRS costs less than 18 euros (approximately \$24). Compare that to the \$30–50 that 1 MB of traffic costs over Mobitex. Service and product popularity create economies of scale that cannot be ignored.

Nonetheless, open standards, an explicit focus on business applications with *Quality-of-Service* (QoS) guarantees in service response times, and narrowband M2M communication may well keep Mobitex going for years to come. Besides, bundling Mobitex with a wireless network that features fast and inexpensive connectivity, for example, WLAN or Bluetooth, might be promising: Large downloads and software updates can be done over the high-speed wireless network and critical messages can always reach the user through the WWAN.

Bundling several functions in a single handheld device is, after all, a major trend in the industry. Vendors scramble to integrate *Personal Information Managers* (PIMs), voice and data communications, as well as entertainment features (digital camera, games, or digital music players) in a single product. This is quite different from earlier mobile devices, which tended to be either single-purpose or tied to a particular set of applications. Even the BlackBerry devices still work, to some extent, in a closed architecture. Enterprise e-mail systems need to be supported by and integrated with BlackBerry servers in order to be accessible over the WWAN. Yet, one of the main objectives in 2.5G and 3G is to allow mobile users to use standard Internet protocols on a mobile radio network at significantly higher bit rates than other systems. In particular, GPRS was designed with certain office applications in mind and can support consumer and enterprise mobile communications alike, without being tied to any given platform or application servers. I expect that functionality bundling and 2.5G and 3G WWANs will allow for more open systems and will expedite the transformation of WWAN operators from integrated application providers to wireless ISPs.

For Further Reading

- [1] Charles Perkins, *Mobile IP Design Principles and Practices*, ISBN 0201634694, Addison-Wesley, 1998.
- [2] Joachim Tsal, *GSM Cellular Radio Telephony*, ISBN 0471968269, John Wiley & Sons, 1998.
- [3] Tero Ojanpera and Ramjee Prasad, *WCDMA: Towards IP Mobility and Mobile Internet*, ISBN B0000660B4, Artech House, 2001.
- [4] R. Ludwig, B. Rathonyi, A. Konrad, K. Oden, and A. Joseph, "Multi-layer tracing of TCP over a Reliable Wireless Link," presented at ACM SIGMETRICS 1999.
- [5] C. Bettstetter, H.-J. Vogel, and J. Eberspacher, "GSM phase 2+ General Packet Radio Service GPRS: Architecture, Protocols, and Air Interface," *IEEE Communications Surveys & Tutorials*, Vol. 2(3), pp. 2–14, 1999.
- [6] Larry L. Peterson and Bruce S. Davie, *Computer Networks: A Systems Approach*, 3rd ed., ISBN 155860832X, Morgan-Kaufmann, 2003.
- [7] Rudi Bekkers and Jan Smits, *Mobile Telecommunications: Standards, Regulation, and Applications*, ISBN 0890068062, Artech House, 1999.

- [8] Ira Brodsky, *Wireless: The Revolution in Personal Telecommunications*, ISBN 089006717, Artech House, 1995.
- [9] Nathan J. Muller, *Wireless Data Networking*, ISBN 0890067538, Artech House, 1995.
- [10] A. K. Salkintzis and C. Chamzas, "Mobile Packet Data Technology: An insight into Mobitex Architecture," *IEEE Personal Communications*, Vol. 4(1), pp. 10–18, 1997.
- [11] M. S. Taylor, M. Banan, W. Waung, and M. Taylor, *Internet-work Mobility: The CDPD Approach*, ISBN 0132096935, Prentice-Hall, 1996.
- [12] Andrew S. Tanenbaum, *Computer Networks*, 4th ed., ISBN 0130661023, Pearson Education, 2003.
- [13] K. Parsa, "The Mobitex packet-switched Radio Data System," presented at IEEE PIMRC '92, 1992.
- [14] Sami Tabbane, *Handbook of Mobile Radio Networks*, ISBN 1580530095, Artech House, 2000.
- [15] J. Rodriguez, W. Schollenberger, M. Anzib, and B. Widyarso, *Mobile Computing: The eNetwork Wireless Solution*, ISBN 0738412856, IBM Redbooks, 1999.
- [16] Research in Motion, "Developer's guide for BlackBerry and RIM Wireless Handhelds—Radio API (DataTAC) Version 2.1," 2001.
- [17] John Agosta and Travis Russel, *CDPD: Cellular Digital Packet Data Standards and Technology*, ISBN 0070006008, McGraw-Hill, 1996.
- [18] P. Sinha, N. Venkitaraman, T. Nandagopal, R. Sivakumar, and V. Bharghavan, "A Wireless Transmission Control Protocol for CDPD," presented at IEEE WCNC '99, 1999.
- [19] A. K. Salkintzis, "A survey of mobile data networks," *IEEE Communications Surveys & Tutorials*, Vol. 2(3), pp. 2–18, 1999.
- [20] L. F. Chang, "Wireless Internet—Networking Aspect," in *Wireless Communication Technologies*, New Multimedia Systems, N. Morinaga, R. Kuhno, and S. Sampei, Eds., Kluwer Academic Publishers, 2000, pp. 215–244.
- [21] Behrouz A. Forouzan, *Data Communications and Networking*, 2nd ed. Update, ISBN 0072822945, McGraw-Hill, 2002.
- [22] Alexander J. Huber and Josef F. Huber, *UMTS and Mobile Computing*, ISBN B000089CJ3, Artech House, 2002.

KOSTAS PENTIKOUSIS, PhD, studied computer science at Aristotle University of Thessaloniki and Stony Brook University. He is an ERCIM Fellow at VTT, The Technical Research Center of Finland, and currently resides in Oulu, Finland. For more about his research and publications visit: www.cs.stonybrook.edu/~kostas. The best way to reach him is via skype. E-mail: kostas@cs.sunysb.edu

Opinion: ICANN, the ITU, WSIS, and Internet Governance

by Geoff Huston, APNIC

This is an opinion piece, intended primarily to provoke thought and comment. The author does not claim to personally hold any of the opinions expressed in this article.

It may have taken some three decades to get here, but there is no doubt that the Internet is now a major public communications utility. That is hardly the most important piece of news you are likely to read today, but the implication of this public role is that there are legitimate issues of public policy to consider when looking at the broad topic of coordination of various aspects of Internet infrastructure. In other words, “Internet Governance” is a matter of significant concern to many.

This opinion piece looks at the various range of views about the *Internet Corporation for Assigned Names and Numbers* (ICANN)^[1] and its rationale and role over its brief history. Of course, no look at Internet Governance would be complete without also looking at the role of the *International Telecommunications Union* (ITU), as well as the broader background to this topic. It is a large topic and it has already been the catalyst for numerous articles.

Data Networking and Public Networks

Whether it was because of its antecedents in the research community, or simply because it was not originally envisaged that the Internet would become a global communications platform in its own right, or for whatever reasons, the administration of the Internet infrastructure was not originally crafted with conventional public network coordination in mind. The retrofitting of a model that incorporates considerations of a public utility role is proving to be a rather complicated process.

For example, the original hierarchical name space for the Internet used a set of generic top-level root zone names of “**edu**,” “**net**,” “**com**,” “**gov**,” and “**mil**.” Adding country codes to the root of the name space was a later modification. Even then the original country code delegations were undertaken to individuals or entities who appeared to have some form of link to the national Internet community, rather than specifically seeking out an appropriate office of the national administration of communications services as the point of delegation. Similarly, IP addresses were structured without any form of national prefix, nor were IP addresses distributed along any national lines. In these respects the Internet was really no different from any other computing networking protocols of the 1980s, such as *DECnet*, the *Xerox Network System* (XNS), *AppleTalk*, or IBM’s *Systems Network Architecture* (SNA), where names and addresses were defined in a limited context of the scope of the network, rather than within some broader public name framework.

There were two notable exceptions to this characterization of computer network protocols, and both were designed with a public communications utility as their primary objective, namely X.25 and the *Open Systems Interconnection* (OSI) model. They can be regarded as offerings from the data services sector of the established telephone industry. X.25, the earlier of these two protocols, had a very obvious relationship to telephony, complete with the notion of a “call” as the means of establishing a data connection and as the unit of a transaction. The addressing scheme used a structured space that drew heavily on the telephone number structure. Like telephony, there was no associated name scheme and endpoints were identified by their numeric X.25 protocol address. OSI represented a later effort to design a packet-switched network architecture that was intended to reflect an increasing level of experience with this technology, but nevertheless continued to draw heavily on telephony design. Much was written about OSI at the time, and it would be a diversion to explore it in depth here. However, the salient observation here is that despite the extensive effort invested into its promotion, OSI was a market failure, and whatever its technical merits it was simply not accepted by the communications industry.

OSI was heavily supported by the ITU, and by virtue of this very active sponsorship of this technology, the implication of the aftermath of OSI was that the ITU was seen as being simply out of touch with data networking. It was often portrayed that the ITU was coming from a mindset that was incapable of engaging with either the data communications industry or the broader consumer market for data services. From the perspective of data networking, the failure of OSI was seen as a failure of the ITU itself.

The ITU and the Internet

The ITU is certainly one of the more venerable institutions in the communications sector. It can trace its origins to May 1865, when the first *International Telegraph Convention* was signed by 20 founding national members, and the *International Telegraph Union* was established to facilitate subsequent amendments to this initial agreement. Two decades later, in 1885, the ITU drafted international legislation governing telephony. With the invention in 1896 of wireless telegraphy, similar coordinating measures were adopted by the *International Radiotelegraph Convention*. In 1932 the Union combined the International Telegraph Convention of 1865 and the International Radiotelegraph Convention of 1906 to form the *International Telecommunication Convention*. The name of the body was changed to *International Telecommunication Union* to properly reflect the full scope of the Union’s responsibilities, which by this time covered all forms of wireline and wireless communication.

In 1947 the ITU, under an agreement with the newly created United Nations, became an agency of the United Nations, with responsibilities in international telephony, telegraphy, and radio communications. Over the next four decades the ITU oversaw a system of international interconnection of telephony and data systems that became an industry in and of itself.

The ITU assumed a role of facilitating what was asserted to be a balanced international environment where the costs of running the international system were fairly apportioned between national service providers. In practice these lofty goals were not achieved very efficiently, and international facilities were priced at levels that were considerably higher than the associated costs of actual service provision. When attempts were made to redress the imbalances between large and small national carriers, the outcomes included collective action on the part of the national carriers that operated in ways not dissimilar to a cartel.

In 1992 the ITU was restructured into three sectors, corresponding to its three main areas of activity, namely the standardization of telecommunications technologies in the ITU-T, the coordination of radiocommunications in the ITU-R, and telecommunication development in the ITU-D. In 1994 the ITU established the *World Telecommunication Policy Forum* (WTPF), a group that encouraged the exchange of ideas and information about emerging policy issues arising from the changing telecommunication environment. The first WTPF was held in 1996 on the theme of global mobile personal communications by satellite, and the second in 1998, on trade in telecommunication services.

The ITU was heavily criticized over the ponderous amount of time taken to generate telecommunications standards, the nature of the process used in developing these standards in a closed set of forums, the marginal relevance of these standards, and the final indignity, that the ITU charged for paper and electronic copies of these standards. As some critics pointed out, perhaps harshly, this was not just a case of paperware about vapourware, it was a case of very expensive paperware about vapourware!

More recently, the ITU has focused on attempting to strengthen the participation of the private sector in the work of the Union, as well as streamlining the ITU's processes to reduce the level of delay and amount of process overhead in standardization of technology and operational practices. The ITU has sponsored the establishment of the *World Summit on the Information Society* (WSIS)^[2], and has been attempting to position itself more centrally in the process of further evolution of the Internet as part of its overall charter.

The Internet has posed a severe challenge to the ITU. Not only was the ITU often perceived as being out of touch with the data communications sector, more critically it had been perceived as being incapable of making the necessary reforms to its mode of operation and policy setting to bring it back into relevance for the rapidly changing communications industry. The inference was being drawn that the ITU was apparently in a state of denial over progressive deregulation of national communications sectors. In many cases the national position had already moved to a position of lightweight regulation, relying on strong competitive pressures in the private sector to enforce regimes of efficiency and effectiveness in the supply of communications services to consumers. The ITU, as an intergovernmental organization, was being seen in some quarters as an anachronistic recalcitrant relic of an earlier era of communications service provision.

It was also evident that this critical view of the ITU was most strongly held within the United States, and in particular those parts of the U.S. administration and industry that were involved with the growth of the Internet. It was perhaps no coincidence that in these growth industries of personal computer technologies and the related Internet industry it was U.S. enterprises that were the “poster children” of this new model of industry-led deregulated communications services. Their consequent rapid expansion into a massive global undertaking of the global Internet was perhaps the most eloquent form of statement about the effectiveness of deregulation, and the degree to which the previous regulatory model had simply not managed to encompass the burgeoning demand for data services in a timely fashion.

From this perspective it should be no surprise to observe that when the transition of the *Internet Assigned Numbers Authority* (IANA) function from a fully federally funded research activity to some form of new foundational base was being considered by the U.S. administration, it appears that the ITU was never seriously contemplated as a viable home for this function. If the Internet was a child of deregulation and industry initiative taking on the outcomes of research activity, then the appropriate progression of the IANA function was also from a research context into an enterprise context. IANA should be responsive to industry needs, and to best achieve this the IANA function itself should be undertaken as a task housed within the deregulated private enterprise sector, rather than establishing yet another public bureaucracy, or using existing bureaucracies for the role. ICANN was the embodiment of this aspiration on the part of the U.S. administration, and to pass the effective levers of control of the Internet to the ITU was seen as denying the Internet any form of a productive, innovative, and successful future.

The Formation of ICANN

Whatever the original motivation in creating ICANN to administer the IANA responsibilities, it is now apparent that ICANN was deliberately structured to provide the industry with an alternative structure of coordination and regulation within national and international communications sectors to that of the ITU. The critical difference is that ICANN had not placed governments at the forefront of visible activity, but instead placed industry needs and the operation of a competitive deregulated international communications sector as being the major thrust of coordination activities.

As with any novel model of public policy determination, ICANN’s acceptance ranged from cautious approval to advanced skepticism. Even within the U.S. administration ICANN has yet to be “unleashed,” and it currently operates under the terms of a Cooperative Agreement with the *National Telecommunications and Information Administration* of the U.S. Department of Commerce under a sole source cooperative agreement. In this light ICANN appears to be a cautious step in a bold direction.

ICANN undertakes activities of management of Internet Protocol infrastructure in the areas of the content of the root of the *Domain Name System* (DNS) and the identification of parties to whom are delegated administrative and operational control of the top-level domains and the associated specification of terms and conditions of this delegation. ICANN, through IANA, also manages the pool of unallocated IP addresses (IPv4 and IPv6 addresses and Autonomous System numbers), and also manages the protocol parameter registries as defined by IETF Standards Actions.

ICANN Mki

The initial structure of ICANN had three “supporting organizations,” focusing on:

- Coordination of the DNS with the *Names Supporting Organization* (NSO)
- Coordination of address policies with the *Address Supporting Organization* (ASO)
- Operation of Internet Protocol parameter registries with the assistance of the *Protocol Supporting Organization* (PSO)

The intended role of these supporting organizations was to provide a venue where interested parties could develop and consider policy proposals, leaving the task of ultimate identification of broad support for particular policy initiatives to the ICANN Board.

As has been evident to any observer of the ICANN process, things did not proceed within the parameters of that plan. The NSO met problems due to the diversity of interests that were encompassed with the DNS domain, including emerging national and regional interests in the country code top-level domains, the operators of the generic top-level domains, the trademark and intellectual property collection of interests, the emerging industry of registrars, and a continual interest of individuals who maintained that they had legitimacy of inclusion by virtue of their representation of interests of end users and consumers, or, to use an emerging ICANN lexicon, the “at large” constituency.

The ASO was formed within the parameters of a different model. The *Regional Internet Registries* (RIRs) had already developed a considerable history of working within their communities, and being widely accepted by these communities as an appropriate means of coordination of activity in the role of number resource administration and distribution. The ASO was formed with membership of the associated council based on processes determined by each RIR. Even then it was unclear as to the relationship between the RIRs’ already well-established open policy development process and the ASO and ICANN. The RIRs were unwilling to pass all regionally developed policies to ICANN for a second round of consideration and potential alteration. They insisted that only those policies that were considered to be “global,” in that they were common to all the RIRs, would be passed into this ICANN sphere.

The PSO was placed under strong pressure to include the ITU-T and the *European Telecommunications Standards Institute* (ETSI), and the *World Wide Web Consortium* (W3C) was also enlisted, in addition to the IETF. If the objective of the PSO was oversight and policy formulation concerning the role of protocol parameter registration of IETF protocols, then this enlarged membership of the PSO was unwarranted. Even within the terms of consideration of the PSO as a source of standards-based technical advice to the ICANN Board, the presence of these additional organizations was somewhat puzzling in terms of the match of resultant structure of the PSO to its intended role. The PSO, however, had a role in seating individuals onto the board of ICANN, and it was likely that this aspect of the PSO had been part of the reason for the interest in broader institutional membership. Uncertainty about the extent of the role of ICANN saw many groups attempting to gain access to board seats.

Missing from this mosaic of diverse interests was the inclusion of various national public communications sector entities who also felt that they had clear legitimacy to undertake an active role within the ICANN policy development process, and, in response, the *Government Advisory Committee* (GAC) was formed.

ICANN Evolution and Reform

If a camel is a horse designed by a committee, then it is unclear whether ICANN was a three-humped camel or a three- and three-quarter-humped camel as a result of all this, but camel it undoubtedly was.

The PSO was dysfunctional and missing any tangible agenda of activity. A fracture was apparent in the relationship between ICANN and the IETF. Attempts to create an agreement between ICANN and the IETF over the IANA function were not recognized by the U.S. administration, who continued to insist that, formally, the IANA function for the IETF was undertaken at the behest of the U.S. Department of Commerce rather than the IETF. This view was not shared by the IETF.

The ASO was criticized by ICANN itself of being insufficiently “representative” of the addressing community, and the ICANN Board established its own temporary advisory committee on addresses, and in so doing alienated the RIR community from the entire ICANN framework.

The NSO was hopelessly wedged into factional-based politics.

The GAC decided at the outset that it would operate behind closed doors, in contrast to ICANN’s continuing efforts to operate in an open and transparent manner.

The “At Large” election process undertaken by ICANN appeared to be of dubious validity because of problems in establishing a reliable constituency of individuals who had an interest in ICANN, and a direct election process was attempted only once.

Not surprisingly, ICANN fell into some disarray under these pressures, and by early 2002 the CEO of ICANN at the time, Stuart Lynn^[3], was warning all who cared to listen that ICANN was paralyzed, dysfunctional, and in danger of an imminent demise. Whether this was a message directed to the ICANN Board or to a fractious set of communities that had some intersection with ICANN, or to the U.S. administration who had been influential in determining the original ICANN structure was not entirely clear to any observer of the process.

However, given that ICANN had been set up as an example of a new form of international coordination of communication infrastructure support activities that was based on private-sector activity rather than governmental fiat, this message of imminent failure was widely interpreted both as a potential failure of ICANN and a sign of failure of this new model of coordination of international activity. ICANN was seen as a point of vulnerability with respect to the U.S. administration's diplomatic efforts to reform this international activity sector. The ITU-T's activities in this same area was reinvigorated, with considerable support from national sectors who saw their national interests being potentially advantaged in a ITU-led international environment.

ICANN MkII

Although still firmly positioned as a private-sector activity, and although still making no concessions in the direction of the ITU, ICANN has managed to reorganize its structure through a protracted evolution and reform process.

With respect to the ASO, The Regional Internet Registries formed its own coordination entity, the *Number Resource Organization* (NRO)^[4], and has proposed this entity to ICANN as the means of interfacing between the addressing community and ICANN's policy-development activities.

The PSO was abolished, to be replaced by a *Technical Liaison Group* that, apart from its function of seating an individual on the ICANN Board, is a group without an obvious role or agenda.

The NSO was forced to recognize the fundamental difference between the generic top-level domains, which fall under a more direct relationship with ICANN and its processes, and the country code domains (ccTLDs), which have from the outset been quite wary of ICANN. From the ICANN reform process emerged the *Country Code Name Supporting Organization* (CCNSO) and the *Generic Names Supporting Organization* (GNSO), as a recognition that these two groupings are so dissimilar that they have almost nothing in common.

In addition, an *At Large Advisory Committee* was formed.

The reform process has had some more tangible outcomes, in that formal open meetings of the ICANN Board of Directors have managed to be progressively refined from efforts at direct dialogue and open debate into highly structured events with many formalisms and appropriate quantities of ceremony.

ICANN Today

Despite the effort to encompass coordination activities in the areas of names, addresses, and protocol parameters, ICANN has been largely captured by the names industry, and ICANN's agenda, activity focus, and outcomes are concentrated mostly in the name domain.

In this activity domain, the track record of ICANN is very mixed. To its credit, it has managed to dismantle the most objectionable parts of the monopoly hold over the *generic Top-Level Domains* (gTLDs), create an operational model that makes a clear distinction between registry operators and registrars, impose price and business controls on the registry operation as a means of controlling the natural tendency for the registry operation to reflect its unique position in the form of monopoly rentals, and assist in the creation of a global network of competitive enterprises, with the expectation that competition will instill operational and price efficiency in the registrar business.

In addition, ICANN has been successful in not only introducing new gTLDs to compete with the established brands of **.com**, **.net**, and **.org**, but also in moving **.org** and **.net** to new registry operations (**.net** is under way at the time of writing of this article). Despite these positive achievements, it is not clear that this new regime has been entirely successful.

True competition in the name space is still some way off, and the recently introduced gTLD brands have failed to gain any leverage within the market. The name market itself remains one where the role of name speculators continues to play a significant role in terms of proportion of registered names. The overarching dominance of **.com** as a brand has continued, and the advantaged position of the U.S.-based registrar of this zone continues.

The obscure nature of the relationships between the IETF, ICANN, and the U.S. administration over the protocol parameter registries remains unresolved. The IETF is clearly not in control of its own protocol parameters, and has abrogated this role to ICANN. Standards making entirely divorced from any effective engagement with deployment tends to result in a standards body of dubious long-term validity, and despite its impressive track record in the past, the IETF is clearly already well-distanced from current technology directions in the industry—and the gap continues to widen.

The DNS *Root Server Operators* continue to operate as an independent group. The recent moves to dramatically increase the number of DNS root servers and improve the overall robustness of DNS resolution through anycasting root servers and distributing anycast instances across the globe has been a well-received initiative. The fact this has occurred without any form of ICANN involvement is an interesting commentary on the ability of ICANN to engage with the operational parts of the infrastructure of the Internet. Comparable activities to improve the DNS in terms of resolution services within the ICANN sphere have become protracted exercises that impose a very heavy burden on the patience of the players.

The moves to introduce IPv6 AAAA records into the DNS root have been anticipated for many years, and the response to the recent ICANN announcement is, in general, of the tenor “why didn’t this happen some years ago?” The continuing frustration to get the DNS root to include *Secure DNS* (DNSSEC)^[5] important information continues to illustrate a perspective that the ICANN process appears to be unresponsive to technical needs and end-user imperatives.

The situation today is that ICANN appears to enjoy a mixed level of success. It has managed to establish itself as a means of administering the infrastructure elements of the Internet Protocol in a manner that is reflective of the deregulated nature of the Internet industry. It has managed to reform parts of the landscape and generate an industry structure that uses open competition as the major control mechanism. ICANN has managed to bring much of the discussion about the administration of Internet infrastructure out into the open. All these are major milestones, and it is to the credit of many dedicated individuals that ICANN has managed these impressive outcomes. However, it has been able to achieve all this with the continued sponsorship of the U.S. administration, and the question of whether it can firmly establish itself in its own right in the coming years remains today perhaps a matter of hope rather than absolute certainty.

There are still the lingering concerns that if ICANN, as a private-sector entity, were to once more explore positioning itself on the brink of imminent demise, the collective task of picking up the pieces and continuing to support the operation of the Internet is one that appears to have a very uncomfortable level of uncertainty. In addition, the perception of ICANN as an entity whose single purpose is to maintain an entrenched advantaged position of the United States and of U.S.-based enterprises in the global Internet has been widely promulgated. It is often portrayed that ICANN offers no viable mechanisms for other national or regional interests at a governmental level to alter this somewhat disturbing picture of international imbalance. Although other aspects of international activity fall under various political or trading frameworks, and national and regional interests and positions can be collectively considered and negotiated, critics of ICANN point out that the message ICANN sends to the rest of the world is that the United States is withholding the Internet from conventional international governance processes. Skeptical commentators interpret the U.S. administration’s use of ICANN as at best a delaying technique to gain time to further strengthen the position of U.S.-based enterprises across a lucrative global Internet market, aided and abetted by a compliant industry body that masquerades as an international standards organization.

Such a critical perspective also points to ICANN’s tenuous lines of authority, its lack of performance in many aspects of the domain name enterprise, its seeming obsession with the registrar sector to the apparent exclusion of any other activity, its burgeoning costs, and its lack of acceptance, particularly as it relates to the acceptance of ICANN by the various country code DNS administrators, to name but a few factors.

Accompanying this strident criticism is the line of argument that the Internet does not actually represent a viable challenge to existing mechanisms for coordination of international activity. At both a national and international level, the Internet should not require novel and untested regulatory mechanisms as a means of expressing public interest and public policies. The line of argument from this perspective is that there is neither the demonstrated need, nor any appropriate level of international support at a governmental level to sustain the argument that a private-sector, nonprofit corporation is the best, or even the only viable model of coordination of Internet activity. If “Internet Governance” is the question, then, the line of argument goes, the model upon which ICANN is based is definitely not the best answer we can devise. This very critical line of reasoning has become particularly prominent in the WSIS process, and lies behind much of the continual fascination of the topic of “Internet Governance” in WSIS meetings.

WSIS and Internet Governance

The WSIS has been a long time coming, and it represents a move on the part of the ITU to formulate a revised role for the ITU to engage with a world richly populated by all manner of information services layered upon a highly diverse and capable communications environment. This summit was planned in two phases. The first summit was held in Geneva December 10–12, 2003, where the foundations were laid by reaching agreement on a *Declaration of Principles* and a *Plan of Action*. The second phase will be held in Tunis, November 16–18, 2005, to implement the agenda leading up to achievable targets by 2015, and to agree on unfinished business, most importantly on the question of Internet governance and of financing mechanisms.

Irrespective of any particular political perspective here, the universal observation is that the Internet has heralded a revolutionary change to the global communications enterprise. Markets for communications services are changing, the technology base is changing, the economic models of communication are changing, and the models of interaction at the provider level are changing. The challenge from the public-policy perspective at a world level is to create a framework that ensures that the benefits of this change, in both social and economic terms, are accessible to all, rather than to a subset of the world’s population. It is within this broad framework that WSIS has been positioned.

These are lofty and ambitious goals, and the task before WSIS is certainly as challenging as any in this environment. The hope is that the myriad of participants in this process includes sufficient resources to engage in the agenda in a meaningful way.

However, the underlying issue is that of the progressive change in the role of communications infrastructure from a predominately public-sector activity to a very diverse spectrum of public- and private-sector activity. We appear to have become increasingly reliant on private-sector investment and private enterprise to support the public communications enterprise. But is this necessarily the appropriate model for the entire world, or even any part of the world?

As many recently privatized industries could attest, private-sector activity has entirely different investment motivations and entirely different service objectives. If the nature of the activity is one that requires long-term investment in infrastructure with low returns, then private-sector activity tends to use the existing infrastructure base without necessarily making adequate longer-term replenishment investments. Private activity also tends to concentrate service delivery to the most lucrative sectors of the market, and, if possible, will deliberately avoid establishing services in areas that are less financially attractive. The task of structural cross-subsidization that makes ubiquitous equity of access possible is not seen as a private enterprise outcome, and aspects of communications such as universal service obligations and equity of access are seen as public regulatory functions rather than natural market outcomes of a deregulated industry.

The Internet today is anything but a level and balanced environment. There are concentrations of investment capability, concentrations of technical knowledge and logistical capability, concentrations of intellectual wealth, and concentrations of power and influence. How to create from this current diverse environment some form of structural cross-subsidization that extends the basic means of access to all is the appropriately lofty goal of the WSIS endeavor. There is also the more focused investigation of "Internet Governance" and the agenda of establishing to what extent the perception of the advantaged position of a small number of national entities in all this can be balanced by measures that allow other national economies to invest in this space on terms and conditions that do not involve a continuing flow of money and a ceding of power to these existing advantaged national interests.

As the WSIS documentation points out, "... building the foundations for an Information Society is a complex task. The digital revolution is already impacting the world in deeply intrinsic ways, perhaps more profoundly than even the industrial revolution itself. Yet, while the digital revolution has extended the frontiers of the global village, the vast majority of the world remains unhooked from this unfolding phenomenon."

The Secretary General of the UN chartered a smaller group to examine Internet Governance, in particular, the *Working Group on Internet Governance*, or WGIG. Its nine-month brief is to glean these issues of public policy in an environment that has very significant private-sector interest. Indeed from an international perspective, where regulatory powers, even of a reserve nature, are in a very real sense ephemeral, the work in WGIG to date with its discussion papers has done little. The discussion papers have illustrated the broad nature of the topics raised in the context of Internet Governance, but their poor depth, visibly poor levels of research, and lack of any real analysis of the selected topics only highlights the complexity of the underlying interplay of public- and private-sector interests within a domain that is also bounded by technical considerations.

At the same time the poor quality of these reports highlights the inability of WGIG to engage directly into the heart of this exercise, given their obvious constraints of time and resources. It is not surprising to observe that, following its February meeting WGIG has decided to abandon this set of discussion papers. If a fresh start is being contemplated for WGIG, then perhaps it is time to note that only half of the group's allocated time remains, and the topic is getting no easier with the passing of the days.

For those interests who wanted the ITU to become engaged in the Internet, hope has now been passed to the WSIS process and the related WGIG study into Internet Governance issues. This is seen as being a means of opening up the control of the Internet into a more conventional international process that dismantles what they see as the current position of global taxation that U.S. national interests have imposed on the rest of the world's population in the adoption of Internet-based services. For those who think the ITU remains an unreformed vehicle for the imposition of anachronistic, inappropriate regulatory measures that stultify any form of innovation and progress in telecommunications, the WSIS process is yet another venue to parade the stark contrast between the rather impressive track record of a deregulated market-driven approach to coordination of telecommunications services, as seen with the Internet, and the ineffectual outcomes from the international public regulatory sector.

Looking Forward

One view of this process is that this is a negotiation of national roles of influence and power over the coming century or more, and that this process requires some considerable care and attention at an international level.

This topic is one that places a model of deregulated private sector-led activity, with its market-based disciplines, into direct contrast with a more traditional model of the balancing of various national interests through common regulatory measures undertaken within each national regime as a regulated public-sector process. The proponents of a deregulated approach argue that the Internet is a child of the progressive position of deregulation of communications markets in many national environments, and it is the dynamic and creative impetus of highly competitive markets that has led to the rapid spread of the Internet and the consequent improvements in the efficiency and effectiveness of national and international communications systems. None of these outcomes would have been achievable, they argue, in a regulated regime where innovation and competition for the consumer were completely stifled by the deadening weight of regressive regulation.

Like many bold innovative experiments in international coordination and the establishment of new world orders, ICANN stands a strong risk of falling foul of an inherent conservatism in international politics, where the careful balancing of national interests is seen as being far more critical an objective than any actual outcomes that may be achieved from the process.

From this perspective, ICANN is critically reliant on its acceptance by all players of its legitimacy to operate in this space, and also critically reliant on acceptance of the proposition that these issues are best addressed in open forums of debate. This task is difficult, and the limited set of outcomes that ICANN can point to as being products of this process do not install a high degree of confidence that this process is stable, scalable, well-founded, and sustaining. Currently the proposition is not that ICANN represents the most appropriate enduring framework here, but that the track record of the alternative has failed in the past and nothing has changed to prevent the historical alternative framework making similar flawed decisions in the future.

The opposite end of the spectrum of views argues that nothing has really changed with the introduction of the Internet, and the international regime remains one where various national interests need to be resolved in a coordinated and equitable fashion. Without some form of common regulatory constraint, there are inevitable market distortions where the expression of vigorous national aspirations results in an advantaged position in the international domain. Public communications is a public-sector activity, they argue, and, ultimately, the only points of control rest within national regulatory regimes, and internationally it is a case where national interests must be balanced through a process that recognizes political realities of coordination and compromise. From this perspective it is asserted that the ITU is the intergovernmental venue for this activity as it relates to the communications sector, and it is to the ITU that national interests must look to redress distortions where one national entity or one region holds a contrived privileged position with respect to international communications.

In looking at these two extremes of perspective, an obvious question is what then is the role of international public policy setting? In this form of market-mediated service supply functions, are international issues being progressively transformed into aspects of international trade? Does such an environment provide adequate protection for developing economies? Are common social priorities being adequately considered in such a framework?

This leads to a more basic question of whether the existing international institutions, such as the ITU, are appropriately positioned to meet these public policy challenges, or should we be considering changes here in order to bring the international institutional framework into better alignment with the emerging information society?

These are certainly difficult positions to attempt to reconcile, and perhaps it is being impatient to expect clear outcomes in the near future, and certainly very difficult to expect that in a few short months WGIG and WSIS will be able to deliver a balanced, considered, and generally acceptable outcome in this space. It is also a natural concern in looking at these rather aggressive schedules for WSIS that short-term political expediency will obstruct genuine attempts to truly understand the fundamental nature of the changes that are happening with the differing model of communications that are heralded by the Internet model.

References

- [1] Internet Corporation for Assigned Names and Numbers (ICANN): <http://www.icann.org>
- [2] The World Summit on the Information Society (WSIS): <http://www.itu.int/wsis/>
- [3] M. Stuart Lynn, "A Unique, Authoritative Root for the DNS," *The Internet Protocol Journal*, Volume 4, No. 3, September 2001.
- [4] Number Resource Organization (NRO): <http://www.nro.net/>
- [5] Miek Gieben, "DNSSEC: The Protocol, Deployment, and a Bit of Development," *The Internet Protocol Journal*, Volume 7, No. 2, June 2004.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector, and has served his time with Telstra, where he was the Chief Scientist in the company's Internet area. Geoff is currently the Internet Research Scientist at the Asia Pacific Network Information Centre (APNIC). He is also the Executive Director of the Internet Architecture Board, and is a member of the Board of the Public Interest Registry. He is author of *The ISP Survival Guide*, ISBN 0-471-31499-4, *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, and co-author of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. All three books are published by John Wiley & Sons. E-mail: gih@apnic.net

Book Review

Unix Network Programming *Unix Network Programming, 3rd Edition*, by W. Richard Stevens, Bill Fenner, Andrew M. Rudoff, ISBN 0131411551, Addison-Wesley Professional, 2003.

It would be difficult to put value on a book that has been a classic text and a reference in academia and in the real world in the context of network programming for over a decade. Richard Stevens published the ever-popular *Unix Network Programming* [UNP] back in 1990, and the second edition followed in 1998. With a dedication to the memory of R. Stevens, the UNP book found itself two new authors, Bill Fenner and Andrew M. Rudoff, who would write the third edition of this book. The third edition has many updates, a new look and feel and many of new chapters that cover the topics more applicable these days. In my opinion, it is still the most valuable and profound text in the context of network programming.

Changes and Updates

For those of us who have the first two editions of this book, the third edition has the following changes:

- IPv6 updates. In the second version of the book, IPv6 was merely a draft, and the sections covering IPv6 have been updated to reflect these changes.
- POSIX updates. The functions/APIs and examples have been updated to reflect the changes to the latest version of the POSIX specification (1003.1-2001).
- SCTP coverage. Three new chapters that cover this new reliable, message-based transport protocol have been added.
- Key Management Sockets coverage. Network security and its applicability and use with IPsec are covered.
- The Operating Systems and machines that are used for the examples have been updated.
- Some topics such as Transaction TCP and X/Open Transport Interface have been dropped.

Many topics and sections have been updated with the authors' comments. These comments even though simple for someone new to the profession, are extremely useful because they are like hints and tips from one developer to the next to help you in your next programming assignment.

Unix Focus

If this is the only edition of the book that you will read, you are in for a treat. Topics in Network Programming are covered in detail, using concrete programming examples that all of us can relate to—all Unix, but what else is there?!

All kidding aside, the topics are covered well enough that they are useful information under any operating system. The concepts don't change; sockets are sockets under any operating system. The function call is different, but one needs to go through the same steps under any environment.

Being the most popular networking protocol, TCP/IP is covered in Part I of the book. You need to have prior understanding of the TCP/IP protocol and the OSI model, however. If this is the first time you are looking at the programming aspects of networking protocols, Part I of this book covers the basics. It begins with a couple of simple examples such as such as daytime client and a daytime server and it builds on that. TCP, UDP, and SCTP (*Stream Control Transmission Protocol*) are covered in brief in Part I, and basic concepts such as the three-way handshake of TCP and the four-way handshake of SCTP are depicted.

Part II of the book covers *sockets* and socket programming. Topics such as the socket Address Structure in IPv4 and IPv6 for TCP, UDP and SCTP are covered and examples (the same daytime client/server) are given to convey the point. It is important to mention here that all the topics and concepts are depicted for the three transport protocols: TCP, UDP and SCTP. Every socket API under the Unix programming environment is covered and examples are given for each function call to show the reader how the function can be utilized. Much attention is dedicated to Socket Options and how they are used or can be used for best results. Hints are given throughout the chapter about the pitfalls and best practices of each option.

After the basics are been covered, various I/O models are depicted in detail and examples are shown to convey the advantages and disadvantages of each I/O model. The five I/O models used through the book (and available under the Unix environment) follow:

- Blocking I/O
- Non-blocking I/O
- I/O Multiplexing (using select and poll)
- Signal driven I/O
- Asynchronous I/O

The *Stream Control Transmission Protocol* (SCTP), a new IETF standard is also covered in detail—from the basics to the advanced. The two interface models of SCTP (one-to-one and one-to-many) are covered in detail, and their differences with TCP are also explained in full. The client/server example used throughout the book is ported to use the new SCTP protocol. The authors then explain in detail the problems that SCTP solves over TCP and where and how it would be useful to use SCTP.

Advanced topics such as IPv4 and IPv6 portability, Unix Domain Protocols, Multicasting and advanced Socket programming for UDP, TCP and SCTP cover the rest of the chapters in this book.

Various options for interoperability between IPv4 and IPv6 are discussed in the last section of the book. Advanced I/O functions bring us a new perspective of how complicated Network Programming can become. Benefits and examples of nonblocking I/O are covered in detail—the authors give examples to show us how, with very few modifications, the performance of a socket application can improve dramatically. Various methods on how to control socket operations are discussed including the use of an alarm along with SIGALRM, the use of select and various timeout options that are available in the API.

The chapters that discuss Multicasting and adding reliability to UDP are my favorite chapters in this book. The Time Server used throughout the book is re-coded to become a multicast application. Some issues that arise when designing multicast applications such as multicast on a WAN are also discussed.

As Good as Ever

The third edition of *Unix Network Programming* is as good as ever. The updates truly reflect solutions to today's challenges in network programming. Bill Fenner and Andrew Rudoff did an amazing job continuing the work of a true legend in the field of Computer Science.

—Art Sedighi
asedighi@tibco.com

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don't have access to it. Contact us at **ipj@cisco.com** for more information.

Internet Pioneers Cerf and Kahn to Receive ACM Turing Award

The *Association for Computing Machinery* (ACM), has named Vinton G. Cerf and Robert E. Kahn the winners of the 2004 *A.M. Turing Award*, considered the “Nobel Prize of Computing,” for pioneering work on the design and implementation of the Internet’s basic communications protocols. The Turing Award, first awarded in 1966, carries a \$100,000 prize, with financial support provided by Intel Corporation. Cerf and Kahn developed TCP/IP, a format and procedure for transmitting data that enables computers in diverse environments to communicate with each other. This computer networking protocol, widely used in information technology for a variety of applications, allows networks to be joined into a network of networks now known as the Internet.

ACM President David Patterson said the collaboration of Cerf and Kahn in defining the Internet architecture and its associated protocols represents a cornerstone of the information technology field. “Their work has enabled the many rapid and accessible applications on the Internet that we rely on today, including e-mail, the World Wide Web, Instant Messaging, Peer-to-Peer transfers, and a wide range of collaboration and conferencing tools. These developments have helped make IT a critical component across the industrial world,” he said.

“The Turing Award is widely acknowledged as our industry’s highest recognition of the scientists and engineers whose innovations have fueled the digital revolution,” said Intel’s David Tennenhouse, Vice President in the Corporate Technology Group and Director of Research. “This award also serves to encourage the next generation of technology pioneers to deliver the ideas and inventions that will continue to drive our industry forward. As part of its long-standing support for innovation and incubation, Intel is proud to sponsor this year’s Turing Award. As a fellow DARPA alumnus, I am especially pleased to congratulate this year’s winners, who are outstanding role models, mentors and research collaborators to myself and many others within the network research community.”

In 1973, Cerf joined Kahn in a *Defense Advanced Research Projects Agency* (ARPA, now called DARPA) project to link three independent networks into an integrated “network of networks.” They sought to develop an open-architecture network model for heterogeneous networks to communicate with each other independent of individual hardware and software configuration, with sufficient flexibility and end-to-end reliability to overcome transmission failures and disparity among the participating networks. Their collaboration led to the realization that a “gateway” (now known as a *router*) was needed between each network to accommodate different interfaces and route packets of data. This meant designating host computers on a global Internet, for which they introduced the notion of an *Internet Protocol* (IP) address.

As a graduate student at the University of California at Los Angeles, Cerf had contributed to a host-to-host protocol for ARPA's fledgling packet-switching network known as ARPANET. Kahn, prior to his arrival at ARPA, led the architectural development of the ARPANET packet switches while at Bolt Beranek and Newman (BBN), and had showcased the ARPANET in 1972, at the first International Conference on Computer Communications. ARPANET had already connected some 40 different computers and demonstrated the world's first networked e-mail application.

In May 1974, they published a paper describing a new method of communication called *Transmission Control Protocol* (TCP) to route messages or packets of data. Like an envelope containing a letter, TCP broke serial streams of information into pieces, enclosed these pieces in envelopes called "datagrams" marked with standardized "to and from" addresses, and passed them through the underlying network to deliver them to host computers. Only the host computers would "open" the envelope and read the contents.

This networking arrangement allowed for a three-way "handshake" that introduced distant and different computers to each other and confirmed their readiness to communicate in a virtual space. In 1978, Cerf and several colleagues split the original protocol into two parts, with TCP responsible for controlling and tracking the flow of data packets ("letters"), and IP responsible for addressing and forwarding individual packets ("envelopes"). The new protocol, TCP/IP, has since become the standard for all Internet communications.

Vinton Cerf and Robert Kahn share a number of awards, including the 1991 ACM Software System Award, the 2001 Charles Stark Draper Prize from the National Academy of Engineering, the 2002 Prince of Asturias Award, and the 1997 National Medal of Technology from President Bill Clinton. They are both the recipients of numerous honorary degrees. ACM will present the Turing Award at the annual ACM Awards Banquet on June 11, 2005, in San Francisco, CA.

The A.M. Turing Award was named for Alan M. Turing, the British mathematician who articulated the mathematical foundation and limits of computing, and who was a key contributor to the Allied cryptanalysis of the German Enigma cipher during World War II. Since its inception, the Turing Award has honored the computer scientists and engineers who created the systems and underlying theoretical foundations that have propelled the information technology industry.

For additional information see:

<http://www.acm.org/awards/taward.html>

New Administrative Structure for the IETF

The *Internet Engineering Task Force* (IETF) is well advanced in the process of making a significant change to the administrative structure that supports the world's leading Internet standards development group. The creation of an *IETF Administrative Support Activity* (IASA) is an important move designed to help the IETF maintain and expand the unique open processes that have enabled the development of Internet standards since 1986.

The new structure will allow the IETF to take full responsibility for managing the resources required to accomplish its work—giving the IETF a solid foundation on which future operations will be based.

This is the first time that all the IETF's administrative and support functions will be managed directly by the IETF as one fully integrated entity. Until now, administration of the IETF has been carried out exclusively by helper organizations and volunteers. The new IASA will be formally structured as an activity within the *Internet Society* (ISOC)—the organizational home of the IETF—and an *IASA Administrative Director* (IAD) will be appointed to provide central management of IETF administration.

The decision to move forward with the new structure was taken after extensive consultations with the Internet community. A number of key prerequisites for efficient administrative operations were identified, including the need for the IETF to have budgetary autonomy. The IETF is currently supported by funding from multiple sources, including meeting fees, donations from interested corporate and non-corporate entities, and donations in kind of equipment or manpower. The IASA will allow the IETF to be able to consider all sources of income, and all expenses involved in running the IETF, as pieces of one budget.

The IASA will also be responsible for defining clear contractual relationships with other organizations that will continue to provide basic services, including meeting organization, secretarial services, IT services, etc. The new structure also gives the IETF flexibility in how it chooses to fund and develop any additional services that may be required.

The IETF is a large open international community of network designers, operators, vendors, and researchers concerned with the evolution of the Internet architecture and the smooth operation of the Internet. It is open to any interested individual. See: <http://www.ietf.org>

ISOC is a non-governmental international organization for global cooperation and coordination for the Internet and its internetworking technologies and applications. Members comprise commercial companies, governmental agencies, foundations, and individuals. ISOC has 82 Chapters in over 60 countries around the world. For more information see: <http://www.isoc.org>

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Technology Strategy
MCI, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
VeriFi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 2005 Cisco Systems Inc.
All rights reserved.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PPSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA
--

The Internet Protocol Journal

June 2005

Volume 8, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
IPv6 and MPLS.....	2
Graph on Path	13
Book Reviews	22
Fragments	26
Call for Papers	31

FROM THE EDITOR

The Internet is a constantly evolving environment which puts pressures on existing and evolving protocols. Any protocol changes must be carefully designed and even more carefully deployed to avoid any disruption to the running system. It is no longer possible to orchestrate a simple overnight switch, so engineers are considering various transition and evolution strategies. In this issue we bring you two examples of this kind of evolutionary protocol development.

Our first example relates to *IP Version 6* (IPv6). A great deal of effort is going into the deployment of IPv6, and good transition strategies can help. Tejas Suthar explains how *Multiprotocol Label Switching* (MPLS) can be used for a transition from IPv4 to IPv6.

Our second example looks at a possible enhancement to the *Border Gateway Protocol* (BGP). BGP in its current form is already nearly ten years old, and calls for its replacement can be heard from network operators. Russ White discusses some possible changes that would not require a wholesale protocol replacement.

It is not every day that a book on punctuation becomes an international best seller, and it is certainly not common for IPJ to review such a non-computer related book. But I think it is appropriate for several reasons. First, accurate punctuation is important not just for computer parsers, it is important for all professionals whether we are sending quick e-mails or writing project reports. Second, this is a really *fun* as well as informative book. And last, but not least, it gives me an opportunity to introduce you to Bonnie Hupton, who provides copy-editing services for this journal. Without her help, IPJ would be far less readable.

Our Website at www.cisco.com/ipj has a new look, but still contains links to our back issues, index files and the IPJ subscription system. Please take a moment to renew or update your subscription. If you have questions or comments, please send them to ipj@cisco.com.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

IPv6—A Service Provider View in Advancing MPLS Networks

by Tejas Suthar, TELUS Communications Inc.

We are all aware of the evolution of the *Internet Protocol* (IP) and its dominance on all aspects of our lives, either directly or indirectly. Currently IP Version 4 delivers critical business application traffic in a so-called new world of the Internet. As the evolution goes on, *IP Version 6* (IPv6)^[5] is becoming a necessary element of the network. IPv6 will enable businesses to expand their capabilities exponentially without having any limitations or restrictions. As technologies evolve and the adoption of IP-enabled devices accelerates, IP will enter a new era as the protocol of choice for communications. Using globally unique IPv6 addresses increases the opportunity for service providers to create new business models and add revenue, and it increases the portfolio of services. However, the major demand for support of IPv6 will be mobile applications; the IT world will also tie in all the systems for transparent operation. The days are not far when permanent IPv6 addresses will be assigned to individuals for their communication purposes—either *Voice over IP* (VoIP), video over IP, video on demand, wireless Internet access, unified messaging, etc. Also, IP smart appliances are becoming more and more popular, and the result will be explosive usage and adoption of IPv6 addresses. Articles outlining the importance of IPv6 and limitations of IPv4 abound. This article is mainly geared toward highlighting the service provider networks that are built or currently being built to support IPv6 in a VPN fashion.

Multiprotocol Label Switching (MPLS)^[4] is widely accepted as a core technology for the Next-Generation Internet that provides speed and functions in packet forwarding. Service providers that offer MPLS/VPN services to their customers are looking forward to adding IPv6 VPN services to their portfolio. Service providers that want to support IPv6 in traditional ways have few options, such as tunneling methods (for example, manual, *Tunnel Broker*, *Generic Routing Encapsulation* [GRE], or *Intrasite Automatic Tunnel Addressing Protocol* [ISATAP], which has scalability problems); or Native IPv6 with dual-stacked MPLS core. However, consider the following:

- For MPLS VPN services, service providers made a significant investment in building the IPv4/MPLS backbone. The return on investment thresholds are probably yet to be achieved.
- Backbone stability is another critical factor; service providers must offer reliable services, especially with regard to voice over MPLS. Most service providers have recently managed to stabilize their IPv4 infrastructure, and they are hesitant to make another significant move when it comes to supporting IPv6 unless the integration is smooth.

Standards bodies with help from vendors and leading service providers are addressing these concerns. Currently service providers have two approaches that they can deploy to support IPv6 without making any changes to the current IP (v4) MPLS backbones, namely 6PE^[1] and 6VPE^[2], originally defined in RFC 2547.

The 6PE approach lets IPv6 domains communicate with each other over an IPv4 cloud without explicit tunnel setup, requiring only one IPv4 address per IPv6 domain. The 6PE technique allows service providers to provide global IPv6 reachability over IPv4 MPLS. It allows one shared routing table for all other devices. Typical applications are IP toll voice traffic and Internet transit services over a common MPLS infrastructure. The 6PE technique does not provide any logical separation because it is for MPLS VPN.

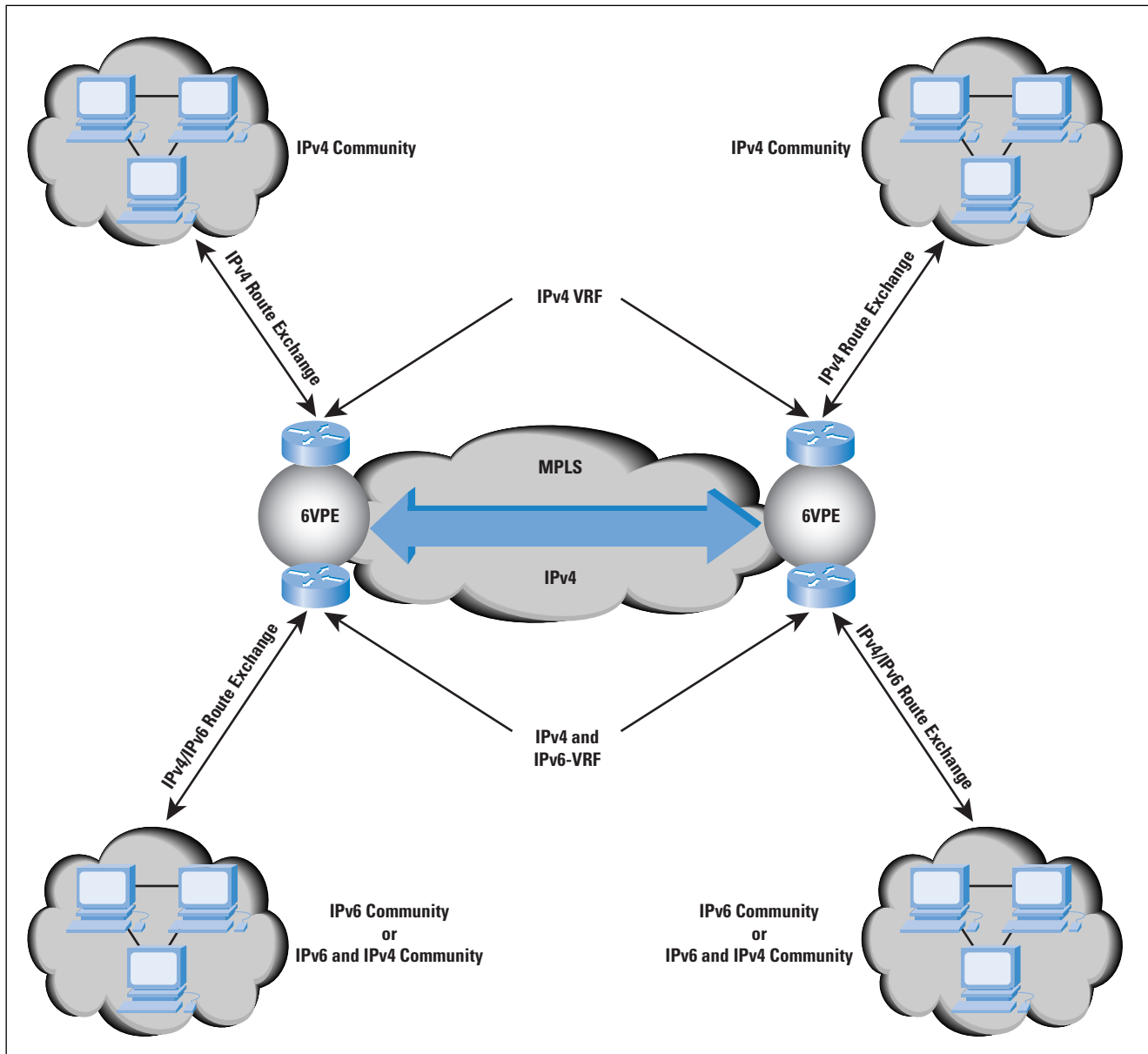
The newest feature to facilitate the RFC 2547bis-like VPN model for IPv6 networks is called 6VPE. It will save service providers from enabling a separate signaling plane, and it takes advantage of operational IPv4 MPLS backbones. Thus there is no need for dual-stacking within the MPLS core. This represents a huge cost savings from the operating expenses perspective and addresses the security limitations of the 6PE approach. 6VPE is more like a regular IPv4 MPLS-VPN provider edge, with an addition of IPv6 support within *Virtual Routing and Forwarding* (VRF). It provides logically separate routing table entries for VPN member devices. This article reviews this approach in more detail because it is the likely approach to succeed in the service provider network.

Under the Hood of 6VPE

Before we look into the 6VPE, it is important to clarify the definition of “dual stack,” a technique that allows IPv4 and IPv6 to coexist on the same interfaces. Today, IPv4 has roots in most of the hosts that run applications. Moreover, stability as well as reliability of new applications over IPv6 is maturing. Therefore, coexistence of IPv4 and IPv6 is a requirement for initial deployment. With regard to supporting IPv6 on a MPLS network, two important aspects of the network should be examined:

- *Core*: The 6VPE technique allows carrying IPv6 in a VPN fashion over a non-IPv6-aware MPLS core. It also allows IPv4 or IPv6 communities to communicate with each other over an IPv4 MPLS backbone without modifying the core infrastructure. By avoiding dual-stacking on the core routers, the resources can be dedicated to their primary function to avoid any complexity on the operational side. The transition and integration with respect to the current state of networks is also transparent.
- *Access*: In order to support native IPv6, the access that connects to IPv4/IPv6 domains need to be IPv6-aware. Service provider edge elements (provider edge routers) can exchange routing information with end users. Hence dual stacking is a mandatory requirement on the access layer as shown in Figure 1.

Figure 1: 6VPE Overview



The IPv6 VPN solution defined in this article offers many benefits. Especially where a coexistence of IPv4 and IPv6 is concerned, the same MPLS infrastructure can be used without putting additional stress on the provider router. Also the same set of *Multiprotocol Border Gateway Protocol* (MPBGP) peering relationships can be used. Because it is independent of whether the core runs IPv4 or IPv6, the IPv6 VPN service supported before and after a migration of the core to IPv6 can be done independent of the customer VPN.

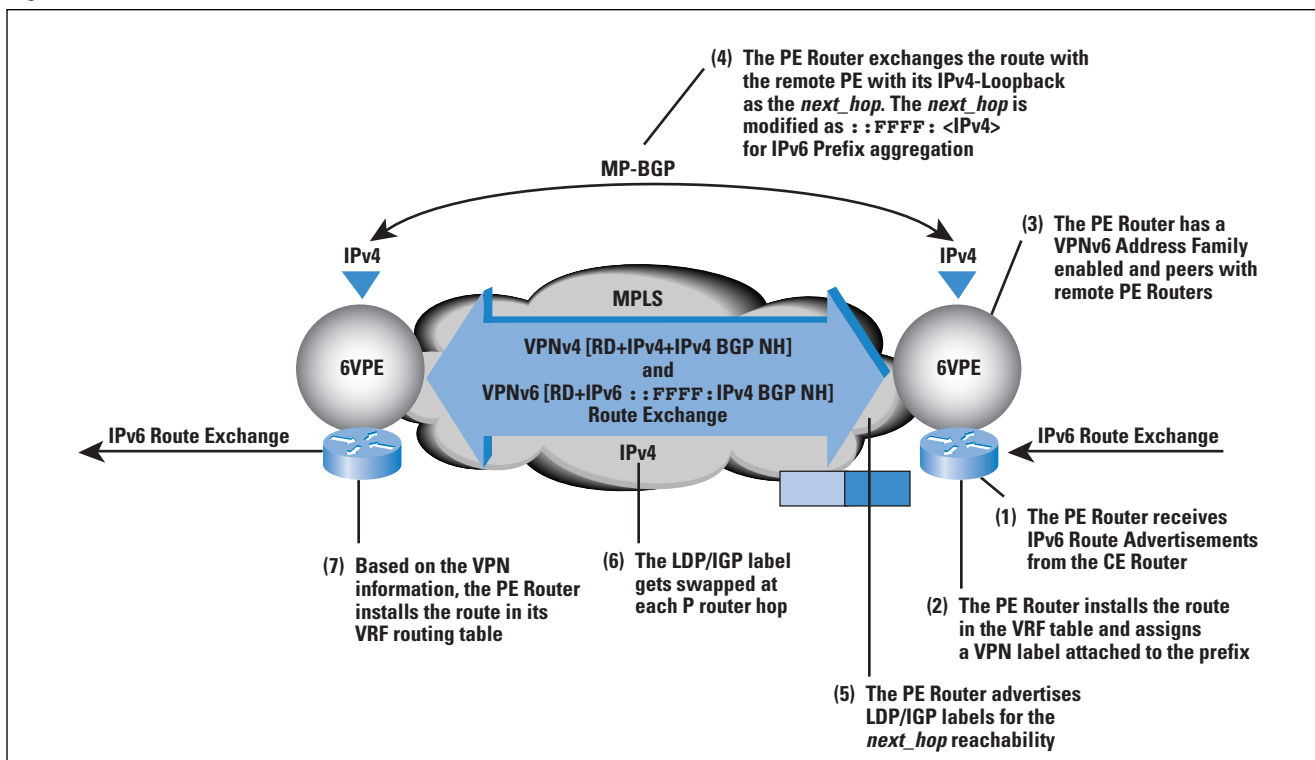
Within the MPLS core, the backbone *Interior Gateway Protocol* (IGP) (*Intermediate System-to-Intermediate System* [IS-IS] or *Open Shortest Path First* [OSPF]) populates the global routing table (v4) with all provider edge and provider routes. As outlined in the draft for IPv4 MPLS VPN (2547-bis), 6VPE routers maintain separate routing tables for logical separation. This allows the VPN to be private over a public infrastructure.

The VRF table associated with one or more directly connected sites (customer edge devices) form close IPv6 or IPv4 speaking communities. The VRFs are associated to physical or logical interfaces. Interfaces can share the same VRF if the connected sites share the same routing information. MPLS nodes forward packets based on the top label. IPv6 packets and IPv4 packets share the same common set of forwarding characteristics or attributes, also known as *Forwarding Equivalence Class* (FEC) within the MPLS core.

6VPE Operation

When IPv6 is enabled on the sub-interface that is participating in a VPN, it becomes an IPv6 VPN. The customer edge-provider edge link is running IPv6 or IPv4 natively. The addition of IPv6 on a provider edge router turns the provider edge into 6VPE, thereby enabling service providers to support IPv6 over the MPLS network.

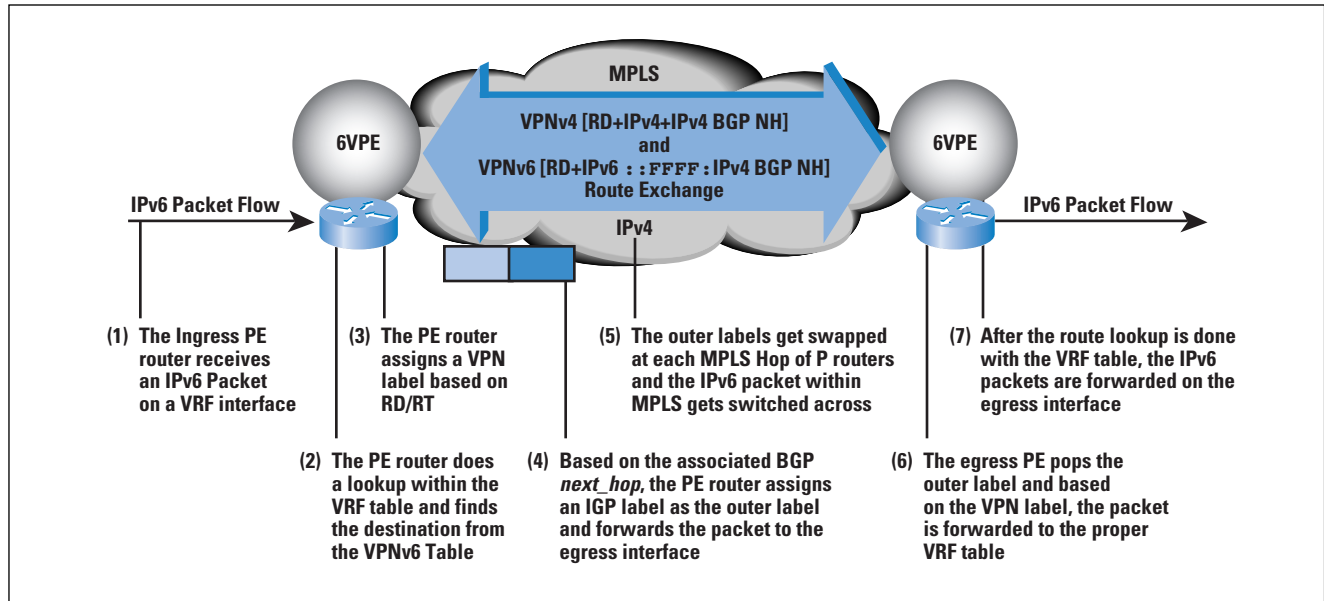
Figure 2: 6VPE Route Advertisement



As outlined in Figure 2, provider edge routers use VRF tables to maintain the segregated reachability and forwarding information of each IPv6 VPN. MPBGP with its IPv6 extensions distributes the routes from 6VPE to other 6VPEs through a direct *internal BGP* (iBGP) session or through VPNv6 route reflectors. The next hop of the advertising provider edge router still remains the IPv4 address (normally it is a loopback interface), but with the addition of IPv6, a value of `::FFFF:` gets prepended to the IPv4 *next_hop*. The technique can be best described as automatic tunneling of the IPv6 packets through the IPv4 backbone. The MP-BGP relationships remain the same as they are for VPNv4 traffic, with an additional capability of VPNv6. Where both IPv4 and IPv6 are supported, the same set of MPBGP peering relationships is used.

MPBGP is enhanced to carry IPv6 in a VPN fashion known as VPNv6, which uses a new VPNv6 address family. The VPNv6 address family consists of 8 bytes—a *Route Distinguisher* followed by a 16-byte IPv6 prefix. This combination forms a unique VPNv6 identifier of 24 bytes. The Route Distinguisher value has a local significance on the router, and the *Route Target* advertises the membership of the VPN to other provider edge routers.

Figure 3: 6VPE Packet Forwarding



In Figure 3, packet forwarding is explained showing end-to-end operation. When the ingress 6VPE router receives an IPv6 packet, destination lookup is done in the VRF table. This destination prefix is either local to the 6VPE (which is another interface participating in the VPN) or a remote ingress 6VPE router. For the prefix learned through the remote 6VPE router, the ingress router does a lookup in the VPNv6 forwarding table. The VPN-IPv6 route has an associated MPLS label and an associated BGP *next_hop* label. This MPLS label is imposed on the IPv6 packet. The ingress 6VPE router performs a PUSH action, which is a top label bind by the *Label Distribution Protocol* (LDP)/IGPv4 to the IPv4 address of the BGP *next_hop* to reach the egress 6VPE router through the MPLS cloud. This topmost-imposed label corresponds to the *Label Switched Path* (LSP). So, the bottom label is bound to the IPv6 VPN prefix through BGP and the top label is bound by the LDP/IGP. The IPv6 packet, now with two labels, gets label-switched through the IPv4/MPLS core router (provider routers) using the top label only (referred to as the *IGP label*). Because only the top label is of significance to the provider core, it is unaware of the IPv6 information in the bottom label.

The egress provider edge router, receives the labeled IPv6 VPN packet and performs a lookup on the second label, a process that uniquely identifies the target VRF and the egress interface. A further Layer 3 lookup is performed in the target VRF, and the IPv6 packet is sent toward the proper customer edge router in IPv6 domain.

In summary, from the control plane perspective the prefixes are signaled across the backbone in the same way as for regular MPLS/VPN prefix advertisements. The top label represents the IGP information that remains the same as for IPv4 MPLS. The bottom label represents the VPN information that the packet belongs to. As described earlier, additionally the MPBGP *next_hop* is updated to make it IPv6-compliant. The forwarding or data plane function remains the same as it is deployed for the IPv4 MPLS VPN. The packet forwarding of IPv4 on the current MPLS VPN remains intact.

6VPE Design Recommendations and Considerations

The following sections identify general recommendations that should be considered when deploying IPv6 in a service provider network:

Working with Enterprise Implementations

Typically *Customer Metropolitan-Area Networks* (C-MANs), also known as *Campus Networks* or *Customer LAN* (C-LAN) elements, form the enterprise network, whereas the 6VPE and customer edge provide the entry point into network access. IPv6 can be supported partially or fully on an enterprise network. In situations where enterprise-wide IPv6 deployment does not exist, network administrators can elect to tunnel the IPv6 traffic toward the provider's customer edge or 6VPE. This can be done with 6-to-4 tunneling methods currently^[7]. So, if a site router within a C-MAN or C-LAN aggregates all IPv6 traffic and tunnels to a provider-managed customer edge or 6VPE router, then integration as well as migration becomes smooth. Therefore, it is important for the vendor and the customer to work together in determining the best approach.

Dual VRF Membership per Interface

RFC 2547 for IPv4 recommends one VRF per interface. When running dual stack on a 6VPE, multiple VRF configurations on a single physical or logical interface are required (IPv4 and IPv6). Each VRF instance configuration on a dual-stacked interface forms IPv4 and IPv6 address families. Each address family within VRF runs a VRF-aware routing protocol—such as static routing (static IPv6 unicast routing for IPv6), BGP (BGP with IPv6 enhancements for IPv6), OSPF (OSPFv3 for IPv6), or *Routing Information Protocol* (RIP) (RIPng for IPv6).

MTU Requirements

One important piece of information within the network elements is the capacity of the interface to transfer the size of datagrams. This is known as the *Maximum Transmission Unit* (MTU). The minimum link MTU for IPv4 packets is 68 bytes, whereas for IPv6 the minimum MTU should be 1280 bytes. While designing and planning for IPv6 support, the network elements should be examined along with interfaces and underlying network technologies to ensure the MTU requirements.

Dealing with Link-Locals

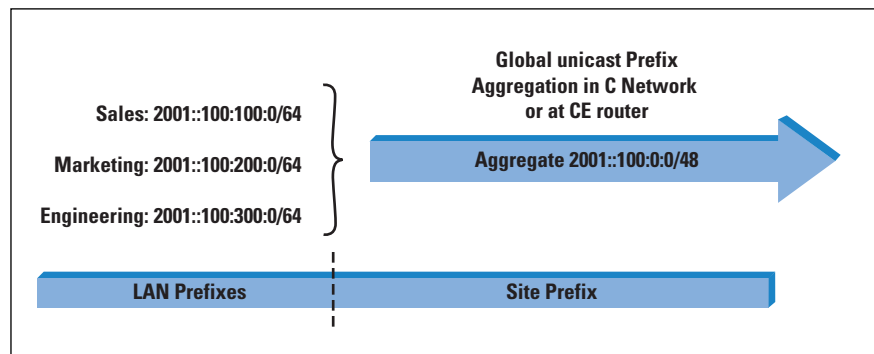
Because link-local scope addresses are defined as uniquely identifying interfaces within a single link, only those may be used on the provider edge-customer edge link.

However, they are not supported for reachability across IPv6 VPN sites and are never advertised with MPBGP to remote provider edges. As outlined in the RFC for IPv6 address assignments, the link locals (**FE80::x**) should not be advertised outside their local scope. Because the link-local addresses are embedded on the IPv6-enabled interface for certain local tasks, the link-local addresses are not and should not be advertised anywhere outside the local link scope, including the customer edge and 6VPE running IPv6. Globally unique aggregatable IPv6 prefixes are defined as uniquely identifying interfaces anywhere in the network. These addresses are expected for common use within and across IPv6 VPN sites. They are obviously supported by this IPv6 VPN solution for reachability across IPv6 VPN sites and advertised through MPBGP to remote provider edges.

Router Capacity Impact

Dual-stacking also introduces another task, namely hardware analysis to determine the resource capacity, that is, CPU and memory usage. Increased memory consumption may occur because of the dual-stack *Routing Information Base* (RIB). It also has implications for the *Interface Descriptor Block* (IDB) and *Routing Descriptor Block* (RDB) limits of hardware. The IDB limit is the capacity of particular equipment to support a number of physical and logical interfaces, whereas the RDB limit is the number of routing protocols and instances supported on such equipment. Typically these values (limits) are very high, but 6VPE is such an important element of the MPLS network that these facts must be considered. From a business case perspective, scalability, high aggregation, and rapid Return on Investment are expected, hence it is important to consider these factors in the design.

Figure 4: Route Aggregation



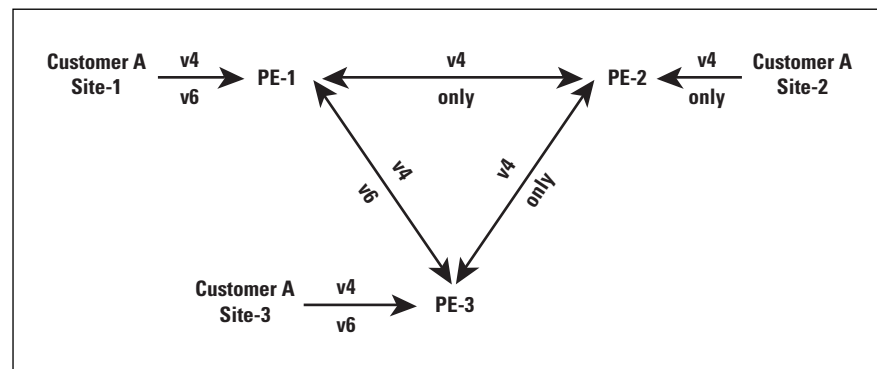
Router Memory Impact

The memory challenges can occur also when large numbers of IPv6 prefixes are advertising toward service provider network elements. In that event, the enterprise on the C-LAN or service provider on the customer edge router may elect to perform route aggregation. IPv6 prefixes can be aggregated to their higher-level significant boundary. Figure 4 shows an example of IPv6 prefix aggregation. Moreover, when a packet arrives on a dual-stacked interface (VRF-aware interface), the 6VPE router determines the packet version number by looking into the IP header. The per-packet header lookup is normally performed (it is a basic router function), but the extra work required by the router is to determine the version number. This additional task creates a longer processing cycle.

The Address Family Identifier and its Importance

All the elements referenced as dual-stacked, such as provider edge and customer edge routers, run IPv4 as well as IPv6 addressing and routing protocols. The 6VPE elements can also mix and match VPNv4 and VPNv6 peering sessions with other 6VPE routers or with route reflectors. What does the term “mix and match” mean here? It was an important enhancement to traditional BGP when MPBGP extensions were introduced. The address family within MPBGP is modular to facilitate distinct peering relationships, and is expressed using the *Address Family Identifier* (AFI). The regular BGP capabilities are exchanged after the peering sessions are turned on. In order for two provider edge routers to exchange labeled IPv6 VPN prefixes, they must use BGP capabilities negotiation to ensure that they both are capable of processing such information. When the service provider network is running VPNv4 peering sessions with other respective elements in the network, it exchanges the VPNv4 AFI capabilities with others. When the VPNv6 peering sessions are turned on, it renegotiates the capabilities and fresh peering sessions are established. The peering sessions established are based on common features if either of the peers does not agree on any of the capabilities.

Figure 5: VPNv4 and VPNv6 AFI



In Figure 5, three provider edge routers out of two need to exchange VPNv6 traffic, but all three provider edge routers need to maintain their existing VPNv4 capabilities. This is possible with the AFI configuration feature, which makes the migration steps very smooth. Service providers can mix and match VPNv4 and VPNv6 provider edge routers as required. Functions of 6VPE can be turned on when and where required. If the customer edge routers are dual-homed to different provider edge routers, the integration of customer IPv4 and IPv6 networks becomes painless. This scenario outlines hybrid environments, but it does not address the IPv4 and IPv6 communication. Consider techniques such as *Network Address Translation* (NAT) or application layer gateways for the IPv4 and IPv6 communication.

Route Reflectors for MP-IBGP

For advertising VPN membership, provider edge routers peer with VPNv4 route reflectors for scalability, thereby avoiding the need for full-mesh MP iBGP sessions among all provider edge routers. The same concept is supported for VPNv6. The same VPNv4 route reflectors can be upgraded to support VPNv6 address families.

Route reflectors can also make addition or removal of a provider edge router from a network simple and flexible. Alternatively, the BGP confederation option can also be deployed to provide MPBGP peering sessions among provider edge routers.

QoS Considerations

Service providers operating customers' MPLS VPN networks and also providing *Quality of Service* (QoS) should account for the new introduction of IPv6 and its impact. QoS and queuing of important application traffic requires distinct policies for IPv4 and IPv6, in turn possibly requiring additional operational tasks where IPv4 and IPv6 networks coexist. Other design considerations should be made to account for each individual network. Both IPv4 and IPv6 have a commonality, which is the 3-bit IP *Precedence* (or *Type-of-Service* [ToS]) field within the IP headers. Alternatively, the *Differentiated Services* (Diff-Serv)-compliant QoS models can also be employed. Irrespective of the technique, QoS is an important factor when low-speed links are concerned. However, there is no additional advantage of QoS on IPv6 versus IPv4. At some point in the future IPv6 can be different by using the flow label in the IPv6 header. QoS within the MPLS core remains *MPLS Experimental Value* (MPLS_EXP)-based and is untouched but still is effective with the addition of IPv6.

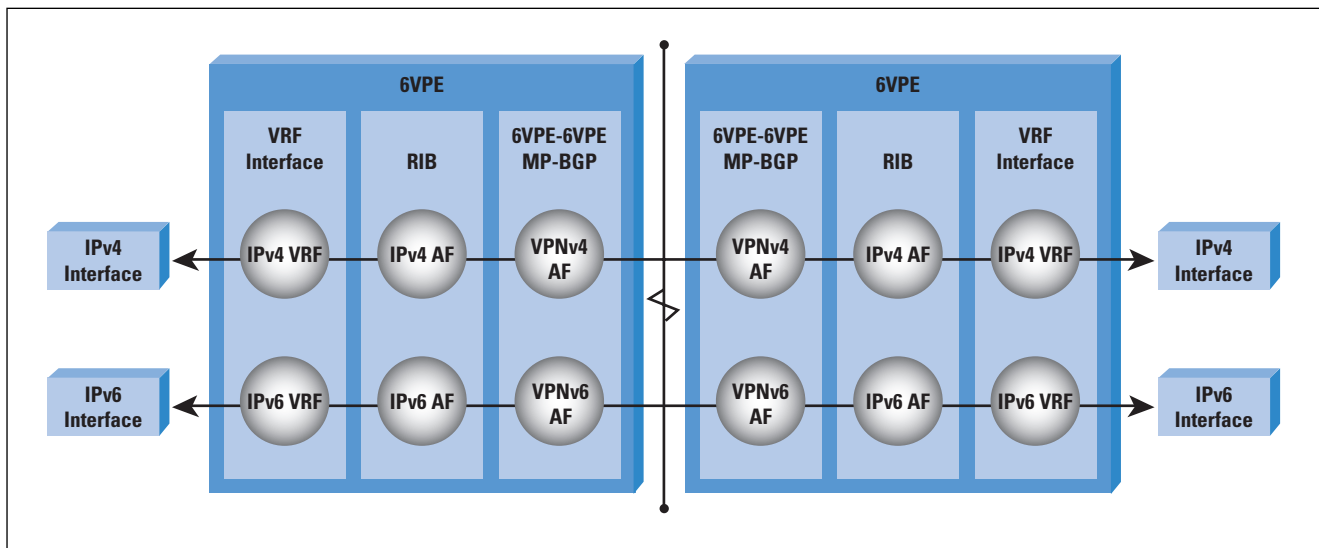
Device Management

Finally, device management is another important aspect that service providers must consider. Device management in a dual-stacked network can be done through an IPv4 or IPv6 address. Where the IPv6 VPN service is supported over an IPv4 backbone, and where the service provider manages the customer edge, the service provider can elect to use IPv6 for communication between the management tool and the customer edge for such management purposes. The management systems, including *Operations-Support-System* (OSS) servers, need to be aware of IPv6 and must run proper *Simple Network Management Protocol* (SNMP) stacks in order to perform IPv6-based management. From the VPN perspective it still remains transparent how the device and services are managed.

Enhancements to the Draft

The current MPLS VPN services that service providers have implemented are based on RFC 2547bis, the Internet Draft required to enhance the Layer 3 VPN approach further to address the IPv6 support. The "BGP-MPLS VPN extension for IPv6 VPN"^[1] is the current Draft that addresses the need for IPv6 support over MPLS networks in a VPN environment. Also, to avoid an extra layer of signaling, the Draft addresses the scalable automatic tunneling of VPN-based IPv6 prefixes. The basic functions remain the same as outlined in RFC 2547. Some of the extensions outlined will require additional work in order to be effective in the service provider network.

Figure 6: Dual Mode 6VPE AFI Model



The standard RFC 2547bis introduces “address family” concepts, as well as MPBGP to carry VPN information across the MPLS network. This enables formation of a full mesh between customer sites. The provider edge routers advertise their VPN membership to other provider edge routers through direct iBGP or value(s). As shown in Figure 6, these new address families are introduced to support IPv6 within VPN, IPv6, and VPNv6. If configured for dual stacking, the interface belongs to multiple VRF instances, IPv4 and IPv6. Each instance maintains its own RIB. MPBGP is now capable of handling the VPNv6 address family to advertise the IPv6 prefix across the VPN.

Summary

“Staying abreast of the best” has always been challenging for service providers when it comes to technology deployment or support. Time to market is another challenge. This article provides a view of the service provider challenges. In this new era where explosive use of IPv6 is envisioned, it is extremely important for service providers to have a simplified, automated, fail-proof, and cost-effective network design. The Internet Draft discussed advances the capabilities to achieve this and allows service providers to take a practical approach in supporting IPv6 for customers’ next-generation applications. The Draft brings service providers closer to the IPv4-to-IPv6 transition with a simple, cleaner, cheaper, and scalable solution.

For Further Reading

- [1] Jeremy De Clercq, Dirk Ooms, Marco Carugi, Francois Le Faucheur, “BGP-MPLS VPN extension for IPv6 VPN,” **draft-ietf-13vpn-bgp-ipv6.06.txt**, February 2005.
- [2] Eric Rosen and Yakov Rekhter, “BGP/MPLS VPNs,” **draft-ietf-13vpn-rfc2547bis-03.txt**, October 2004.
(See also RFC 2547, March 1999, by the same authors.)
- [3] Mallik Tatipamula, Patrick Grossetete and Hiroshi Esaki, “IPv6 Integration and Coexistence Strategies for Next-Generation Networks,” *IEEE Communications Magazine*, Vol. 42, No. 1, January 2004.
- [4] Bates, Chandra, Katz, and Rekhter, “Multiprotocol Extensions for BGP4,” RFC 2858, June 2000.
- [5] Deering, S. and R. Hinden, “Internet Protocol, Version 6 (IPv6) Specification,” RFC 2460, December 1995.
- [6] Rekhter and Rosen, “Carrying Label Information in BGP4,” RFC 3107, May 2001.
- [7] Carpenter, B. E., Moore, K., Fink, R., “Connecting IPv6 Routing Domains Over the IPv4 Internet,” *The Internet Protocol Journal*, Volume 3, No. 1, March 2000.

TEJAS SUTHAR holds CCIE # 8423. He is working as a Service Architect at TELUS Communications Inc. in Toronto. He focuses on Converged Network designs for customers in various industry sectors. He is very active in IP-related deployments. E-mail: **tejas.suthar@gmail.com**

Graph Overlays on Path Vector: A Possible Next Step in BGP

by Russ White, Cisco Systems

Over the past several years, much research and thought has gone into a replacement for the current interdomain routing protocol, *Border Gateway Protocol* (BGP)^[1]. For instance:

- In 2002, the *Internet Research Task Force* (IRTF) published a set of requirements for a next-generation interdomain routing protocol. In fact, several sets of requirements documents have been published in this area.
- In December 2001, *The Cook Report* noted that BGP needs to be replaced^[2]:
- In October 2003, the *Workshop on Internet Routing Evolution and Design* (WIRED) presented papers arguing that BGP needs to be replaced^[3].
- In December 2001, the IETF published RFC 3221^[5], authored by Geoff Huston, which provided some background information toward finding a replacement for BGP.

There are probably thousands of references in magazine articles, conference proceedings, and research papers, all stating that BGP should be replaced. Of course, all these discussions wind up at the same place: It is almost impossible to replace BGP, wholesale, in the public Internet, or even in any of the private networks running BGP today.

The basic problem is you cannot take the network down, and you cannot replace the routing protocol without taking the network down. Many very clever ideas have been proposed to get around this problem—complex transition schemes, moving partitions, and all sorts of other concepts. But, in the end, the idea of transitioning from one routing protocol to another on something as large—and as distributed in both geography and ownership—as the Internet, has been a hard wall against which all the proposals for new interdomain routing protocols pile up. In an article^[4] here in *The Internet Protocol Journal*, Geoff Huston states:

“Another approach is to consider the feasibility of decoupling the requirements of inter-domain connectivity management with the applications of policy constraints and the issues of sender- and receiver-managed traffic-engineering requirements. Such an approach may use a link-state protocol as a means of maintaining a consistent view of the topology of inter-domain network, and then use some form of overlay protocol to negotiate policy requirements of each Autonomous System, and use a further overlay to support inter-domain traffic-engineering requirements.”

In this article, we propose building on this concept, but in a novel way: rather than replacing BGP, or attempting to solve all the currently perceived problems with BGP at once, we attempt to address two problems in a way that does not heavily modify day-to-day BGP operation. Rather than replace BGP, enhance it to account for new requirements by providing new capabilities. If done right, this avoids the problem of deploying a new routing protocol altogether, because BGP is already deployed throughout the Internet.

Problems with BGP

No discussion of replacing BGP would be complete without a discussion of why so many people think BGP needs to be replaced. We need to consider three main points in this area: *convergence speed*, *policy*, and *security*. Each of these is covered in the following sections.

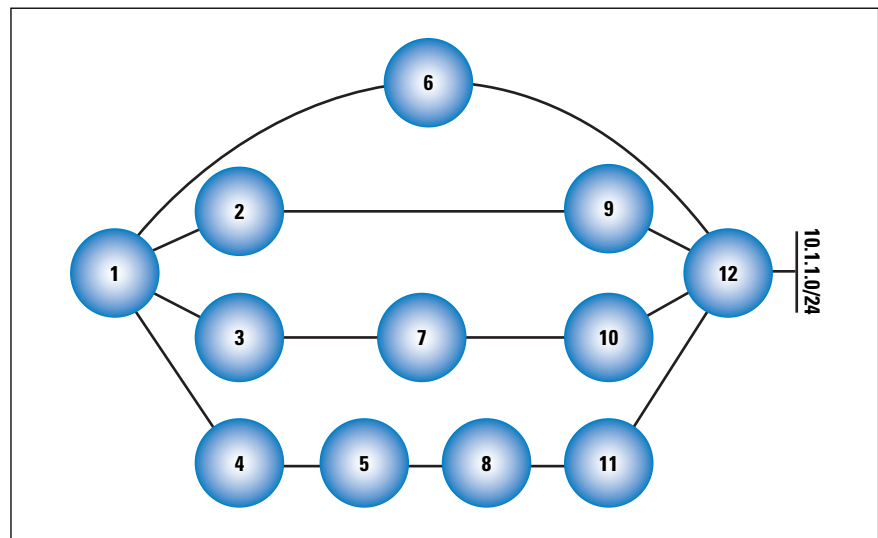
BGP Convergence Speed

Through various studies, and through examining the way in which BGP works, it has been shown that BGP, in an interdomain environment, always converges roughly in:

$$(\text{Maximum AS_PATH} - \text{Minimum AS_PATH}) \times \text{Minimum Advertisement Interval}$$

To understand why this is so, let's examine the following small internet-network as it converges.

Figure 1: An Example Internetwork Using a Path Vector Protocol



Let's assume autonomous system (AS) 12 is advertising some destination, 10.1.1.0/24, and that every other autonomous system in the internetwork chooses the path to the right to reach that destination. So, for instance, AS4 chooses the path {5,8,11,12} to reach 10.1.1.0/24, AS3 chooses the path {7,10,12} to reach 10.1.1.0/24, AS2 chooses the path {9,12} to reach 10.1.1.0/24, and AS6 chooses the path {12} to reach 10.1.1.0/24.

At this point, let's examine what happens if AS12 loses its connection to 10.1.1.0/24. AS12 sends out a withdraw, which reaches AS6, 9, 10, and 11 at about the same time. These autonomous systems then send out withdraws, with the second set of withdraws reaching AS1, 7, and 8 at about the same time.

When AS1 receives this first withdraw, it examines its local table, and finds the next best path to reach 10.1.1.0/24 is through AS2, with the path {2,9,12}. AS1 does not realize that AS2 has received a withdraw for 10.1.1.0/24 at the same time it received the first withdraw for this destination from AS6. So, AS1 switches over to its next best path, and continues forwarding traffic to 10.1.1.0/24.

AS2, 7, and 8 now also send withdraws to each of their peers, including AS1, 3, and 5. AS1 now receives another withdraw, again for the path it is currently using to reach 10.1.1.0/24. AS1 examines its local tables and finds it has another path, through {3,7,10}, to 10.1.1.0/24, so it switches to that path, without knowing AS3 has just received a withdraw for this same path. AS3 and 5 now send withdraws to each of their peers, AS1 and 4. AS1 has again received a withdraw from the peer it is using to reach 10.1.1.0/24, so it examines its local tables, and finds it still has a path through {4,5,8,11,12} to reach this destination. It switches to this path, without realizing AS4 has just received a withdraw as well.

AS4 now sends the final withdraw to AS1, removing AS1's final path from its local tables. AS1 now removes all reachability information for 10.1.1.0/24, and the network is converged. Note that the actual convergence in this situation would be a bit more complicated, with AS1 sending updates at each stage, and all the other autonomous systems re-converging at each step along the way, but we have used only the simplest set of messages through the network, to illustrate the basic procedure BGP follows when converging.

This short example illustrates why BGP has the convergence characteristics described previously. BGP "hunts" through each possible autonomous-system path, from shorter ones to longer ones, until it finally converges. The rate at which it can hunt through each possible autonomous-system path is determined by the minimum advertisement interval, the rate at which new routing information is allowed to flow through the system.

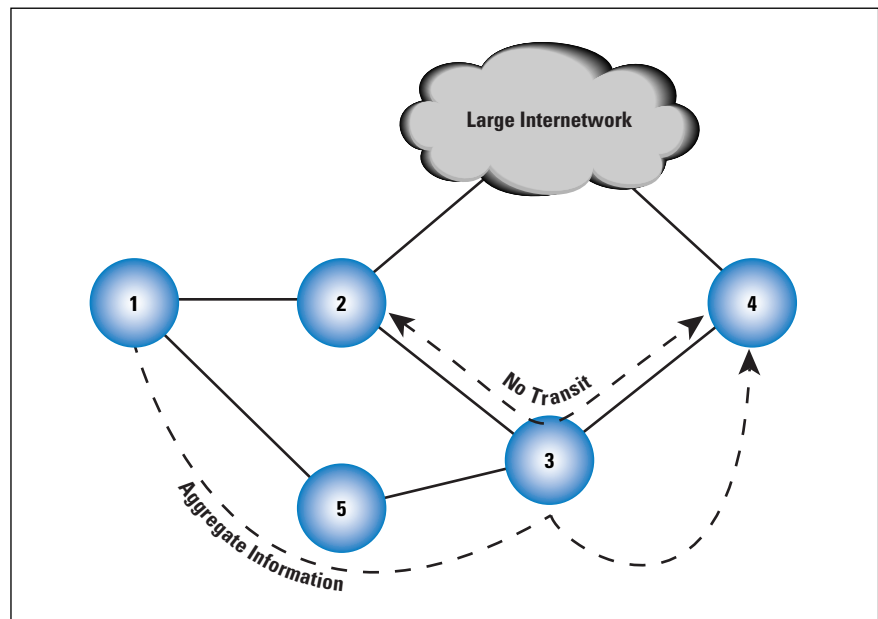
This problem has several obvious solutions. The first is to simply increase the rate at which routing information flows through the system, by reducing the minimum advertisement interval. But, this plays against route flap dampening, and network stability in general, so, beyond some lower possible bound, reducing the minimum advertisement interval is not possible (without further modifications to BGP).

Another obvious solution is to simply add a “reason code” to the original withdraw. If AS12 originally stated it was withdrawing reachability to 10.1.1.0/24 because it had lost local connectivity to it, then all paths with AS12 in the path could have been discarded immediately, at the first step. The problem here is making certain the original withdraw message actually makes it through the network, from AS12 all the way to AS1. Because BGP is a very efficient protocol, many control messages of this type are actually removed from the network, through implicit withdraws, aggregation, and other mechanisms.

Policy

The second problem we encounter with BGP is its rather rough sense of policy. For instance, let’s examine the following small network, and look at one specific example of where policy transmission and enforcement are problematic in BGP.

Figure 2: Issues with Policy Transmission in a Path Vector Protocol



Here AS2 has a policy that AS3 should never be used for transit. In other words, traffic originated in AS4 should always pass through the large internetwork rather than through AS3 to reach AS1. This type of situation is very common in the public Internet, such as when AS3 is actually AS2’s customer. How can AS2 communicate this policy to AS4, however?

AS2 could simply mark the routing information it sends to AS3 so AS3 cannot readvertise it to AS4, but this is problematic. Simple mechanisms, such as marking the routes with the NO_EXPORT community, are easy for AS3 to simply strip off the routing information it receives. We could conceive of some way to cryptographically sign the included policy, so AS3 cannot disturb the policy and AS4 can see the policy when it receives the information from AS3, but this is problematic as well.

Suppose AS3 is receiving aggregated routing information directly from AS5, which includes some of the same destinations AS2 has advertised to AS3, but has blocked AS3 from advertising to AS4. AS3 could, conceivably, readvertise this routing information to AS4, and AS4 could prefer this shorter prefix aggregate to reach the destinations in AS1, rather than the paths through the large internetwork. AS4 would then forward traffic to AS3, which would then rely on its longer prefix routes, received from AS2, to forward this traffic to these destinations in AS1. AS3 is, contrary to AS2's policy, transiting traffic through AS2 to AS1. There is no simple answer to this problem.

Security

It has been widely acknowledged that BGP is an insecure protocol, with many areas where attackers can hijack, inject false routing information, and perform other attacks. The IETF's *Routing Protocols Security* (RPsec) working group is working on a set of documents describing vulnerabilities of BGP, and creating recommendations for systems to secure BGP. For the latest information about these Drafts, refer to the RPsec homepage at: <http://www.rpsec.org>

What sort of requirements are likely to come out of such an undertaking?

- Any proposed mechanism must be able to show that a specific autonomous system is authorized to originate specific routing information.
- Any proposed mechanism must be able to show that the AS Path carried in received routing information corresponds to a real path in the internetwork, beginning with the origin AS and ending in the advertising peer.

There will be many other requirements that proposed mechanisms for providing security for BGP will need to, or should, meet, but these two will be the largest areas of concern for our purposes.

Solving the Problems

Now that we have an idea of the three areas we want to solve problems in, how can we actually solve them? The most elegant solution would be a single mechanism that does not change the current semantics of BGP itself too greatly, would provide greater benefits as it is deployed throughout a large-scale internetwork, and would rely on existing—and understood—techniques within routing.

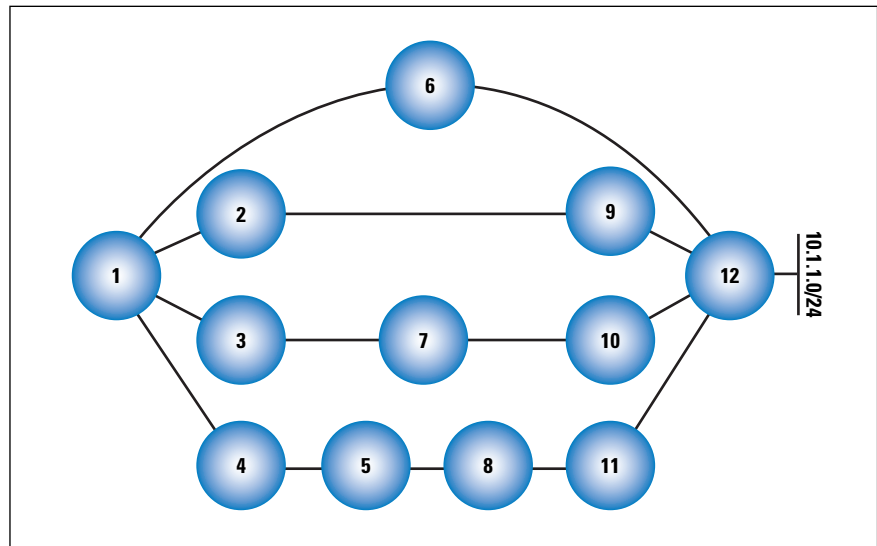
One perfect example of such a mechanism would be to simply overlay a link state-like graph of interconnectivity over the BGP protocol. This graph would provide information about the interconnections between autonomous systems, rather than between routers, and would be used to convey information about the topology and policies in the internetwork, rather than to find loop-free paths through the internetwork.

Let's go back through our three examples, and see how overlaying an internetwork connection graph would be able to solve some of the problems currently facing BGP.

Convergence Speed

Looking at our small sample internetwork again:

Figure 3: An Example Internetwork Using a Path Vector Protocol



What is the one thing we said would resolve the problems with BGP hunting through every possible longer autonomous-system path alternative to finally converge around loss of reachability to 10.1.1.0/24? Could AS12, somehow, communicate directly to every autonomous system in the internetwork that it has directly lost this connection, rather than waiting for AS1 to try every possible path to 10.1.1.0/24, and discover each one, in turn, withdrawn?

If we had a topological graph of the network, AS12 could simply remove 10.1.1.0/24 from its connectivity information. AS12 would then flood this information, on an interdomain basis, to all the other autonomous systems in the internetwork at roughly the same time. Thus, in the worst case, AS1 would receive this information at about the same time it received the first withdraw for 10.1.1.0/24, from AS6.

When AS1 receives this updated topology information from AS6, it will discover that AS12 is no longer connected to 10.1.1.0/24, and, therefore, it can remove every possible path to 10.1.1.0/24 containing AS12. This would allow AS1 to remove the paths {2,9,12}, {3,7,10,12}, and {4,5,8,11,12} at the same time. The internetwork now converges as soon as AS1 computes the new connectivity graph, and acts on it by examining each entry in its local tables and discarding the ones with AS12 in the autonomous-system path.

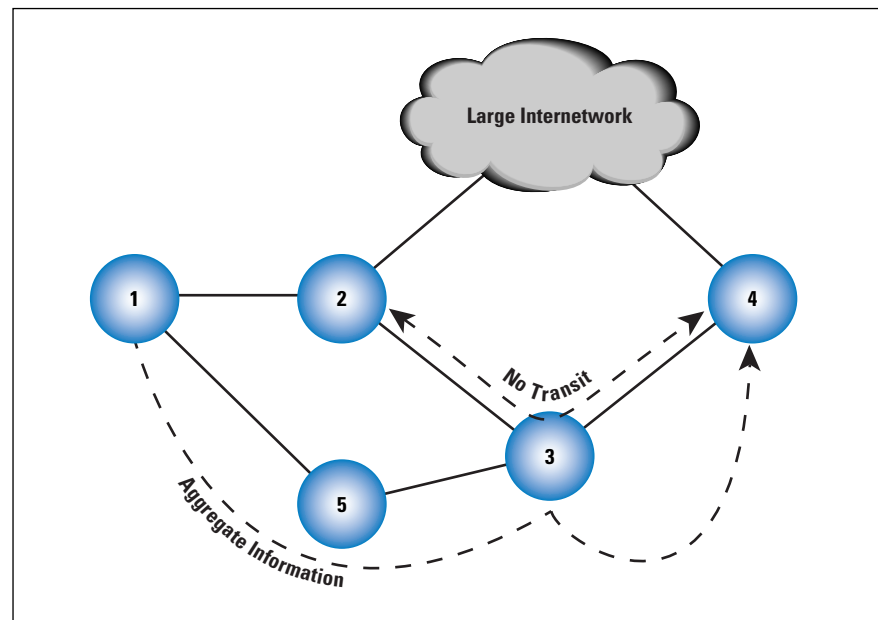
We have not changed the way BGP finds paths through the network—the path still is not valid unless we receive an advertisement from our connected peers. We have also not changed the format of any BGP updates, any peering state machines, or anything else. We have simply overlayed an interconnection graph on top of the current protocol mechanisms, which we can use to our advantage to speed up network convergence.

What about partial deployments in this situation? Suppose only autonomous systems 6, 7, 8, 9, 10, 11, and 12 are running this new extension. Would it still help us to speed up network convergence? When AS12 withdraws 10.1.1.0/24, AS6, 7, 8, 9, 10, and 11 would immediately discard any routes passing through AS12 to reach 10.1.1.0/24. At this point, they could each withdraw those routes, meaning AS1, 2, 3, and 5 would all receive a withdraw at about the same time. This short-circuits the number of possible paths for AS1 to hunt through, decreasing the amount of time the internetwork takes to converge. Even without a full deployment, we see some positive impact from this new technique.

Policy

Let's examine our policy problem after placing our interconnection graph on top of the internetwork.

Figure 4: Injecting Policy on an Interconnection Graph



Here, we see that AS2 could actually place its policy for AS3 not to transit traffic in the interconnection draft. AS4 would then be able to independently verify what AS2's policy toward AS3 transiting traffic is. AS4 could then examine the routing information it receives from AS3, and determine if it should install—or not install—routing information received from AS3, based on this policy.

Objections to an Interconnection Graph

When a link-state protocol has been proposed as a possible replacement for BGP in the past, two primary objections have been raised:

- Providers are reluctant to accept the wholesale replacement of a known working system with a new one.
- Many providers wish to hide their policies and connectivity to other providers or customers for policy reasons.

This article does not propose replacing BGP, just augmenting it, so the first argument is, to some degree, not valid against this approach. The second objection, that of using a link-state protocol for interdomain routing specifically, also does not apply, because we are not proposing changing the way BGP finds loop-free paths through the network. The proposed interconnection graph is not used for finding paths through the network, it is used only for faster signaling of path failure (by short-circuiting the slower withdraw mechanisms), and for providing a place to hang policy and security information.

Concentrating on a few smaller spaces allows us to design a smaller solution set that can be incrementally deployed in a simple way.

The second objection is harder to meet, simply because the concepts of policy within a routing system are hard to define and understand in all possible cases or respects. In fact, there are policy requirements not met by BGP today, but rather are met through contracts, packet filters, and other mechanisms (even sometimes by violating the BGP specification).

Consider two facts about this proposal that work around many of the specific objections we have heard in this area:

- The interconnection graph can be partial, in different parts of the internetwork. For instance, a given service provider might provide different views of who they are connected to to different peers, depending on their policy of revealing this information.
- The interconnection graph only contains autonomous system-level connectivity information, not specific peering-point information. For instance, two autonomous systems may be connected in a large number of places, or as few as one. The interconnectivity graph does not care about such details, only whether at least one connection exists. Such an interconnectivity graph would not reveal actual connection points between peering autonomous systems, how rich that connectivity is, nor any other information about the business relationship between the two peers.

In fact, the types of interconnectivity information an interconnection graph could provide is already available by examining the autonomous-system paths of routes retrievable from various route view servers. Some mechanism would be required to collate this information into a usable graph, but a good deal of current research on the scaling and convergence properties of large-scale internetworks actually depends on the ability to build an interconnection graph before beginning any other work, so mechanisms to collate this data already exist, and are in use today.

Security

The internetwork interconnection graph can actually show whether a path exists from the origin to the advertising peer, through *signed certificates*. For example, soBGP^[6] (<ftp://ftp-eng.cisco.com/sobgp/index.html>) uses this specific mechanism to validate the autonomous-system path carried in received routing information. Other research is currently being pursued in this area as well.

Summary

We have proposed a single step forward that could be used to resolve some of the problems facing BGP in the near term, and possibly provide the networking community with a path forward on other fronts as well. The concept of simply making incrementally deployable changes to BGP to solve pressing problems can provide us with options outside the normal lines of thinking: either making very small changes to BGP, making BGP more and more complicated, or simply replacing the BGP protocol, with all the deployment problems this would entail.

References

- [1] Yakov Rekhter, Tony Li, “A Border Gateway Protocol 4 (BGP-4),” RFC 1771, March 1995.
- [2] <http://www.cookreport.com/10.09.shtml>
- [3] <http://www.net.informatik.tu-muenchen.de/wired/position/bruce.html>
- [4] Geoff Huston, “Scaling Inter-Domain Routing—A View Forward,” *The Internet Protocol Journal*, Volume 4, No. 4, December 2001.
- [5] Geoff Huston, “Commentary on Inter-Domain Routing in the Internet,” RFC 3221, December 2001.
- [6] Russ White, “Securing BGP Through Secure Origin BGP,” *The Internet Protocol Journal*, Volume 6, No. 3, September 2003.

RUSS WHITE works for Cisco Systems in the Routing Protocols Deployment and Architecture (DNA) team in Research Triangle Park, North Carolina. He has worked in the Cisco Technical Assistance Center (TAC) and Escalation Team in the past, has coauthored several books on routing protocols, including *Advanced IP Network Design*, *IS-IS for IP Networks*, and *Inside Cisco IOS Software Architecture*. He is currently in the process of publishing a book on BGP deployment, and is the co-chair of the Routing Protocols Security Working Group within the IETF. E-mail: riw@cisco.com

Book Reviews

A Brief History of the Future

A Brief History of the Future—The Origins of the Internet, by John Naughton, ISBN 0-75381-093X, 2000, Published by Phoenix,
<http://www.orionbooks.co.uk>

This is a well-written book by a well-known Irish academic and journalist, which charts the growth of the Internet from a 1950s military project to the pervasive networking infrastructure that dominates the IT world today. It is relevant to the readership of this journal because it charts the growth of the technology that underpins the IP world—and it gives a sound understanding of the culture and approach that led to the development of the Internet as we know it.

Naughton takes the reader from the inception of the *Advanced Research Projects Agency Network* (ARPANET) through most of the major developments such as packet switching, mail, TCP/IP, and the Web, not only covering the technology, but also providing insights into the background of the Internet pioneers and the political environment.

Organization

The book is divided into three major sections, the first of which is largely concerned with scene setting and is aimed at bringing those less familiar with the subject area up to speed. In the first chapter, Naughton likens the evolution of the “Net” to that of amateur radio, moving on in succeeding chapters to cover basic technology and to provide some perception of scale and rate of growth.

The second part of the book covers the growth of the Internet up to the early 1990s. This starts by looking at the origins of the ARPA project, noting the influence of MIT and important figures such as Vannevar Bush, Norbert Weiner, and J.C.R Licklider. Naughton describes how ARPA was initiated and its relationship with NASA and academia, highlighting the desire to provide time-sharing systems and the breakthrough concept of the *Interface Message Processor* (IMP) as a solution to the “n-squared” problem. This is followed by two chapters that discuss the adoption of packet switching as the underlying technology, following its initial proposal by Paul Baran and further development by Donald Davies’ team in the UK.

Naughton next examines how e-mail became the first “killer application” that drove up Internet usage, even telling the reader where the use of the ubiquitous “@” symbol comes from. He then considers the maturing network during the 1970s, discussing the formulation of the first *Request For Comments* (RFCs), the development of the gateway concept, and the evolution of TCP/IP. The discussion leaves the network area, concentrating on the evolution of UNIX and its impact, stressing the role of AT&T’s regulatory situation. Then Naughton considers how this accelerated the development of USENET.

In a chapter called “The Great Unwashed,” Naughton discusses the popularization of computing and networking, through the availability of the PC and the evolution of readily available file transfer tools such as X-Modem and the creation of bulletin board systems such as fido-net. He then considers the development of Open Source, telling the story of Linux and its derivation from MINIX.

The third section of the book deals with the emergence of the World Wide Web, tracing it back through the original ideas of Vannevar Bush and Ted Nelson, to its ultimate development by Berners-Lee at CERN. He links this to the subsequent development of Mosaic at NCSA and shows the dramatic impact this had on Internet growth.

Naughton concludes his book by looking at the prognosis for the “Net.” Here he refuses to try to predict the future; instead he analyzes the forces that will drive the future of the Internet and discusses their impact in the past and hence their potential impact. At the end of the book, he provides notes and references for each chapter, a short section on the sources he consulted, and a comprehensive glossary.

Synopsis

I found this book provided excellent insights into the development of the Internet, adding a lot of perspective to the engineering field I currently work in. Naughton places appropriate emphasis on the technical, personal, commercial, and political factors that have steered its evolution. He is not afraid to disturb the reader’s preconceptions by looking at things from unusual angles, and he emphasises the importance of *timing*. This is apparent when he points out that according to many sources, most of the important inventions around the Internet have come from graduate students, rather than the professors they work for. He similarly recounts the story that AT&T turned down the opportunity to run the “Net” in the early 1970s and reflects the view that if the Internet had not existed we could not invent it now.

This is an excellent read (it was nominated for the *Aventis Prize* in 2000), which helps the reader understand the How, When, Where, and Why of the Internet’s development. It covers most of the major milestones in the evolution of our discipline and is very well-written.

The Author

John Naughton is Professor of Public Understanding of Technology at the Open University, and he writes a weekly column in *The Observer* Business Section, covering important developments and trends in the IT industry. He describes himself as a “Control Engineer with a strong interest in systems analysis and computer networks” and is a Fellow of Wolfson College, Cambridge.

—Edward Smith, BT, UK
edward.a.smith@btinternet.com

Eats, Shoots & Leaves *Eats, Shoots & Leaves*, by Lynne Truss, ISBN 1-592-40087-6, Gotham Books, 2003.

Eats, Shoots and Leaves is a book about punctuation, but boring it is not. Informative and delightful it is. Lynne Truss includes in the book—which she says is not about grammar—wonderful examples of misused and misplaced punctuation marks. She claims to have written the book to unite us sticklers who do care about the written word, and how we communicate through it. We sticklers cringe with many misuses of punctuation, and we are cringing more and more often it seems.

Truss defines punctuation as a tool to clarify the written word, and who can argue with helps for clarification? She suggests that punctuation is dying, but then asks what would happen without it? Just imagine all the words in the first paragraph with no punctuation marks and no capital letters. You might be able to figure out its meaning with some work, but it would not be easy. Also consider, she suggests, the following:

A woman, without her man, is nothing.

A woman: without her, man is nothing.

Punctuation makes all the difference!

The book begins with a discussion of the apostrophe. Meaning “omission,” the apostrophe was first used in the 16th century. The most common egregious misuse of this tool is found in the word “it’s.” It’s translates “it is,” but it is often used as a possessive word, as in “The keyboard is useless; some of it’s keys are missing,” when it should be “The keyboard is useless; some of *its* keys are missing.” As a test, if you cannot substitute the words “it is” or “it has,” it should be “its;” if you can, it is correctly “it’s.” And the same is true for you’re and your. You’re translates “you are,” and your is the possessive (“It’s your turn”).

Another amusing example Truss gives is: Member’s May Ball. Of course it should be Members’ May Ball, because who would just one member dance with? Truss asks.

In her discussion of the comma, we learn that commas were first used 2000 years ago by Greek dramatists to show the actors where to pause or breathe. Then when printing was invented and used increasingly in the 14th and 15th centuries, a Mr. Aldus Manutius (1450–1515) developed italics, the semicolon, the comma, the colon, and full stops (we call them periods in the U.S.).

Truss is a master of the metaphor. She calls the comma the “sheepdog” of words. The comma organizes words, phrases, and groups of words that fit together. Consider one of her comma examples, a properly placed comma: No dogs, please.

Now think about that sentence without the comma: No dogs please. Now consider this: But many dogs *do* please. Thus the importance of the properly placed comma.

Truss addresses all the other marks, including semicolons, quotation marks, brackets, hyphens, parentheses, the four attention-grabbers: *italics*, the exclamation point, the dash —, and the question mark, and finally the ellipsis (the three dots ...). She tells us that, amazingly, someone actually did a PhD thesis on the ellipsis!

One chapter discusses the fact that proper use of punctuation steadily declined in the 20th century, many blaming the decline on television; and that it will continue to decline in the 21st century because of the Internet. E-mail messages cry for brevity, and brevity they get. For example, “**CU B4 8.**” “Netspeak” is, no doubt, here to stay. Language usage also is trending toward the deletion of spaces between words, so that now we say healthcare, chatroom, and the like.

And finally, Truss discusses the newest job that punctuation marks have assumed: emoticons. Examples include the smiley face :-), the sad face :-(, and many others, all made with common punctuation marks.

I thoroughly enjoyed this book, and recommend it to anyone who wants to learn while being entertained. It is a wonderful read.

—Bonnie E. Hupton, Editor
bhupton@sbcglobal.net

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. Contact us at ipj@cisco.com for more information.

Paul V. Mockapetris Wins 2005 ACM SIGCOMM Award

Paul V. Mockapetris, Chairman and Chief Scientist at Nominum Inc., is the winner of the 2005 *ACM SIGCOMM Award*. The SIGCOMM Award is widely recognized as the highest honor in computer networking. The Award recognizes lifetime achievement in and contributions to the field. It is awarded annually to a person whose work, over the course of his or her career, represents a significant contribution to the field and a substantial influence on the work and perceptions of others in the field. The SIGCOMM Award is presented to Dr. Mockapetris “in recognition of his foundational work in designing, developing and deploying the *Domain Name System* (DNS), and his sustained leadership in overall Internet architecture development.”

Paul Mockapetris created the original DNS protocol, wrote its first implementation, and worked with others to spread the DNS across the Internet. The design of DNS, which was the first major datagram protocol of the Internet, established a number of principles for key Internet infrastructural services. Its simplicity of design and fitness for purpose have stood the test of time. The strength of its design lies in a novel combination of hierarchy and caching that gives each organization absolute control over part of the namespace while simultaneously relying on caching to make the entire system efficient. Its success can be seen from the fact that DNS now handles many orders of magnitude more names and traffic than when it was first deployed, and yet the design and structure have remained intact. As a result the DNS design and caching mechanisms are often cited as two of the cornerstones on which the success of the Internet is built.

In addition to his work on DNS, Dr. Mockapetris’ career has included pioneering work on multiprocessor operating systems, virtual machines, and ring LAN technology. Further, Dr. Mockapetris played an important role in the deployment of networking technologies internationally. Starting during 1990–1993 as a program manager at ARPA, Dr. Mockapetris fostered the international deployment of multimedia conferencing, multicast, and QoS. His strong leadership in development of Internet architecture continued as Chair of the Internet Engineering Task Force during 1994–1996, as member of the Internet Architecture Board during 1994–1996, and then as member of the Federal Networking Council. Dr. Mockapetris is also a recipient of the *IEEE Internet Award* and is an ACM Fellow.

In summary, through his sustained effort in support of the Internet architecture, beginning with DNS and continuing through work at ARPA, IETF, and industry, Dr. Mockapetris has made far-reaching and influential contributions to computer networking. The 2005 SIGCOMM award recognizes Dr. Mockapetris for this lifetime record of achievement.

SIGCOMM is the *Special Interest Group (SIG) on Data Communication* of the *Association for Computing Machinery (ACM)*. SIGCOMM is a professional forum for the discussion of topics in the field of communications and computer networks, including technical design and engineering, regulation and operations, and the social implications of computer networking. The SIG's members are particularly interested in the systems engineering and architectural questions of communication. For more information please visit: <http://www.acm.org/sigcomm/>

Voice over IP (VoIP) And Government Policy

Voice over IP technology has the potential to provide much cheaper telephone service, particularly internationally. More importantly, it can enable exciting new services, such as voice-enabled Web pages and integrated phone, voice-mail, and e-mail. Unfortunately, some national governments are trying to limit its use. In late April, 2005, the *Advisory Committee on International Communications and Information Policy (ACICIP)* of the U.S. Department of State issued a very useful paper describing how VoIP works, the benefits it can provide, and what governments around the world are doing to promote or hinder its development.

Michael Nelson, the Internet Society's Vice President for Policy, represents ISOC on the Committee, and is helping draft "Version 2.0" of the paper, which will report on recent developments in additional countries. If you would like to make suggestions about the paper, please submit them to Michael Nelson at mnelson@isoc.org

For more information, see:

<http://isoc.org/pub/pol/pillar/voip-paper.shtml>

ISOC Commentary on the Status of the Work of WGIG, April 2005

When the first phase of the *World Summit on the Information Society (WSIS)* called on the UN Secretary General to set up the *Working Group on Internet Governance (WGIG)*, it was in the context of supporting the *WSIS Action Plan*. The Plan calls for concrete actions to advance the achievement of internationally agreed development goals by promoting the use of ICT-based products, networks, services and applications, and to help countries overcome the digital divide. This is, by the way, something the Internet community has worked hard to achieve since the very first days of the Internet.

These goals include those described in the *Millennium Declaration*. The 8th goal of that document is to develop a global partnership for development, which would make available the benefits of new technologies—especially information and communications technologies—in cooperation with the private sector for the benefit of all. This is the context (making the benefits of ICT available to everyone) in which we initially engaged in the WSIS and WGIG efforts. The Internet has a huge potential as an enabler bringing these benefits to people everywhere and we remain excited about the WSIS mission. However, it is not clear how WGIG's actions to date have helped support achieving such goals.

The *Internet Society* (ISOC) believes that the best way to extend the reach of the Internet is to build on those aspects that have worked well, for example, the long established open, distributed, consensus-based processes and many regional forums for the development and administration of the Internet infrastructure. Decision-making about issues such as resource allocation or IP Address Policy has always been in the hands of the Internet community, in order to be as close to those who require and use the resources as possible. It is this participative model, close to the end users, that led to the phenomenal, stable growth of the Internet. The Internet community and its bottom-up processes are constantly evolving in response to changes in needs and availability. For example, in response to moves by the African Internet community, the African countries now have their own *Regional Internet Registry* [RIR] (AfriNIC) that helps coordinate users' needs and IP Policy in that region. Latin America has the same story to tell. Support for the development of both these RIRs (educational, financial and boot-strapping of various processes) came from the global Internet community and primarily came from the other RIRs.

Developing and maintaining the Internet infrastructure are just two aspects of what has come to be referred to as *Internet Governance*. WGIG has pointed out that there are many others, and has recognized the fact that Internet Governance encompasses a much wider range of topics than IP address and domain name administration. However, much of WGIG's focus has been on Internet infrastructure, thereby missing an opportunity to focus on those aspects of the Internet's development that are less developed and that could benefit from improved, lightweight mechanisms facilitating an exchange of information between policymakers and the Internet community. Examples here are issues concerning inappropriate usage of the Internet—cybercrime and spam being just two examples. Much work has already been done on technical solutions to these issues, and many legal frameworks already exist for handling criminal activity such as fraud. The challenge today is to bring the lawmakers and policymakers together with the Internet community to discuss the most appropriate mechanisms to ensure the continued development of the Internet.

Many players have a role, and this clearly includes governments and intergovernmental organizations. WGIG had a clear mandate to not only develop a working definition of Internet Governance, but also to develop a common understanding of the respective roles and responsibilities of governments, existing intergovernmental and international organizations and other forums, as well as the private sector and civil society encompassing both developing and developed countries. Unfortunately an inordinate amount of time has been spent focusing on challenging current structures (those that brought us the Internet and its rapid, stable growth), rather than looking forward to the potential benefits of extended cooperation with (and based on the proven success of) existing models and structures. WGIG seems to have lost sight of this larger goal.

Also, many of WGIG's premises seem to start with an assumption that the Internet needs a hierarchical top-down governance model, thereby ignoring the decentralized, distributed structure on which the Internet was so successfully built. Not only does this "governance hierarchy" model prevent an accurate understanding of the Internet's infrastructure and development (forcing key organizations to be classed in prescribed categories that do not fit with the reality of their actions or their role in developing and supporting the Internet) but it also will very likely lead to conclusions that will harm the Internet's development and growth.

While WGIG appears to ascribe the growth of the Internet to deliberate regulatory decisions to liberalize telecommunications, in reality regulatory measures have been a relatively small factor. A more significant factor in the growth of the Internet has been the fact that the Internet architecture has enabled many tens of thousands of users to develop their own applications independent of the underlying architecture, thereby empowering people to add true value to the global Internet network. The continued expansion of the Internet to developing countries though will be greatly aided in the future by a more competitive telecommunications environment. We urge WGIG to recommend more concrete and aggressive action in this direction.

Further, WGIG has put great focus on comparing the relative merits of established treaty bodies and intergovernmental organizations to undertake a central role in the development of Internet infrastructure while very largely overlooking areas where attention and support are required and where national governments more naturally have a role to play, areas such as misuse of the Internet (cybercrime and spam to name a few). The limited perspective of this approach displays an obvious bias in the characterization of the issues and seems to pre-suppose a solution. In conclusion, we would urge WGIG to spend more time looking at what is actually being done to enable more people around the world to take greater advantage of the power of the Internet. This includes a focus on the many regional and global education activities that different Internet-related organizations are undertaking to "connect the unconnected."

These same organizations are also working to make the Internet more secure, more accessible, more reliable, more affordable, and more versatile. The development of the Internet, as well as many well-established capacity-building efforts could be jeopardized by applying a too heavy-handed approach to the operation and administration of this unique network of networks. Decentralized, lightweight governance has clearly proven itself to be a positive feature not a weakness. We want to encourage WGIG and WSIS to work with the Internet community within the already well-established Internet model to improve co-operation between policy makers and the Internet community.

In the spirit of meeting the international development goals highlighted by WSIS, any review of today's Internet model or structures must be carried out in the context of how well they have worked in the past, how well they meet the needs of the people who depend upon them today, and how well they will adapt to changing requirements in the future; and not simply focus on a comparison to other historical telecommunications or governance models. These historical models have not been demonstrated to be well suited to the Internet. For more information, see:

<http://isoc.org/>

<http://wgig.org/>

<http://www.itu.int/wsis/>

An interview with the new IETF Chair

IBM Distinguished Engineer and former ISOC Chairman Dr. Brian Carpenter has just taken over the role of IETF Chair. In a recent interview, Brian describes the future challenges facing the IETF and the Internet in general. The full interview is available here:

<http://resources.isoc.org/20503>

Upcoming Events

The *Internet Corporation for Assigned Names and Numbers* (ICANN) will meet in Luxembourg City, Luxembourg, July 11–15, 2005 and in Vancouver, Canada November 30–December 4, 2005. For more information see: **<http://www.icann.org>**

The *South Asian Network Operators Group* (SANOG) will meet in Thimpu, Bhutan, July 16–23, 2005. More info at:

<http://www.sanog.org>

The *Internet Engineering Task Force* (IETF) will meet in Paris, France, July 30–August 5, 2005 and in Vancouver, Canada, November 6–11, 2005. For more information, visit: **<http://ietf.org>**

ACM's *SIGCOMM 2005* will be held in Philadelphia, PA, August 22–26, 2005. For more information visit:

<http://www.acm.org/sigs/sigcomm/sigcomm2005>

The *North American Network Operators' Group* (NANOG) will meet in Los Angeles, October 23–25, 2005. For more information see:

<http://nanog.org>

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, Sr. VP, Technology Strategy
MCI, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 2005 Cisco Systems Inc.
All rights reserved.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA

The Internet Protocol Journal

September 2005

Volume 8, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
IPv4 Address Space Consumption	2
SSH Tunneling.....	20
Book Review.....	27
Fragments	30
Call for Papers	35

FROM THE EDITOR

Protocol transitions are never easy, particularly not when they involve something so fundamental as the *Internet Protocol* (IP). Organizations considering a move to IPv6 must consider many factors when deciding on the timing for such a deployment. One of the first questions that arises is: “When will the IPv4 address space actually run out, forcing us to use IPv6 instead?” That question is not a new one; it was being asked in the early 1990s when the IPv6 effort was started. Several factors, such as the deployment of *Classless Interdomain Routing* (CIDR) and *Network Address Translation* (NAT), have “delayed the inevitable,” and perhaps led to some complacency on the part of network operators. In this issue we examine the topic of IPv4 address space depletion in more detail. Our main article is by Tony Hain, and it is followed by a response from Geoff Huston and a roundtable discussion with Tony, Geoff, Fred Baker, and John Klensin. We would also like to hear from our readers on this important topic. Please send your comments to ipj@cisco.com.

As an old-time network and UNIX user, I am a big fan of tools that allow simple terminal access to remote host computers. My “Internet career” started in Norway in 1976, where I used *Telnet* to access machines in California through the ARPANET. Today, I still access remote servers through a simple terminal interface, but Telnet has been replaced by the *Secure Shell* (SSH) *Protocol* for all the obvious security reasons. SSH is used not just for terminal traffic—it also can be configured to provide secure tunnels to a variety of services such as Webpages and file transfers. Ronnie Angello explains the details in our second article.

In order to better serve our readers, we will be conducting an IPJ Reader Survey in the near future. Details will be available on our Website at www.cisco.com/ipj. We appreciate your cooperation in completing the survey.

Finally, let me remind you to visit the IPJ Website and update or renew your subscription.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

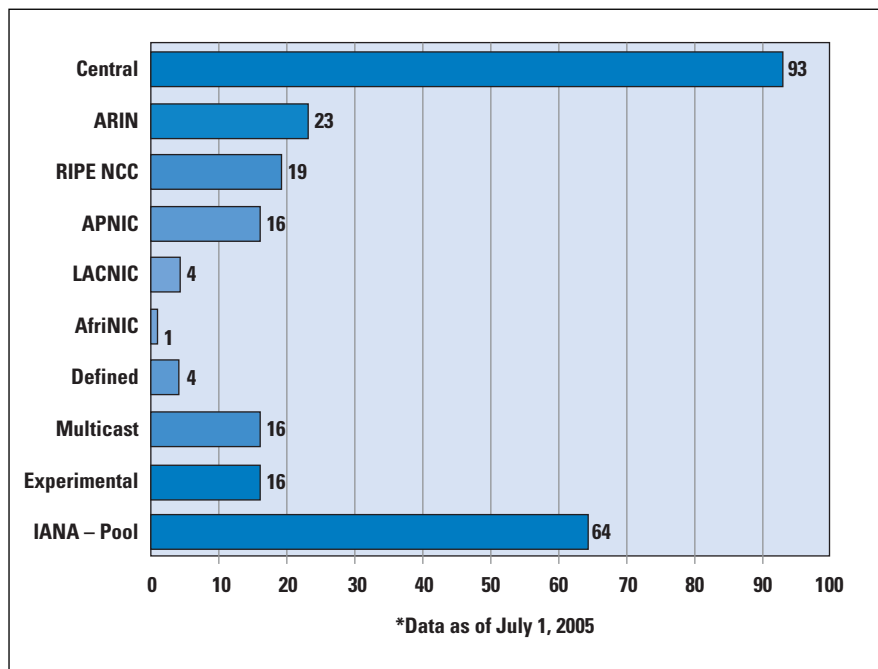
You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

A Pragmatic Report on IPv4 Address Space Consumption

by Tony Hain, Cisco Systems

When I interact with people from all around the world discussing IPv6, there continue to be questions about the projected lifetime for IPv4. This article presents consumption rate and lifetime projections based on publicly available *Internet Assigned Numbers Authority* (IANA) data. In addition, there is discussion about why the widely quoted alternative projection may be flawed, thus leading everyone to believe we have much more time than we might.

Figure 1: IANA /8 Allocations



Allocations

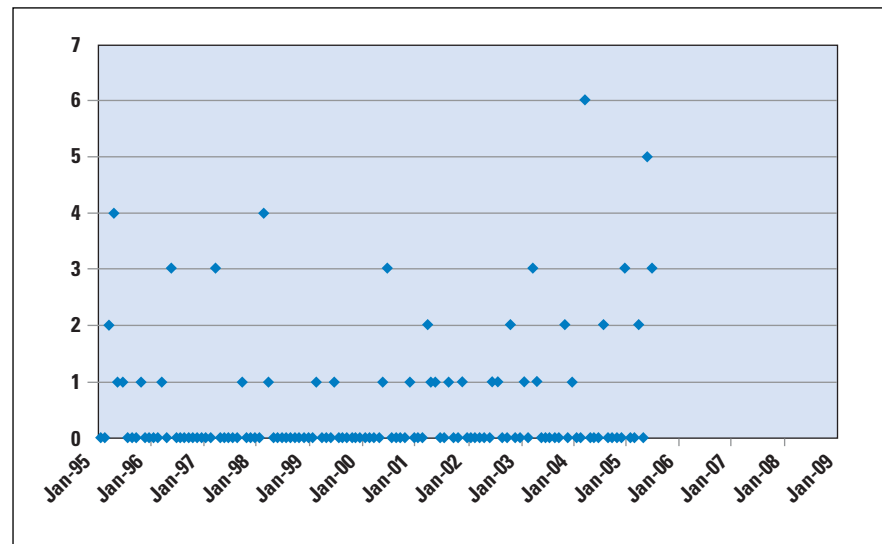
The chart in Figure 1 shows the distribution of all 256 IANA /8 allocation units in IPv4^[1] as of July 1, 2005. The Central registry represents the allocations made prior to the formation of the *Regional Internet Registries* (RIRs). ARIN (North America)^[2], RIPE NCC (Europe)^[3], APNIC (Asia/Pacific)^[4], LACNIC (Latin America)^[5], and AfriNIC (Africa)^[6] are the organizations managing registrations for each of their respective regions. RFC 3330^[7] discusses the state of the Defined and Multicast address blocks. The Experimental block (also known as *Class E*—RFC 1700^[8]) was reserved, and many widely deployed IPv4 stacks considered its use to be a configuration error. The bottom bar shows the remaining useful global IPv4 pool. To be clear, when the IANA pool is exhausted there will still be space in each of the RIR pools, but by current policy^[9] that space is expected to be only enough to last each RIR between 12 and 18 months.

The projection published at <http://bgp.potaroo.net/ipv4>^[10] is often quoted as the definitive reference for IPv4 consumption. This report presents a viewpoint consistent with that author’s long-standing position that we do not need to change from IPv4 to IPv6 anytime soon, thus showing an extended lifetime for IPv4.

The approach used in the potaroo report is to take the simple exponential fit to the allocation data since 1995. As discussed later in this article, this approach includes the effects of the policy shift to *Classless Interdomain Routing* (CIDR) and subsequent digestion of prior allocations, the lull in IANA allocations to the RIRs for two full years, as well as the fact that the model used does not generate a particularly close fit to the actual run rate over the 10-year period.

Although this author agrees that over very long timeframes (20–50 years) there will be substantial variations in the consumption rate for any number of reasons, the opportunity for events that would reduce the recent rate in the timeframe of the remaining IANA IPv4 pool is not evident. That said, there are numerous things that could increase the consumption rate and exhaust the pool even sooner than this projection.

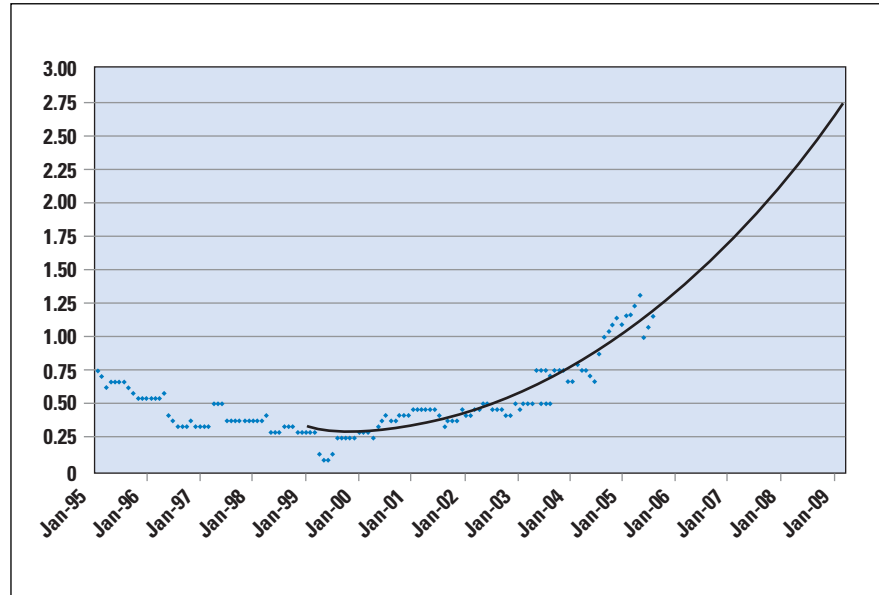
Figure 2: IANA Allocations to RIRs —
Raw /8 Allocations per Month



The graph in Figure 2 shows the raw per-month IANA allocations since 1995. In raw form it is difficult to discern the trend, or develop an expectation about the overall lifetime of the remaining pool.

Taking a closer look at Figure 3, smoothing the data with a 24-month sliding window (averaging over 12 months back and 12 months forward) exposes the underlying reality that the combined rate and quantity of /8 allocations has been steadily accelerating since 2000 (the graphs for 12-, 18-, and 24-month sliding windows show the same fundamental trend). Though a few of the allocations may arguably have been “one-time” events, those are lost as statistically insignificant in the extended and continuing overall growth rate.

Figure 3: IANA Allocations to RIRs —
Sliding-Window 24-Month Average



Taken by itself, the most recent allocation rate (22 /8s over the 18 months leading up to July 1, 2005) suggests that the remaining pool of 64 /8s will be exhausted in about 5 years, even if growth abruptly flattens out to hold around 1 /8 per month. Unfortunately at this point there is no reason to believe the allocation rates will slow or that they will turn downward again. All the gain of CIDR absorbing the pre-1995 allocations has already been incorporated, and there is no obvious economic bubble that might burst to lower demand within the time window of the remaining pool.

To the contrary, the following URL shows potential demand (to bring developing countries up to just 20-percent connectivity, which is half of what the existing Internet world enjoys today) that will swamp the remaining pool, even in the face of much stricter allocation policies.

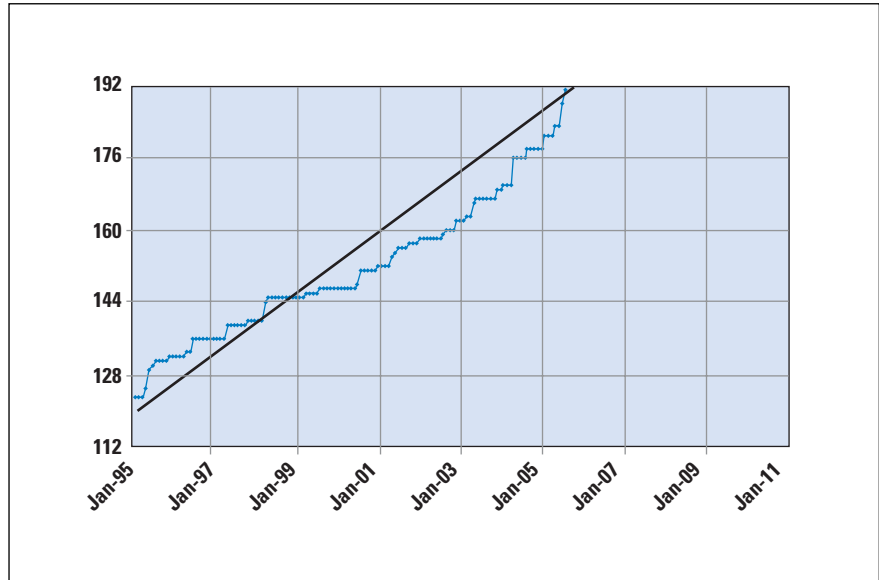
<http://www.nav6tf.org/documents/e-Nations-data.pdf>

So this view of the sustained trend in allocation growth rate suggests that the lifetime of the remaining central IPv4 pool is 4 years +/-1.

Projections

Differing from recent articles and section 5 of the report at **<http://bgp.potaroo.net/ipv4>** that hint at linearity in growth, Figure 4 shows that the raw data after 1995 is clearly nonlinear. It starts with a decelerating rate through mid-1998 as the pre-1995 allocations were absorbed (precipitated by the allocation policy shift from class-based to CIDR), followed by a 2-year lull (only 1 /8 per year), then a return to accelerating growth from mid-2000 onward.

Figure 4: IPv4 Lifetime Projection —
Non-Linear Nature of Raw Data



This suggests that using the past 10-year IANA data is likely to skew the projection toward a much longer period than the recent allocation data would support. Although a longer lifetime projection helps to avoid short-term panic, it can mislead people into believing there is substantial time to worry about this later, resulting in a much bigger problem when reality blindsides everyone sooner than they expected.

Figure 5: IPv4 Lifetime Projections —
Order-N Polynomials, Post-2000
History Basis

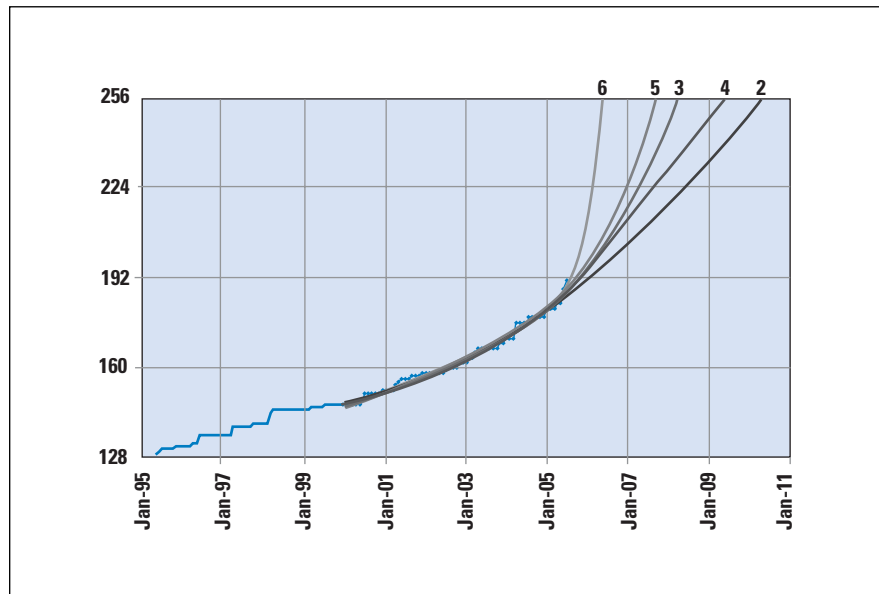
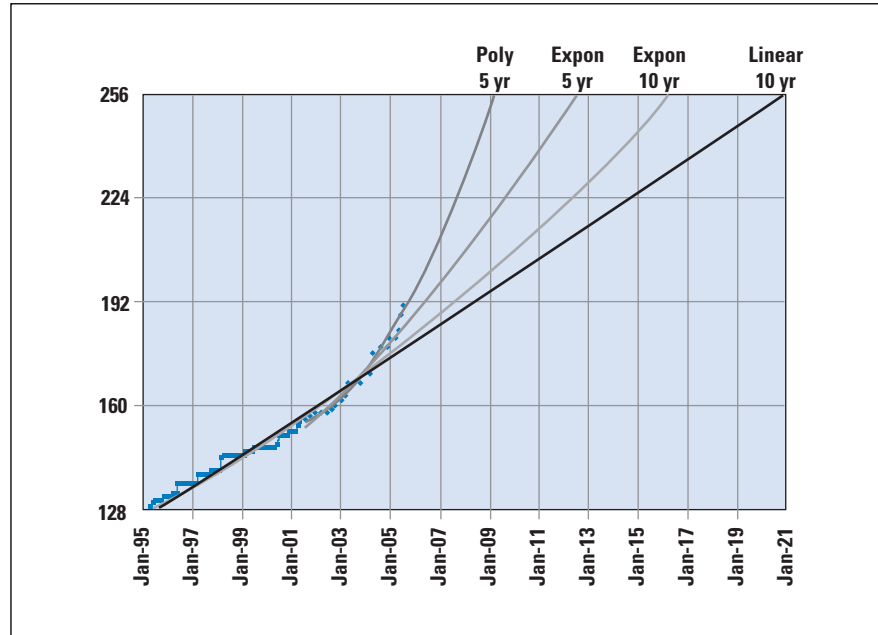
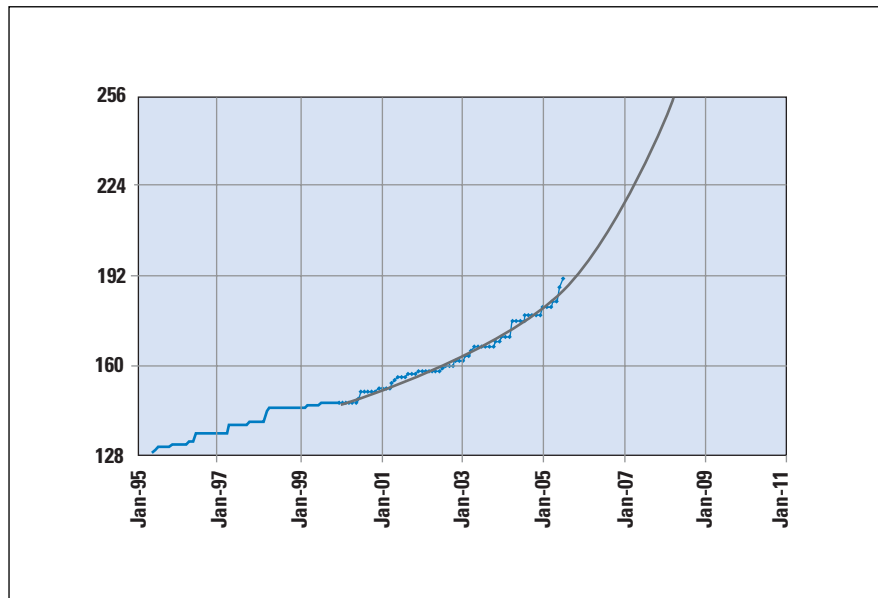


Figure 6: IPv4 Lifetime Projections —
Polynomials and Exponentials

As in any statistical endeavor there are many ways to evaluate the data. The various projections in Figures 5 and 6 show different mathematical models applied to the same raw data. Depending on the model chosen, the nonlinear historical trends in Figure 6 covering the last 5- and 10-year data show that the remaining 64 /8s will be allocated somewhere between 2009 and 2016, with no change in policy or demand (though as discussed previously there are already reasons to err toward 5-year-based nonlinear models).

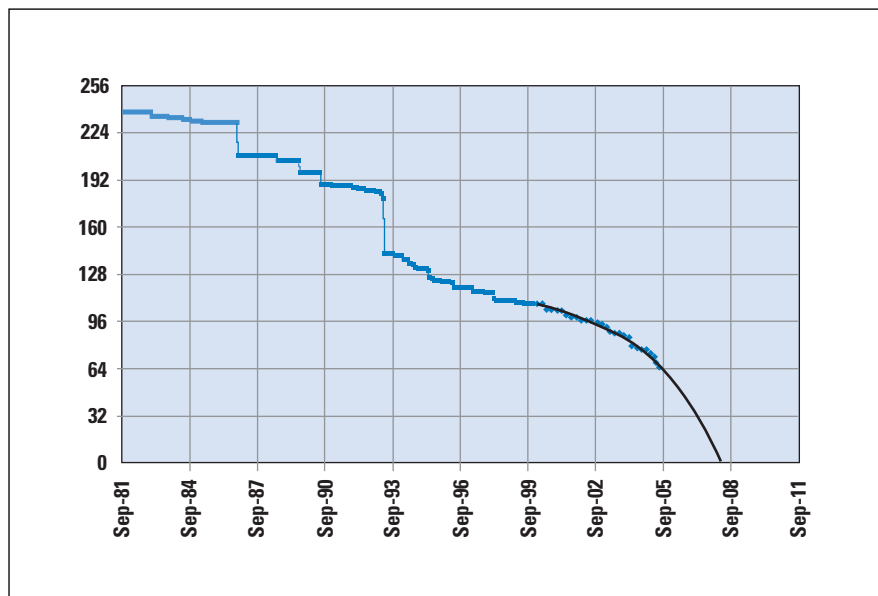
Adding to that, policy is continually changing. ARIN, for example, has recently clarified its policy allowing organizations that demonstrate they have exceeded the capacity of the private space defined in RFC 1918 to acquire IPv4 address blocks from the remaining public pool, even when it is clear these allocations will never be announced to the global Internet. The other regions already have similar policies or are likely to follow suit because the most vocal members of the RIR community have adamantly commented against expanding the private IPv4 range. This policy approach coupled with persistent demand means the actual run rate is going to continue increasing as the large organizations begin consuming public space where they had been using private to support their network growth. For example, one large enterprise has steady growth over 1 percent per month, which currently requires an efficiently managed /12 per year for its expanding network. The enterprise is less than a year from exhausting all the space provided in RFC 1918, so it was very interested in the ARIN policy that allows the enterprise to continue growing through public space. Additionally, multiple commercial service providers expect to reach the capacity of the 1918 space within 12 to 18 months, just supporting management addresses on their existing devices. This does not take into consideration their pending deployment of new services, which they expect will use several new IPv4 addresses per device with marketing targets measured in multiple millions of units.

Figure 7: IPv4 Lifetime Projection —
5-Year History Basis



The graph in Figure 7 hints at the likely outcome as word spreads about the perception of policy liberalization and the demonstrable exhaustion of the remaining global IPv4 pool landing within the *return-on-investment* (ROI) period for new equipment. It is based on the same raw historical data as the frequently quoted long-term projection on potaroo's Figure 2.4, but the more aggressive fit on the most recent data set describes a significantly higher consumption rate and shorter lifetime for the remaining pool.

Figure 8: IPv4 /8 Pool —
5-Year History-Based Projection



The graph in Figure 8 provides the exhaustion perspective, showing the entire address pool from the publication of IP Version 4^[11] (note that data prior to 1995 is accurate as to where it was allocated, but with very coarse granularity as to exactly when). The projection curve is based on the IANA allocations from January 2000 onward.

Only time will tell which projection is correct, but it will already take a fairly significant stalling event to slow consumption and put the actual allocation curve back on the extended track in potaroo's Figure 2.4.

Reserved Space

There are occasionally arguments that the 16 /8s reserved in the experimental space could be used. Although this is likely to be possible for some IP stack implementations, for others it is not. At a minimum, some quick tests show that Windows 95 through Windows 2003 Server systems consider that block to be a configuration error and refuse to accept it. The operational ability to restrict the space to a select stack implementation is limited, and the amount of space there does not really help even if deployment and operations were trivial. Assuming the sustained growth trend in allocations continues, by the time the remaining 64 /8s in the IANA pool are finished the rate would be approaching 3 /8 allocations per month, so the entirety of the old Class E space would amount to about 6 months of run rate.

Reclaiming Allocations

Another debate occasionally resurfaces about reclaiming some of the early allocations to further extend the lifetime of IPv4. Hopefully this article has shown that the ROI for that approach is going to be extremely low. Discussions around the Internet community show there is an expectation that it will take several years of substantive negotiation (in multiple court systems around the globe) to retrieve any /8s. Then following that effort and expense, the likelihood of even getting back more than a few /8 blocks is very low. Following the allocation growth trend, after several years of litigation the result is likely to be just a few months of additional resource added to the pool—and possibly not even a whole month. All this assumes IANA does not completely run out before getting any back, because running out would result in pent-up demand that could immediately exhaust any returns.

Summary

Network Address Translation (NAT) and CIDR did their jobs and bought the 10 years needed to get IPv6 standards and products developed. Now is the time to recognize the end to sustainable growth of the IPv4-based Internet has arrived and that it is time to move on. IPv6 is ready as the successor, so the gating issue is attitude. When CIOs make firm decisions to deploy IPv6, the process is fairly straightforward. Staff will need to be trained, management tools will need to be enhanced, routers and operating systems will need to be updated, and IPv6-enabled versions of applications will need to be deployed. All these steps will take time—in many cases multiple years. The point of this article has been to show that the recent consumption rates of IPv4 will not be sustainable from the central pool beyond this decade, so organizations would be wise to start the process of planning for an IPv6 deployment now. Those who delay may find that the IANA pool for IPv4 has run dry before they have completed their move to IPv6. Although that may not be a problem for most, organizations that need to acquire additional IPv4 space to continue growing during the transition could be out of luck.

References

- [1] <http://www.iana.org/assignments/ipv4-address-space>
- [2] <http://www.arin.net/>
- [3] <http://www.ripe.net/>
- [4] <http://www.apnic.net/>
- [5] <http://www.lacnic.net/>
- [6] <http://www.afrinic.net/>
- [7] <http://www.rfc-editor.org/rfc/rfc3330.txt>
- [8] <http://www.rfc-editor.org/rfc/rfc1700.txt>
- [9] <http://www.rfc-editor.org/rfc/rfc2050.txt>
- [10] <http://bgp.potaroo.net/ipv4>
- [11] <http://www.rfc-editor.org/rfc/rfc791.txt>
- [12] Geoff Huston, “The Myth of IPv6,” *The Internet Protocol Journal*, Volume 6, No. 2, June 2003.
- [13] Geoff Huston, “IPv4: How long do we have?,” *The Internet Protocol Journal*, Volume 6, No. 4, December 2003.

Another Perspective

Ed.: We asked Geoff Huston to provide some feedback on this article and he responded with the following:

Dear Editor,

There are, of course, many ways to undertake predictions, and over the millennia humanity has explored a wide diversity of them. In every case the challenge is to make predictions that end up being closely correlated to the unfolding story, and of course hindsight is always the harshest judge of such predictions.

Tony’s work takes a different base point for making the projection from earlier work that I did in this area. Tony looks at the rate of allocation from the IANA to the RIRs, and bases his predictions on the trends visible in that time series of data. By contrast, I used the assumption that assigned addresses are destined for use in the public IPv4 Internet, and I used the trends visible in the amount of advertised address space as the basis for the predictions of consumption.

One of the more interesting data artifacts is the first-order differential of the rate at which the span of addresses announced in the IPv4 public Internet has increased over time.

(Figure 4.4 of <http://bgp.potaroo.net/ipv4/>)

One interpretation of this data is that there are two phases of recent activity: prior to March 2003 and post-March 2003. Prior to March 2003 the longer-term address growth rate was the equivalent of some 3.5 /8 blocks per year.

Post-March 2003 we see a different consumption growth rate, fluctuating between 5 and 8 /8s per year, with a mean value of some 7.5 /8s per year. There is no strongly obvious longer-term compound growth rate visible in this view of the data. Given some 64 /8s remaining in the IANA pool as of July 2005 and a base consumption rate of a mean of 7.5 /8s per year, the simple division yields 8.5 years, or 2014 as the time of forecast exhaustion of the IANA address pool. At that point the RIRs will be holding about 25 /8 blocks in their unallocated pools, and a further two years of allocations could be made from these pools.

So I would offer the view that the post-2003 data offers a perspective of exhaustion of the unallocated address pools in 2016, with the caveat that such a prediction assumes that the current address demand levels will continue, the actions of industry players are invariant, and the current address allocation policies will continue as they are at present.

Of course these three caveats represent relatively major assumptions about the future—and are perhaps unlikely to happen. It is likely that there will be changes in all these factors in the coming years, and these will obviously impact these predictive models.

To summarize, I observe that these different predictive approaches yield slightly different outcomes, but not beyond any reasonable error margin for predictions of this nature. Sometime in the forthcoming 5 to 10 years the current address distribution policy framework for IPv4 will no longer be sustainable for the current industry address consumption model because of effective exhaustion of the unallocated address pool.

When looking at this prediction from the perspective of the service provider enterprise, the prediction can be re-expressed as a problem relating to investment lifecycles. The ISP industry and the enterprise sector have already made considerable investments in IPv4-based infrastructure in equipment, infrastructure, and operational capability, and we are seeing some considerable reluctance to add to this with additional investment into IPv6 capability at this time. The direction of the use of various forms of NAT-based approaches and increasing use of application layer gateways in the public and enterprise environments can be seen as an effort to extend the lifetime of the existing infrastructure investment. In a volume-based market with relatively low revenue margins, this position certainly has some sound rationale from a business management perspective. But I agree with Tony here that such business approaches are ultimately short-term in nature, because they do not allow IPv4 to encompass indefinite further decades of Internet growth in a silicon-dense world.

However, in terms of understanding the next few years of a process of industry transition of protocol infrastructure into IPv6 deployment, perhaps the real issues here are more centered on competitive business factors and sector investment profiles than they are about detailed introspection of trends within various number series.

The numbers all indicate that this is not a matter that can be deferred indefinitely. Tony's call for some timely attention to the need to commence investment in IPv6-based service infrastructure is one that I hope the industry is listening to attentively.

—Geoff Huston
gih@apnic.net

A Virtual Roundtable

Ole: Let's open this discussion on the point of measurement methods. We invited John Klensin and Fred Baker to join Geoff and Tony in the discussion at our virtual round table. (We often all see each other at IETF meetings, but there is seldom enough time to gather everyone around a real table, hence this discussion took place with a few rounds of e-mail).

Geoff: As I said in my response letter, Tony's work takes a different base point for making the projection from the earlier work that I did in this area. My work has focused on the trends from the addresses used in the public IPv4 Internet, and then deriving projections on consumption based on this data. It assumes that the influencing factor for address consumption is the use of addresses in the public IPv4 Internet.

Tony: As Geoff noted, he and I have discussed over time that we are looking at different parts of the data set and coming to different conclusions. One specific point that distorts the approaches is the time delay between IANA allocation to the RIRs and the appearance of that space for public use. In particular, his comment about 5 to 8 /8s per year is based on the delayed public use data that will eventually catch up with the fact that IANA has allocated 13 /8s just since the beginning of 2005. If the allocation rates had close to linear growth, the delay would not be a big factor. Another point of distortion is the potential for some of the allocations to never show up as publicly routed.

Ole: So when do we actually run out?

Geoff: There are many specific milestones that will pass in sequence. The unallocated address pool held by IANA will exhaust first, and then the RIR pools of unallocated data will drain. At that point there is no stream of "new" addresses to fuel further growth, and that is probably a reasonable point in time to say that we have "run out." Assuming that the current business influential factors and allocation policies remain in place, then the projection models from recent data indicate that this "run-out" date is around 2016, or some 11 years from now. Of course these are unlikely assumptions as the prospect of exhaustion draws nearer, and there may be a "last-minute rush" of address allocation requests from the service provider industry that could draw in that projected "run-out" date. Such additional consumption pressures are difficult to factor in to trend-based predictive models, of course. It is also conceivable that the industry could shift its attention almost entirely to IPv6-based protocol infrastructure in the coming years, in which case the "run-out" projection for IPv4 would extend out further in time simply because of the translation of the consumption activity to the IPv6 address pool.

Tony: As I noted early on in my article, there will still be pool available at each of the RIRs when the IANA pool that I focused on is exhausted. In the past I have said we would never completely run out because nobody could afford that last address, but in light of the accelerating consumption of IPv4 coupled with the less-than-aggressive deployment of IPv6, I can see how the pool might actually run dry.

John: In practical terms, the point at which one has “run out” of address space is not tied to being the last applicant to the RIRs for an address pool. I have suggested that point will never arise: the RIRs (and, to the extent to which the *Internet Corporation for Assigned Names and Numbers* [ICANN] can make decisions, the IANA), will continually recalibrate policies to prevent “running out.” Of course the inevitable consequence of those recalibrations is that, although one does not need to worry about approaching an RIR and being told “no space left,” the combination of monetary, justification, and general aggravation costs is such that one does not even want to contemplate being the applicant for the next-to-last available block. That reasoning says that looking at the date on which near exhaustion is reached is relatively uninteresting. The more important question is when one enters the end game for IPv4 space because, as soon as the end game begins, the space is essentially exhausted.

I suggest that the criterion for entrance into the end game is not measured statistically but by looking at the point at which one needs to start designing networks and subnets, not in a way that is optimal from a network architecture or network management and growth standpoint, but in order to conserve address space and/or to avoid extended discussions with applicable RIRs (or one’s ISP that deals with the RIR). From that point of view, we have already run out, and probably ran out a couple of years ago. Every time someone who has multiple machines is pointed to private address space because of a presumed shortage, it is an indication that we have already run out of space. Every time China manages to make a successful political point—regardless of the country’s actual internal dynamics and economics—about its inability to get addresses for its population, it is an indication that we have already run out of address space. Every time an ISP decides to use private space to manage its backbone, it is an indication that we have already run out of address space.

Fred: I have made the same point, from a point of view of economics. In essence, when a commodity is common and demand is low, there are calls to squander it because it costs nothing—something one hears a lot of in the IPv6 community. When supply and demand are comparable, a market develops, and I need to tell you that I certainly pay for the IPv4 addresses at *my* house. When demand outstrips supply, we enter a regulated market of some kind, and our current allocation policies certainly reflect a regulated market. The step after a regulated market is a black market, and it is not too hard to find that either.

John: Actually, in our present situation, there is an intermediate step before things deteriorate completely into a black market. Although it is unlikely that any significant fraction of the early IPv4 academic, research, or commercial allocations could be recovered and reused, there are governmental allocations that might be recovered under significant political pressures. Unfortunately, in addition to politicizing the allocation process much more than we have seen so far, such moves might push the present users of those allocations toward NATs in ways that would make the ultimate transition to IPv6 more difficult while not gaining very much additional time for the IPv4 space.

Tony: Political pressure or not, simple logistics argues against this. Given the rate of growth in consumption, any reclaimed government space would be consumed in substantially less time than it would take to rebuild their network and release it. Even a small network sitting on a /16 would take at least a year to release that much space, and at the current spot on the escalating curve that /16 represents around 2 hours of IANA run rate. Getting back a whole /8 would logistically take several years, and then at that point on the curve the result would be about a week of run rate. If several of these government organizations have a mesh of direct interactions and head down the same path, the resulting overlap in the private address space would require creating a complex NAT system worthy of a Nobel Prize. Reclamation is a nice bar-room debate topic, but the return on investment is extremely low. If an organization were to consider rebuilding its network to release an IPv4 allocation, it would make much more sense for that organization to rebuild it as IPv6 than to move publicly addressed nodes behind a NAT.

Geoff: It would be strongly preferred by all, I would suggest, that the “black market” option be avoided. If the consequence of the exhaustion of the unallocated pool of IPv4 addresses is the trading of already-allocated IPv4 addresses, then a responsible way for the industry to support that scenario is to encourage such a market to operate with the support of some form of “clear title” that could legitimate trading transactions. Without structure and stability in a trading market, the value of the trade is meaningless, and in this case the potential for chaos in the network itself is undeniable.

Fred: We are in fact starting to see networks designed to be IPv6-only or IPv6-dominant (the latter being a network that might use IPv4 internally but offer only IPv6 services to some or all of its customers) in China, Japan, and other places. The economic argument is the one these operators are primarily giving—they state that they see a roadmap to the number of addresses that they need in IPv6, while in IPv4 they are significantly constrained. This sounds to me a lot like John’s comments about network design, but the other way—rather than designing their networks to what they perceive as IPv4 addressing policy limitations, they are choosing a path that they perceive as giving them options.

We also see evidence of networks designing themselves to the limits of address allocation in IPv4, usually using multiple layers of NATs. For quite a while, for example, China Unicom used multiple layers of NAT in order to work around what the company felt was a deficiency in its ability to get IPv4 addresses from its national registry. As I understand it, the company has changed its strategy to include getting IPv4 address allocations directly from APNIC, and at the same time to deploy an IPv6 network in parallel to move away from IPv4 dependence.

John: There is another factor at work in this. Transitions are never free. If we are going to design and build out a substantially new network, we are rapidly reaching the point—some would say that we have reached it already—at which it is cheaper to design and build that network for IPv6, making whatever arrangements are needed at its interconnection points with IPv4 networks, than to build in IPv4 and face a transition later. As those decisions are increasingly made, it may both reduce pressure on new IPv4 allocations and create free pools of IPv4 space that could be recovered and reused. For example, the U.S. Department of Defense (DoD) has announced a fairly aggressive schedule for moving to IPv6. If they meet that schedule and were then willing to free up the IPv4 space that they would presumably no longer be using, it would free up the equivalent of several /8s. While I agree with Tony that this hypothetical case would be unlikely to make any significant difference in the long run, it illustrates another difficulty with trying to make assertions about what is happening by statistical projections alone.

Ole: It is frequently stated that North America is immune to the address exhaustion problem.

Tony: Well despite persistent rumors and press statements to that effect, ARIN continues to consume about 30 percent of the annual allocation from IANA. If the past allocations were sufficient to stave off global exhaustion, why the continued consumption? In any case, when the central pool is exhausted the North American region will be in the same situation as everyone else—unable to expand or acquire new IPv4 addresses.

Geoff: We are seeing growth in Internet-based services in all regions of the industry, including North America. And network growth needs to be fueled by network addresses. We are seeing a combination of a continued demand for further addresses, and the use of various forms of network configurations that attempt to make the most efficient use of already-allocated addresses. There is little data to suggest that any region, including that of North America, is in a position of immunity from these growth-related factors.

Ole: There is widespread opinion that NAT will solve the problems for a long time to come.

Geoff: The ISP industry certainly has made considerable investments there, and many millions of end users today use the Internet behind NAT devices. Given the size of this investment and the factors of inertia in large-scale service markets, it is reasonable to predict that NATs will be around for quite some time. But NATs add cost to network services. If we are talking about a network that is restricted to servicing the communications needs of people, then this is a relatively high-value activity, and the additional costs of the deployment of NATs are being absorbed within the cost base of the network service economy. And for such human activity-based services this may well continue for some time, given the existing levels of industry investment in service infrastructure that includes the use of NATs. Certainly any new application that is adopted by the Internet user population needs to work across a wide variety of NAT configurations. From this perspective it is likely that IPv4 and NATs will continue to be part of the Internet landscape for a long time to come.

But although this approach has the potential to service a portfolio of service markets for some time to come, it cannot service all forms of service markets—not in the future nor even today. It does not solve all the “problems” and certainly does not encompass all the opportunities that the Internet offers. The potential of IPv6 is one that includes an address span designed to match the full potential of the volume-driven silicon industry, both now and in a future that extends out for many decades to come. One likely scenario for IPv6 is in servicing a truly massive device-dense environment. This scenario encompasses far more than services that are primarily directed at human end users. And the associated service market will be more akin to that of a relatively undifferentiated commodity market, where simplicity and low cost are the dominant service provider discriminants. Because of their additional complexity and associated incremental cost, NATs are marginalized in such commodity markets directed at servicing device density, and it is there that the true leverage of the IPv6 address span becomes a major influential factor.

Tony: As Geoff notes, NAT has been widely available and deployed globally over the same timeframe as the recent consumption. Yet the accelerating growth trend continues, consuming to the point where only 25 percent of the total IPv4 space remains available. Although NAT does slow the rate of public address consumption from what it might otherwise be, it creates more problems than it solves. Geoff also raises the economic investment in NAT to date, which is an interesting contrast to many complaints I hear about the cost of deploying IPv6. Most people who look at what it will take to deploy IPv6 in their network are very quick to dismiss this investment in the array of costs associated with NAT. Often they insist on a demonstration of value for the IPv6 investment while at the same time they refuse to allow consideration of removing their development, and ongoing operational support costs for IPv4 NAT.

Although I agree that in the interim overlap period the costs are additive, in the long term staying on the IPv4/NAT path those costs only compound, whereas on the IPv6 path they disappear. The duration of that overlap is somewhat self-controlled as a direct trade-off between the costs for running both protocols in parallel versus the costs associated with aggressively moving the end systems and applications to IPv6.

Ole: Another area frequently discussed on various lists is that the U.S. DoD and Federal Government mandates for service availability in 2008 are just another instance of the *Government OSI Profile* (GOSIP) and that they too will disappear.

Tony: What these discussions miss is that the situation is entirely different now. In the early 1990s the U.S. GOSIP effort was directed by a strong desire to consolidate the array of protocols in use at that time toward a common one. Other governments had similar efforts that led them collectively toward a suite that was developed with international governmental input. IPv4 was an alternative to the mandate with applications already supporting it, while the OSI protocols existed in some router products but did not have many applications available.

At this point the existing government networks are already consolidated, and there is no alternative. Yes, IPv6 still has fledgling application support, but the IPv4 pool is no longer a sustainable resource to draw on, and there is no other option. So the government networks either stop growing or, as the U.S. DoD and Government agencies have announced, they will move to IPv6. This implies preparing the application community to meet the impending reality.

Geoff: Although the strategic directions of one single—but relatively large—market player does have some bearing on the direction of the global market in Internet-based service provision, I do not see evidence that this will be sufficient to influence the entire market in any particular direction. This was certainly evident in the case of GOSIP some years ago, and continues to be an aspect of the market today. The global communications sector carries the impetus and burden of massive investment in infrastructure, process, technology, services, and consumer product portfolios. The sector has already undergone a revolutionary change with the advent of the Internet over the past decade. Doubtless there is considerable reluctance on the part of many sector players to continue to invest in further change in the protocol infrastructure of Internet-based services. On the other hand, the upheavals in the service provider sector have also eliminated much historical complacency about the stability of these markets and the adequacy of the associated service portfolio. It is reasonable to suggest that this sector is now very attentive to the prospect of expanded markets and new service opportunities that can take advantage of the existing infrastructure to create new revenue streams. So I think it is the current dynamics of the service provider sector and the potential for new service markets that would be the most persuasive factor for service providers to invest in an IPv6 protocol infrastructure.

Ole: Closing thoughts?

Tony: As I said at the end of my article, now is the time to recognize that we have reached the end of sustainable growth in IPv4. For most existing organizations that can foretell they have as much space as they will need for the next decade, this is not really an internal problem. Where these organizations will have a concern is when they deal with newcomers or others that have been forced into IPv6 because of exhaustion of the pool. Those organizations that foresee expansion and growth should evaluate Geoff's analysis as well as mine and weigh their plans against the risks of either or both of us being wrong.

In any case it only makes sense to start IPv6 capability discussions with the product vendors now. Product development cycles can be lengthy, and the only way for the vendor community to mesh with an organization's deployment plans is to have sufficient notice about those plans and timeframes. It would also be wise for the organization's network architects to start thinking about the impacts of an IPv6 deployment. Both protocol versions are packet-based and the names start with IP, but there are enough differences in the details that it is worth taking a fresh look to see what might be easier or cheaper than just blindly deploying IPv6 identically to the IPv4 deployment.

Geoff: The Internet continues to present challenges to the communications sector, and I would suggest that the underlying influential factor is the combination of the silicon and software industries that continue to fuel the demand side with fascinating, innovative, and compelling uses of communications that continue to surprise us with their continual re-statement of the size of the domain in which we operate. We appear to be moving beyond servicing devices that are activated and influenced primarily by direct human activity, such as e-mail and Web use, and we are now looking at various command, control, and monitoring functions that embed themselves deeply in other devices and in other elements of our infrastructure. This encompasses larger concepts such as "smart buildings" and "smart traffic control," and they reach all the way down to the level of embedding into consumer devices and even identification tags. This is not a world that can readily be serviced by an IPv4 protocol infrastructure, and we are already seeing various levels of network indirection in both NATs and various forms of overlay networks to attempt to compress this new scale of basic network addressing demands into the IPv4 environment. This appears to be a complex, and therefore costly task. But the expectation here is that the service industry is heading toward a commodity utility function, where the essential attributes of the underlying network are simplicity and efficiency. These factors suggest that the market characteristics that arise from the propulsion of the silicon and software industries are inexorably tugging the communications service industry to embrace simple, scalable, and efficient networking technologies. It is in this space that the essential attribute of IPv6, that of the size of the address pool, has its most effective leverage. Here the "run out" of IPv4 will inevitably focus our common attention on how best to engage with future needs and roles. And in this perspective the IPv6 technology has a critical and central role.

John: Tony, I think we need to assume that, when it comes down to translating the projections into an answer to the “when do we need to get serious about IPv6?” question, both you and Geoff are, to a considerable extent, wrong. Geoff’s articles and projections have been interpreted by some people as containing a “there is no problem, we can continue with IPv4 until we all retire” message. Viewed from that direction, yours can be seen as “we cannot be quite *that* complacent.” Instead, I think we should all be looking at going directly to IPv6 in newer network installations rather than concentrating on whether we can get enough IPv4 space for them. We also need to be examining—now, not a few years in some projected future—the applications and services for end networks and end users, not just backbone and ISP services and operations. One of my particular concerns is that we have enterprise and customer support people and protocols all over the world who are used to thinking about things in an IPv4 world, including the support advantages of “all NAT-based end networks look the same” architectures. The need to retrain them to think about things differently, and to design and build new tools for their use, may suggest a more time-consuming and expensive transition than changing over the networks themselves.

Fred: What is clear to me from this discussion, Geoff’s prior analysis, and Tony’s analysis here, is that there is a timeline. We are *not* debating whether IPv4 address availability is limited or whether it can be “saved” by address allocation policy, nor are we debating the economic or technical impacts of more or less draconian allocation policies. We *are* debating what constitutes the end game, when and why that end game will become important, and whether perhaps we are already seeing the first steps of it. We are also not debating whether perhaps some new architecture would be preferred over the one in IPv6; if we had an alternative on the table today we could discuss that, but experience tells us that the proposals being considered by the *National Science Foundation* (NSF) and others are sufficiently “researchy” to not be ready for wide-scale deployment in the necessary timeframe.

As such, from my perspective, there is a present call to action.

What U.S. DoD and recent congressional hearings have recommended is in keeping with the IETF’s recommendation and with the IPv6 address allocation strategies of the RIRs. The simplest transition strategy involves presently procuring equipment, operating systems, and applications that are IPv6-capable in preference to systems that are limited to IPv4. At some point in the future, perhaps in the 2008–2010 timeframe, we should plan to turn on IPv6 networking capabilities throughout our networks, and this means gaining experience with IPv6 on a smaller scale in 2005–2007 in our networks, in server applications, and in user systems. Turning down IPv4 capabilities, which is the endpoint of such a transition, is a business decision that does not need to be made hastily; we should presume that coexistence will be important for a decade, and probably more.

Ole: Thank you, gentlemen!

TONY HAIN is currently the Senior Technical Leader, IPv6 technologies, with Cisco Systems. In addition to providing guidance to the various internal product teams, he was also co-chair of the IETF working group developing IPv6 transition tools. His IETF participation since 1987 includes a term on the Internet Architecture Board from 1997 to 2001. Named an *IPv6 Forum Fellow* in 2004, he is currently serving as Technology Director on the forum's North American IPv6 Task Force steering committee. Prior to joining Cisco in 2001, he spent 5 years at Microsoft, where his roles included Program Manager for IPv6 as well as Network Analyst for the CIO's office. Prior to Microsoft, he was the Associate Network Manager for the U.S. Department of Energy's Internet effort, ESnet. With this range of roles, spanning the space between the implementation technologists and senior management, he brings a real-world viewpoint to the deployment decision process. E-mail: ahain@cisco.com

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for the past decade, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector, and has served his time with Telstra, where he was the Chief Scientist in the company's Internet area. Geoff is currently the Internet Research Scientist at the Asia Pacific Network Information Centre (APNIC). He served as a member of the Internet Architecture Board from 1999 until 2005, and currently co-chairs the Site Multi-homing and Routing Operations IETF Working Groups. He is author of *The ISP Survival Guide*, ISBN 0-471-31499-4, *Internet Performance Survival Guide: QoS Strategies for Multiservice Networks*, ISBN 0471-378089, and co-author of *Quality of Service: Delivering QoS on the Internet and in Corporate Networks*, ISBN 0-471-24358-2, a collaboration with Paul Ferguson. All three books are published by John Wiley & Sons. E-mail: gih@apnic.net

JOHN KLENSIN is an independent consultant based in Cambridge, Massachusetts. He has been involved in the design, development, and deployment of ARPANET and Internet applications, and occasionally lower-layer technologies, since the late 1960s and early 1970s. He has also been intermittently involved with Internet administrative and policy issues since the early 1980s. His current work primarily focuses on internationalization of the Internet on both technical and policy dimensions. E-mail: klensin@jck.com

FRED BAKER has worked in the data communications industry, building network elements such as switches and routers, since 1978. His involvement with Internet technology started in 1986, and with the IETF in 1989. He has contributed to the development of OSPF, QoS, PPP, SNMP MIBs, and a variety of other technologies. He has also held a variety of management positions, including chairing various working groups, participating in the IAB, and chairing the IETF. He currently serves on the Technical Advisory Board of the U.S. Federal Communications Commission and as the Chairman of ISOC's Board of Trustees. E-mail: fred@cisco.com

Practical Uses of SSH Tunneling in the Internetwork

by Ronnie Angello

While the growing popularity of broadband Internet services and elevated concerns with securing *Wireless LANs* (WLANs) have become major concerns for network administrators today, *Secure Shell* (SSH) *Protocol* tunneling has proven to be a secure and effective solution for addressing various needs and concerns of both network users and administrators. Making the transition from traditional dialup remote access to a broadband solution can bring along with it some roadblocks when trying to preserve functions and security. WLANs can be difficult to secure in the enterprise, mainly because of the various client types that must connect to the network. SSH tunneling can help alleviate both of these issues.

SSH tunneling, also known as SSH *port forwarding*, is the process of forwarding selected TCP ports through an authenticated and encrypted tunnel. These tunnels can be constrained to within two points of the company's enterprise network, or it can originate on a small office or home office (SOHO) computer on a given provider's network, and transit the Internet to a server on the enterprise network. Some practical uses for SSH tunneling are outlined in this article.

A Look Back at Traditional Remote Access

Remote access is the method of connecting from a SOHO computer that resides on a remote foreign network, or has no permanent network connection, to the enterprise network or central office. Usually this involves traversing the Internet. This can be for the purpose of telecommuting, providing on-call support from home, checking e-mail while away from the office, or for the old-fashioned workaholic who must work from home. Remote access used to involve simply accessing a network through an analog phone line or possibly ISDN. In either case, the user was authenticated by an access server that resides on the enterprise network and given authorization to certain resources.

When connected to the access server, users had the feel of being connected to their company's enterprise network. They were free to browse internal Web pages and access various Windows domain resources. They could connect to the network neighborhood and transfer files to and from the work computer. They could connect directly to internal UNIX servers with SSH and use a local X-server application to access UNIX applications from the SOHO.

PC remote-control applications such as VNC, etc. could be used to access files and applications that reside on a host computer on the enterprise network without extensive configuration on the home PC. In addition to the ease of configuration for the administrator or user, fewer applications need to be installed on the home computer to accomplish work tasks from home. This approach saves software licenses in addition to valuable company resources.

Most network administrators cannot let PC configuration consume a great deal of their time because they are busy enough as it is. From a function standpoint, users felt like they were working from their office at work. It was too slow though, so it did not really matter. Then broadband services were introduced, and they offer high bandwidth, but getting the same functions is a bit more challenging. Users benefit from the extra added bandwidth, but of course the administrator has to make sure that everything works as if nothing ever changed.

Broadband Services Emerge

Many users are now migrating from their traditional dialup connections for Internet access to a technology that offers more bandwidth such as cable or DSL. Broadband wireless services are now emerging in some areas as well. These services may even be cheaper than what the company or individual was previously paying for ISDN service, and it is “always on.” Most users are no longer dialing a company access server to access the resources that are vital to their job. They are now permanently connected to a foreign provider’s network, and often the only choice for secure remote access to the enterprise is through a VPN. Strict policies, however, may need to be enforced on the remote SOHO computer for it to be a comfortable solution for security administrators to implement.

For those organizations without the time, money, or manpower to implement and support VPN, Linux login servers can be opened up to the Internet to authenticate users that employ SSH to access the enterprise network from these remote networks. These servers are no more than relay points to access internal systems. They should be placed in the DMZ or on a “screened” network protected by a firewall. The other internal systems are not directly accessible from the remote networks. In cases where remote access is considered a valuable resource to the organization, more than one of these servers should be implemented for load sharing and redundancy.

However, certain functions are lost. Initiating an application from a UNIX computer and displaying it to your SOHO computer with a local X server has been proven to be slow and inadequate from some remote networks. In addition, internal domain PCs and network shares are no longer accessible through the network neighborhood, and file transfer is not available without an additional secure, standalone application. The remote-control applications that access the internal PC will no longer work without opening holes in the firewall. There is a simple solution to all this that is free, secure, and effective: SSH tunneling.

Securing Broadband Remote Access

The functions described in this section can be achieved with any SSH client capable of tunneling, any Web browser that supports HTTP and *Secure Sockets Layer* (SSL) proxies, and any PC remote-control application. The first step is always to connect to the remote login server that has been made accessible to the SOHO user. When connected to this login server, the user can use SSH to access any other internal machine, or take advantage of SSH port forwarding to accomplish their other tasks.

A proxy server may already be configured on your enterprise network. This server is configured to accept connection requests for Web pages and allow the clients to view them with little network overhead. The SSH client on the SOHO computer is configured to forward the specified local source HTTP port (such as 8080) to port 80 on the remote destination HTTP proxy server. It can also be configured to forward the specified local source SSL port (such as 4433) to port 443 on the remote destination SSL proxy server.

The browser on the client machine is configured to use the HTTP or SSL proxy server **localhost** on the specified local port(s). When the browser attempts to download a page, the SSH client forwards the request to the specified remote proxy server on your enterprise network through the established tunnel. Internal Web pages that would normally be available only on the enterprise local intranet are available without latency and without compromising security.

The same concept can be followed for tunneling PC remote-control application data through SSH. The remote-control host service is not changed, and it is waiting for a connection attempt from a remote computer as it normally would. A new remote-control connection is configured on the SOHO computer pointing to **localhost**. Using any additional encryption offered by the remote-control application is possible, but not necessary. Additional encryption will add latency, and SSH provides strong encryption itself with *Triple Digital Encryption Standard* (3DES), Blowfish, etc. The SSH client is configured to forward the local source ports used for the remote-control data (that is, port 3389 for RDP) to destination ports on the host computer on the enterprise network.

Once again, all the functions that the user had when dialing up the enterprise network directly are now available. With SSH, an additional layer of security is provided. Because the desktop of the internal computer is available on the SOHO computer's desktop, users have access to all applications, files, and network resources that they would if they were physically working from their office at work. No additional software applications need to be installed on the office computer to satisfy requirements of working from home, and minimal software needs to be installed on the users' personal home computers. Some of these remote-control applications also provide a file transfer tool that can be used to transfer or synchronize files between the two PCs.

SSH Tunneling for WLAN Security

Securing WLANs has become a monumental problem today for most network administrators. Many organizations are resorting to proprietary solutions or are simply avoiding the implementation of WLANs entirely. An entire article could be dedicated to the importance of securing wireless and the details of accomplishing such a feat.

In addition to the uses described in the previous sections, SSH tunneling can also be used to supplement or replace weaker, more vulnerable encryption found in other network applications. Consider *Wired Equivalent Privacy* (WEP) encryption, for example.

Although other alternatives such as *Wi-Fi Protected Access* (WPA) are available, most WLANs have been implemented with either no encryption or with static WEP only. Static WEP has been highly criticized because of vulnerabilities in the protocol that have been discovered and widely documented. Even when implemented at the 128-bit level, there are tools circulating the Internet that exploit a well-known vulnerability that allows a hacker to crack WEP keys. Even with a WPA solution in place, there will be clients that support only static WEP. These traditional clients can be secured in the meantime by restricting network access with an *Access Control List* (ACL) and tunneling insecure protocols through SSH. Once again, the same functions can be achieved with a VPN solution, but some organizations have neither the money nor resources to implement it.

Summary

In conclusion, SSH tunneling can be used well beyond the scope of the methods explained in this article. The particular uses outlined in the previous sections have been practical in my experience and have been very successful implementations. When users decide to change to a provider that offers broadband, I have found that simply providing a procedure for configuring tunneling has been successful for getting them operational from home.

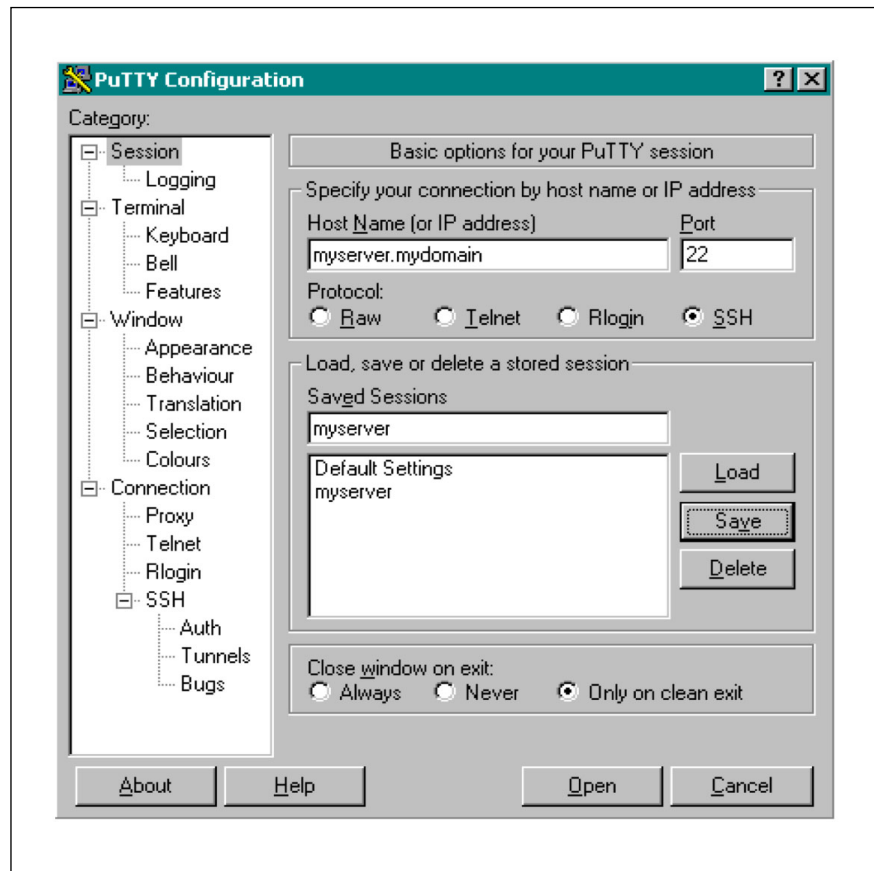
SSH tunneling should be of interest to any organization that wishes to allow its users secure access to all the resources that they may need to accomplish their job functions—especially from a remote location. While exploring possibilities to make a particular application or protocol secure, always consider SSH tunneling an option. SSH provides authentication and encryption that has been proven to be effective for any application.

Securing Remote Access to Internal PCs, Web Pages, etc.

The following is a short example procedure for configuring tunneling for this specific function. It does not include detailed instructions for configuring specific applications, but it outlines the important steps that must be followed in order for it to work properly.

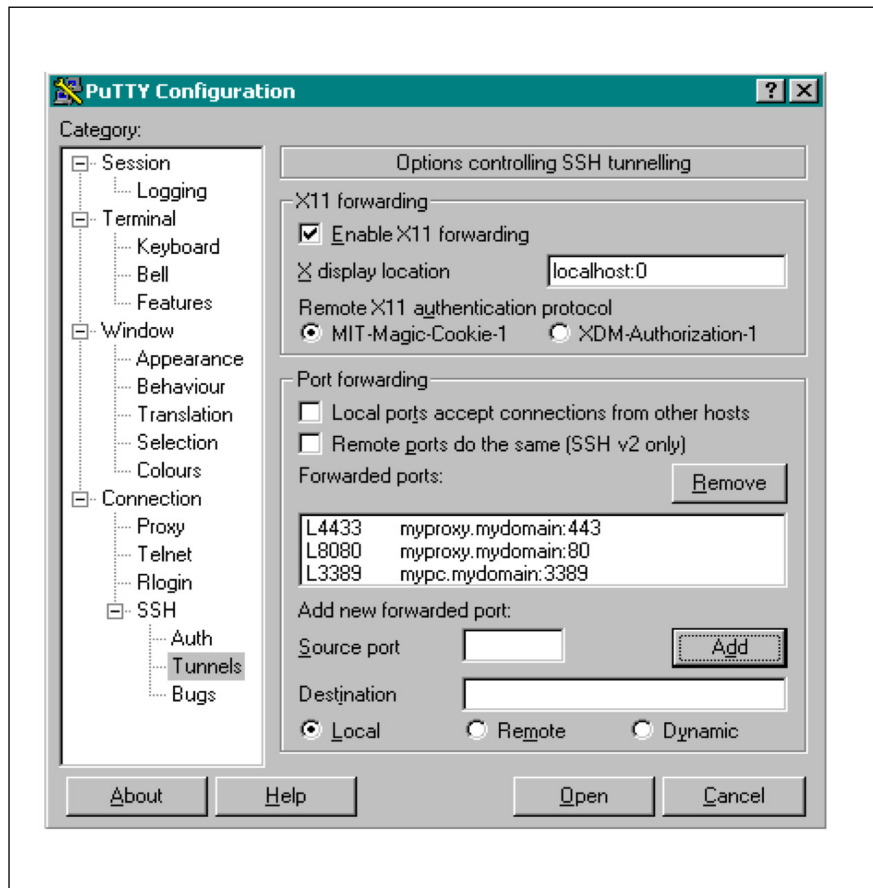
- Any SSH client that supports tunneling can be used. You can download the PuTTY SSH client (**putty.exe**) from:
<http://www.chiark.greenend.org.uk/~sgtatham/putty/download.html>
- Make sure that you select port 22 (SSH). (See Figure 1.)

Figure 1: PuTTY Configuration
Screen — Sessions



- Choose your preferred encryption cipher; enable compression and X forwarding if desirable. Click “tunnels” in the tree menu. Add the local source port(s) and the remote destination port(s) for the ports that you would like to forward through the tunnel. (See Figure 2.)

Figure 2: PuTTY Configuration
Screen —Tunnels



- Make sure that the LAN settings in your Web browser are configured to use the HTTP/SSL proxy server **localhost** on the local port that you specified.
- Make sure that your remote-control connection is pointing to the computer “LOCALHOST.” If you have trouble connecting, make sure that the host service is running on the host PC.

For Further Reading

- [1] The SSH (Secure Shell) Remote Login Protocol, SSH-1 Specification, T. Ylonen, November 1995.
- [2] SSH-2 Specifications IETF Secure Shell working group, June 2003.
- [3] O'Reilly Network Using SSH Tunneling:
<http://www.oreillynet.com/pub/a/wireless/2001/02/23/wep.html>
- [4] SSH Tunneling:
<http://www.ccs.neu.edu/groups/systems/howto/howto-sshtunnel.html>

- [5] SSH Tunnel Tiny HOWTO:
<http://www.frozenblue.net/tools/howtos/?v=ssh-tunnel>
- [6] Secure Email Through SSH Tunneling:
<http://www.slac.com/~mpilone/projects/kde/kmailssh/>
- [7] Mac OS X SSH Tunneling:
<http://info-center.ccit.arizona.edu/~consult/macx-tunnel.html>
- [8] PuTTY Links:
<http://cdot.senecac.on.ca/software/putty/links.html>
- [9] William Stallings, "SSL: Foundation for Web Security," *The Internet Protocol Journal*, Volume 1, No. 1, June 1998.

RONNIE ANGELLO, CCNP, CQS-CWLANSS, CCNA, holds an A.A.S. Degree in Information Systems Technology (Specialization in Operating Systems and Network Operations) and is currently completing degree requirements for the Bachelor of Science Degree in Information Science (Concentration in Networking and Communications) at Christopher Newport University in Newport News, Va. He recently passed the CCIE Routing and Switching Qualification Exam and is preparing for the CCIE Lab Exam. E-mail: **angelo@jlab.org**

Book Review

Network Algorithmics *Network Algorithmics: An Interdisciplinary Approach to Designing Fast Networked Devices*, by George Varghese, ISBN 0120884771, Morgan Kaufmann, 2004.

This is not a generic algorithms book (that is, it does not overlap much at all with Sedgewick or Coleman as an introduction to algorithms), nor is it a typical introduction to TCP/IP networking book (for example, there is no chapter defining the TCP/UDP/IP header fields, thank goodness). It might best be described as an algorithms analysis book set in the context of networking and also in the context of implementations that mix hardware and software solutions. For those familiar with Radia Perlman's book *Interconnections*, I found aspects of the writing style and approach to be similar. George Varghese—in addition to having been a networking professor for many years—has had a lot of industry experience from licensing algorithms to networking companies, to consulting with Procket Networks in the company's early days of architecting its core router, to starting a security company that was recently acquired by Cisco Systems. I have been doing architecture work at Cisco for several years and can say that George's book has real grounding in how systems are built and analyzed today.

Organization

Chapter 2 presents abstractions for networking protocols, hardware design, routers, memory technology, and Internet end nodes (servers). This is a great introduction into "systems" thinking. In section 2.2.7, "Final Hardware Lessons," one thing I thought George should have mentioned along with metrics of chip size, speed, I/O, and memory is *power*. Power is becoming a major systems concern in many platforms and deserves mention as an optimization constraint.

Chapters 3 and 4 go through a list of 15 implementation principles to use in approaching algorithmic design in systems and then give examples of these principles in action. What I find interesting about this section is that from working with George in the past, he really does believe and practice "principle"-based architecture thinking. I remember discussing several of the principles with him several years ago, and you can see how his many years of experience working in the networking field have shaped these principles. Many have probably employed some of these, but as George says in the chapter introduction, having them explicitly documented with examples is useful to help clarify our thinking. Some of the principles (and both the short examples in this chapter as well as examples cited in more detail in later chapters) are really fundamental, and I think reading through examples helped clarify in my mind when to use them.

Chapter 5 covers copying data, for example, in a server design. I really like this type of chapter, in which a subject (in this case the effect of packet copying on Web server performance) is explored in detail but with a focus on where algorithms and systems design play an important part.

My biggest question about this chapter is that I was unsure how applicable this is to, say, modern server design using Linux and with latest Gigabit Ethernet *network-interface-card* (NIC) designs. I know there was a lot of interesting work in the late 1990s, but this chapter without any data is more along the lines of an extended example of how to apply implementation principles.

Chapters 6 through 9 are not what I would consider the meat of the book; they treat the topics of implementation and analysis for servers, timers, parsing/classification of packets, and buffer management (memory allocation).

Chapter 10 covers exact match lookups. There is not a lot of meaty algorithmic discussion, but the history of scaling performance of bridges is used to elegantly show an evolution of algorithmic approaches to exact matching.

Chapter 11 is an awesome overview of the state-of-the-art in longest prefix match (used for destination address matching in routers and switches). A good read of this chapter will yield an understanding of the trade-offs in all major published algorithms, although there may be variations or tuned versions of these algorithms in use at companies like Cisco. I believe this chapter covers all the major categories of solutions.

Chapter 12 extends the prior chapter into more general packet classification (which is used in applications like extended access lists). Like the lookup chapter, this chapter addresses one of George's prime core competencies. There is good discussion on leading published approaches (Grid-of-Trie, cross producting, geometric, and decision tree-based approaches). I strongly recommend this chapter.

Chapters 13 and 14 cover packet switching (that is, architecture of fabrics like crossbars for connecting line cards in a router or switch) and then packet scheduling. These topics get a good academic treatment (after all, George is one who introduced *Modified Deficit Round Robin* (MDRR) to the industry as well as academia), and although there are gaps between what many networking markets are defining as requirements for packet scheduling and what is in this chapter, the chapter is still useful.

Chapter 15 is a short chapter that tries to treat at a high analytic level the algorithmic problems involved with routing protocols. It covers this topic without getting very specific into nonrelevant (to the analysis) networking details.

Chapter 16, which addresses measuring network traffic, was probably one of my least favorite chapters. Some of it is academically interesting but requires network level changes that I just do not think will occur. There are some cute tricks relative to counters and such, but I think they are similar to approaches already being used.

Chapter 17 is a network security chapter and seems to serve as an early introduction to the topic of algorithms in network security; this is not a major focus area of the book.

Areas for Improvement

There is always room for improvement, and I list here three areas in which this book could have been improved:

1. There is a running thread in the book of prefacing technical discussions in some cases with an example from the “normal world,” like comparing packets to envelopes in the postal system. I estimate this is less than 1 percent of the content of the book and fairly easy to ignore if it annoys you.
2. I would have enjoyed better (more detailed) figures. A well-done, detailed figure can incorporate multiple concepts in the text around it and make it much clearer. On the positive side, there are numerous figures in the specifications, even if they do tend to be simple and high level.
3. Another area that I would have enjoyed seeing more on is empirical data (tables of data and graphs). I enjoy detailed empirical data of the type that Hennessy and Patterson so effectively use in their *Computer Architecture* book. There are many places (for example, Web server optimizations in Chapter 5) that I think could have benefited from detailed empirical data. However, I think folks often rely on empirical data too much when a simple analysis like the type done throughout the book could be done to help optimize the problem.

Recommended

Many chapters in this book are directly relevant to the development of networking equipment and software, as well as what is “under the hood” of networking equipment. The book is fun to read and I believe succeeds in trying to convey an organized systems approach to thinking about problems in the networking space.

—Will Eatherton
will@cisco.com

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. Contact us at **ipj@cisco.com** for more information.

Internet Governance Report Available

The *Computer Science and Telecommunications Board* (CSTB) of the National Academies has recently published a report entitled “Signposts in Cyberspace: The Domain Name System and Internet Navigation.”

A summary report, as well as links to the full report can be found at:

<http://www.cstb.org/dns/signpost.html>

From the summary: “The *Domain Name System* (DNS) enables user-friendly alphanumeric names to be assigned to Internet sites. Many of these names have gained economic, social, and political value, leading to conflicts over their ownership—especially names containing trade-marked terms. Congress, in Public Law 105-305, directed the Department of Commerce to request the *National Research Council* (NRC) to perform a study of these issues. When the study was initiated, steps were already underway to address the resolution of domain name conflicts, but the continued rapid expansion of the use of the Internet had raised a number of additional policy and technical issues. Furthermore, it became clear that the introduction of search engines and other tools for Internet navigation was affecting the DNS. Consequently, the study was expanded to include policy and technical issues related to the DNS in the context of Internet navigation. This report presents the NRC’s assessment of the current state and future prospects of the DNS and Internet navigation, and its conclusions and recommendations concerning key technical and policy issues.”

The report was produced by the Committee on Internet Navigation and the Domain Name System: Technical Alternatives and Policy Implications, National Research Council.

First Protocols for Policy Makers Forum to be held October 28

The Internet has achieved the same global economic significance that propelled issues of international trade and finance onto the front pages of newspapers and the forefront of international policy thinking twenty years ago. This change is raising the profile of specialized issues and “obscure” policies for a rapidly expanding circle of public and private-sector stakeholders. Increased general understanding will be vital to assuring that Internet’s growth, development, and coordination mechanisms continue to serve important public interests.

In recognition of this growing need for public education, Packet Clearing House is organizing a series of day-long roundtable fora to encourage sharing of technical and institutional know-how between prominent Internet architects, policy makers, and leading opinion leaders from related sectors. With the support of the *American Registry for Internet Numbers* (ARIN), the forum, to be called *Protocols for Policy Makers* (PfP), will meet for the first time on October 28, in conjunction with the NANOG 35 and ARIN XVI Internet operations and policy meetings in Los Angeles, California.

See **<http://nanog.org/arinnattend.html>**

PfP will explore themes of competition, coordination, and possible conflict between new alternative Internet naming and addressing systems which are challenging the status-quo, such as the national registries recently proposed by the International Telecommunications Union and competitive private-sector “alternate roots.” What outstanding problems are these new mechanisms intended to solve, and what goals might they achieve? How will these innovations contribute to the advancement of Internet public interests? What risks, costs, and complications may be imposed on the Internet by the emergence of multiple divergent systems? At PfP, these issues will be examined through a day of structured round-table discussions, interspersed with comments from leading experts on the Internet’s current naming and addressing systems and prominent advocates of the current restructuring proposals. A complete agenda and list of speakers will be published shortly at <http://www.pch.net>

PfP will be open to the public, but space is very limited. For more information, or to request an invitation, please e-mail pfp@pch.net. Expressions of interest from potential speakers, meeting hosts, and institutional co-sponsors are also welcome. Plans for future PfP meetings are already underway, with a second meeting, tentatively titled “When Voice Goes to Bits” to focus on technical, commercial, and regulatory implications of the migration voice telephony to IP. Suggestions for future meeting themes, venues, and contributions should be directed to PfP Forum Chair Tom Vest at pfp-sponsor@pch.net

Jun Murai Recognized with Postel Award

Professor Jun Murai is this year’s recipient of the Internet Society’s prestigious *Jonathan B. Postel Service Award*. The award recognizes Professor Murai’s vision and pioneering work that helped countless others to spread the Internet across the Asia Pacific region.

The Postel Award was presented during the 63rd meeting of the *Internet Engineering Task Force* (IETF) in Paris, France by Daniel Karrenberg, chair of this year’s Postel Award committee, and Lynn St. Amour, President and CEO of the Internet Society.

“Jun Murai has always encouraged, inspired and helped others, particularly his students and his colleagues in other parts of the Asia Pacific region,” said Karrenberg. “He has also played a key role in creating structures for Internet coordination in the region (particularly the *Asia Pacific Network Information Centre* [APNIC]), and he is widely recognized for his recent pioneering work in IPv6 implementation.”

Jun Murai is currently Vice-President at Keio University in Japan, where he is a Professor in the Faculty of Environmental Information. In 1984, he developed the *Japan University UNIX Network* (JUNET), and in 1988 established the WIDE Project (a Japanese Internet research consortium) of which he continues to serve as the General Chairperson. He is President of the *Japan Network Information Center* (JPNIC), a former member of the Board of Trustees of the Internet Society and a former member of ICANN’s Board of Directors.

The Jonathan B. Postel Service Award was established by the *Internet Society* (ISOC) to honor those who have made outstanding contributions in service to the data communications community. The award is focused on sustained and substantial technical contributions, service to the community, and leadership. With respect to leadership, the nominating committee places particular emphasis on candidates who have supported and enabled others in addition to their own specific actions.

The award is named after Dr. Jonathan B. Postel, who embodied all of these qualities during his extraordinary stewardship over the course of a thirty-year career in networking. He served as the editor of the RFC series of notes from its inception in 1969, until 1998. He also served as the ARPANET “Numbers Czar” and the *Internet Assigned Numbers Authority* (IANA) over the same period of time. He was a founding member of the *Internet Architecture Board* (IAB) and the first individual member of ISOC, where he also served as a trustee.

Previous recipients of the Postel Award include Jon himself (posthumously and accepted by his mother), Scott Bradner, Daniel Karrenberg, Stephen Wolff, Peter Kirstein and Phill Gross. The award consists of an engraved crystal globe and \$20,000.

ISOC is a not-for-profit membership organization founded in 1992 to provide leadership in Internet-related standards, education, and policy. With offices in Washington, DC, and Geneva, Switzerland, it is dedicated to ensuring the open development, evolution and use of the Internet for the benefit of people throughout the world. ISOC is the organizational home of the IETF and other Internet-related bodies who together play a critical role in ensuring that the Internet develops in a stable and open manner. For over 13 years ISOC has run international network training programs for developing countries and these have played a vital role in setting up the Internet connections and networks in virtually every country connecting to the Internet during this time. For more information visit: <http://www.isoc.org>

Internet Root Servers Deployed in India

APNIC recently announced that three new Internet DNS root name servers are now operational in India.

These servers, launched in an official ceremony in New Dehli, India, on 25 August 2005, are the first root name servers deployed in India and South Asia and are already bringing significant improvements in speed and reliability to Internet users in India and the surrounding region.

APNIC has coordinated these deployments with the *Department of Information Technology* (DIT) and the respective root server operators.

F-root, operated by *Internet Software Consortium* (ISC) has been installed in Chennai; I-root, operated by Autonomica, has been installed in Mumbai; and K-root, operated by RIPE NCC, has been installed in Noida, near Delhi.

The installation of the root servers in India has been made possible by DIT, the *National Internet Exchange of India* (NIXI), and the *Internet Service Provider Association of India* (ISPAI), with financial and logistical support from APNIC. The three deployments in India bring the total number of root DNS servers in the Asia Pacific region to 24, 16 of which have been made possible with APNIC's support.

"We are pleased that India is able to contribute to the deployment of the first root name servers in South Asia," said Mr Pankaj Agrawala, Joint Secretary of DIT. "These three root servers will not only benefit the Indian Internet community, but also Internet communities in the surrounding region."

Paul Wilson, Director General of APNIC, added, "The deployment of these three root name servers in India is a positive example of Internet community coordination. The installation has involved the private sector, not-for-profit organizations, and government bodies working together to improve DNS stability and Internet response times for developing countries in South Asia."

Amitabh Singhal, Acting CEO of NIXI, said, "India is among the top ten countries in Internet usage, with over 35 million current subscribers and a five year target for 40 million, translating into more than 200 million total users by 2010. Sustainable infrastructure capacity building is imperative. As a budding intellectual capital of the world, with conducive socio-economic and political environments, India is justifiably proud of hosting three root servers, visibly putting our country, as well as the South Asian region, firmly on the world Internet route map."

More information about the participants can be found below.

- APNIC is one of five Regional Internet Registries currently operating in the world. It provides allocation and registration services which support the operation of the Internet globally.
<http://www.apnic.net>
- Autonomica AB is responsible for **[i.root-servers.net](http://www.i.root-servers.net)**, the first root name server to be installed outside the United States of America. **[i.root-servers.net](http://www.i.root-servers.net)** has been operational since 1991 and is now anycast from more than 25 locations around the Internet.
<http://www.autonomica.se>
- DIT operates under the Ministry of Communications and Information Technology, *Government of India* (GOI).
<http://www.mit.gov.in>
- ISC operates one of the 13 root DNS servers as a public service to the Internet. ISC has operated F-root for the IANA since 1993.
<http://www.isc.org>
- NIXI is joint effort between the GOI and the ISP industry to localize Internet traffic in India. NIXI has nodes in Delhi, Mumbai, Chennai and Kolkatta. **<http://www.nixi.in>**
- The RIPE NCC is one of five Regional Internet Registries currently operating in the world. It provides allocation and registration services which support the operation of the Internet globally.
<http://www.ripe.net>

IETF Journal Announced

The Internet Society (ISOC) is pleased to announce the *IETF Journal*, a new publication produced in cooperation with the IETF Edu team. Our aim is to provide an easily understandable overview of what is happening in the world of Internet standards, with a particular focus on the activities of the IETF *Working Groups* (WGs). Each issue of the journal will highlight some of the hot issues being discussed in IETF meetings and in the IETF mailing lists.

The focus of this first issue will be a look back at the accomplishments of the recent 63rd meeting of the IETF in Paris.

We trust that this publication will give all those with an interest in the increasingly important Internet standards development process an opportunity to keep abreast of many of the topics being debated by the IETF. Articles will cover issues such as:

- Reports from the IETF and IAB Chair
- News from the IETF Edu Team
- Update from the IASA and the IAD
- Summary of the plenary discussions
- Highlights of IETF developments related to topics such as Routing, DNS, and IPv6
- Recently published RFCs.

The journal will be available shortly at the following URL:

<http://www.isoc.org/pubs/IETF-Journal>

Upcoming Events

The *North American Network Operators' Group* (NANOG) will meet in Los Angeles, October 23–25, 2005. For more information, see:

<http://nanog.org>

The *American Registry for Internet Numbers* (ARIN) will meet (jointly with NANOG) in Los Angeles, October 26–28, 2005. For more information, see: **<http://arin.net>**

The *Internet Engineering Task Force* (IETF) will meet in Vancouver, Canada, November 6–11, 2005. For more information, visit:

<http://ietf.org>

The *Internet Corporation for Assigned Names and Numbers* (ICANN) will meet in Vancouver, Canada, November 30–December 4, 2005. For more information, see: **<http://www.icann.org>**

The *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) will be held in Perth, Australia, February 22–March 3, 2006. For more information, see: **<http://www.2006.apricot.net>**

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 2005 Cisco Systems Inc.
All rights reserved.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA

The Internet Protocol Journal

December 2005

Volume 8, Number 4

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Anti-Spam Efforts	2
Another Look at Spam	15
Testing Routing Protocols	20
Book Review	28
Letters to the Editor	30

FROM THE EDITOR

Perhaps the greatest challenge facing the Internet is the ever-increasing amount of unwanted e-mail, commonly known as *spam*. It is tempting to compare electronic mail to its paper counterpart, but there are some important differences. First, “junk-mail” is relatively self-limiting in scope because it costs real money to print and distribute even the most modest flyer. Second, advertisers in the real world are interested in *targeting* their audience. It makes little sense for a supermarket in Boston to advertise weekly specials on produce to consumers in Tokyo. Bulk mail—when delivered by the local postal service—is also quite carefully regulated. It is somewhat rare that you cannot locate the sender of paper-based advertising. None of these observations can be applied to spam. Sending spam is more or less “free,” spammers often target “the entire world,” and spammers can easily hide behind fake or transient addresses.

To date, spam has been tackled largely by applying sophisticated filtering techniques for incoming e-mail, but this does nothing to decrease the amount of actual spam sent. Anti-spam legislation has been passed in some countries, but it remains difficult—if not impossible—to pursue spammers through legal means, especially in an international context. It is therefore natural to look at technological solutions to the spam problem. If we can secure our network and authenticate its users, would it not be possible to allow only “authorized and verified” senders to send e-mail? Dave Crocker examines this problem in our first article.

Of course, no simple technical solution for spam exists, and not surprisingly there are divergent views on how the problem should be tackled. Our second article, by John Klensin, looks at spam from a different perspective and suggests some possible avenues towards a solution.

Our final article looks at routing protocol testing. Russ White examines testing mechanisms and discusses guidelines for realistic testing.

Many of you have already responded to the *IPJ Reader Survey*. There is still time to participate. If you received an e-mail invitation to take the survey, simply follow the link in the message. You can also take the survey by following the survey link on the IPJ home page: **<http://www.cisco.com/ipj>**. If you prefer to just drop us a line with your comments and suggestions you can do so by sending e-mail to: **ipj@cisco.com**.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Challenges in Anti-Spam Efforts

by Dave Crocker, Brandenburg Internet Working

It is said that the Internet teaches us one lesson. That lesson is “scaling.” The Internet comprises perhaps one billion users, millions of machines and many tens or hundreds of thousands of independent service operators. It operates in, and between, virtually every country on the planet. It is used for personal, organizational and governmental services. Therefore, it must be compatible with many different cultures, many different styles of communication and many different methods of administration. The Internet has no central point of control and operates according to no set schedule. Hence, changes must be gradual and voluntary—when we agree on what those changes should be.

In the early 1990s, the Internet grew from a small research community into a global mass market. Imagine a small town changing into a large, undisciplined city. In a large city, most people are strangers, and the strangers have a diverse range of values and behaviors. Hence, people must use much more caution with each other. In other words, the problems are not with the original way the town operated, but with changing requirements. So, spam is merely an unfortunate—but frankly predictable—example of the Internet’s success, not its failure.

This article explores the system-level complexities of the spam problem, as the intersection of social diversity, complexity of e-mail technology and operations, and specific lines of attack that seek to control spam. On the question of control methodologies, most prior work has been on analytic tools that are used by sites receiving spam, to evaluate the mail content, associated addresses or traffic flow. Recent efforts focus on assignment and assessment of an accountable identity that is responsible for individual messages or for the transit of aggregate message traffic.

The Nature of Spam

People agree that spam is a serious problem, but they have difficulty agreeing on its definition. *Unsolicited Bulk E-mail* (UBE) is probably the most useful.^[1] A spammer sends a large number of messages to many different recipients who have not requested the content. (Interestingly most spammers do not care whether a particular addressee receives the message; they merely seek to get a sufficient percent of their postings delivered to some of the addressees.)

Spam can conform to Internet technical standards and can contain no technical differences from legitimate—desired—messages. Hence, spam that violates standards or has other peculiarities might be common today, but detection efforts that are based on these anomalies offer no long-term benefits. Spammers are highly adaptable and use the easiest method that works. However what spam *always* violates are our *social* conventions. Therefore, any long-term, proactive, technical responses to it, such as formulation of standards, must follow, rather than lead our social decisions about it.

Like other social problems, we probably can control spam, even if we cannot eliminate it. This means that we must adjust to having spam as a permanent part of our social landscape, even as we seek to limit it to tolerable levels. Efforts to detect and eliminate spam have been underway for quite a few years. Some techniques have shown useful, localized results, but most only for a short time. In other words, none of the many spam control attempts, over the years, has yet reduced the amount of global spam! So we must be cautious about our expectations for any new anti-spam proposal. It also is likely that controlling spam requires an array of complementary techniques and continued efforts to adapt them, as spammers continue to adapt their own methods. This means that we need to assess any new proposal in terms of its likely *incremental* benefit, rather than as a candidate to be the *Final Ultimate Solution to Solve Spam* (FUSSP).

Changing a global infrastructure takes a long time and is very expensive. Some proposals require complex technology, while others require substantial, on-going administrative effort. Worse, some impose onerous requirements on end-users. Therefore we need to ensure that the mechanisms we deploy will have significant, long-term benefit, even after spammers try to adapt to their presence. They also must have reasonable development cost, require limited, on-going administration and be sufficiently easy to use. In evaluating the likely efficacy of a proposal, a useful heuristic is to ask whether it would be desired even if spam were not a problem. If the answer is yes, then it provides general, strategic benefit, so that counteracting spam merely adds urgency to its adoption.

The Internet provides us all with vastly better access to each other. For collaboration, or the formation of specialized communities or for personal interaction, this is wonderful. For intrusions into our privacy and threats to our online security, this is problematic. Unfortunately, the benefits and the detriments are tightly coupled. Our efforts to control e-mail's problems need to be made cautiously, lest we also reduce its benefits. Worse, our efforts need to limit the damage that might be done to innovative benefits that we have not yet envisioned.

The sender of spam incurs almost no incremental cost for a single message. It is easy to think that we should simply make e-mail be the same as sending letters or making phone calls, by directly charging the sender for every message. This cost provides a barrier against abusive, bulk use. In reality e-mail is a different kind of service, with an extensive history, and it is subject to different choices. Telephones and postal service have highly centralized, formal operational authorities, and the fees charged for their use are based on offsets to direct, real expenses. By contrast, e-mail is a highly decentralized service, with correspondents' private systems contacting each other directly, rather than having to be mediated by state-regulated utilities. If additional fees are charged, they also need to be based on the costs of real services; an arbitrary "tax" will simply create its own problems. For example, who gets the money, and why?

To retain its flexibility and its ability to support new human communication uses, we must retain the current, open model of spontaneous e-mail exchanges. Therefore, over time, it is likely that Internet mail will evolve into two logical subsets. One comprises trusted, accountable participants and the other includes everyone else. Trusted participants may be subject to less stringent checks and filtering. Perhaps more importantly when there is a problem, it is likely that mail from a trusted identity will still be delivered, while the origination agent is consulted, rather than rejecting the mail automatically.

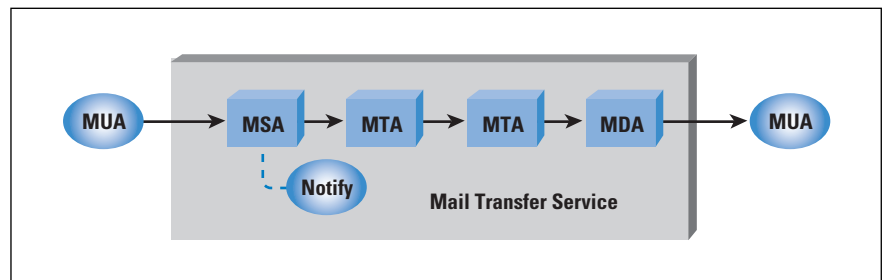
E-mail Architecture

Internet mail is based on a simple model. It distinguishes the world of users from the world of transmission. Anyone may send a message to anyone else. The basic service does not have a central authority and does not require authentication by the Originator, the Recipient or the operators. (It is worth noting that the telephone and postal services usually do not authenticate those sending letters or making calls.)

As shown in Figure 1, this model has grown to distinguish:

- *Mail User Agents* (MUA), which represent end-users
- The *Mail Transfer Service* (MTS) comprising a sequence of one or more *Mail Transfer Agents* (MTA), using the *Simple Message Transfer Protocol* (SMTP)^[2,3]
- Posting new mail via a *Message Submission Agent* (MSA)^[7]
- A *Notification Handler* or *Bounce Handler*, is an MUA that processes returned transmission reports such as a notice about failure. The Handler's address is specified by the MSA, during message posting.^[11]
- Delivering mail via a *Message Delivery Agent* (MDA), possibly with user-specific delivery behaviors^[8, 9]

Figure 1: Internet Mail Architecture



The purpose of e-mail is to exchange messages among MUAs. For users, their e-mail client—the MUA—is all they directly experience. For most network administrators, the MTS software is their scope of concern.

The core e-mail message object also has a simple framework. Its *content* comprises:

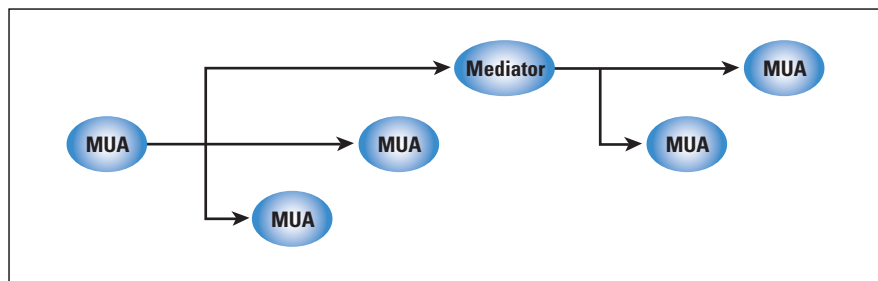
- Structured, textual meta-information, called the *header*, including *fields* for addressing, posting date, unique message identifier and a free-form description of the content^[4,5]

- Lines of free-form ASCII text, called the *body*, which has evolved to support a potentially complex, structured set of multi-media, multi-character set attachments^[12]

Figure 2 demonstrates a simple user-to-user example, with a message sent to three addressees, one of which is a special MUA that re-mails it to two additional recipients. The purpose of the Figure is to emphasize the user-to-user nature of e-mail and to provide a basis for considering the combinatorial explosion that marks the aggregate interactions of Internet mail components even in very simple uses. It further introduces another architectural construct:

- A *Mediator* is an MUA that re-posts messages, such as for a mailing list.^[10] It preserves much or all of the original message, including author address, but can make substantial changes or additions to the content, which an MTA cannot. Therefore, a Mediator's role is user-level content responsibility, rather than MTS-level transit responsibility.

Figure 2: Simple Multi-Recipient Scenario

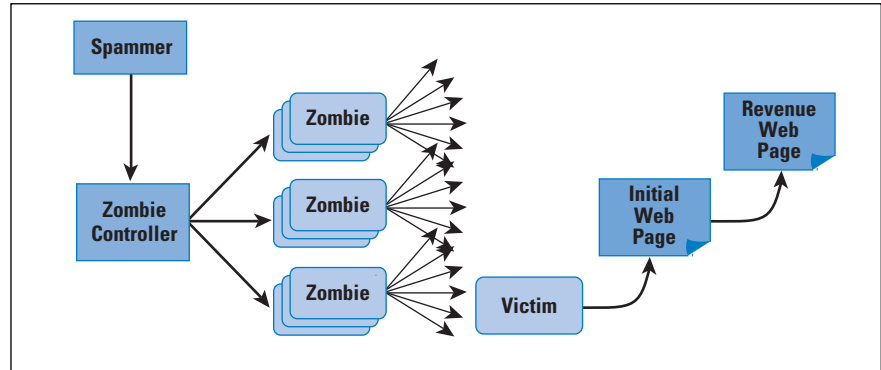


Spamming Architecture

Some spammers are legitimate businesses, engaged in overly aggressive marketing efforts, because there are no formal limits on their actions. In spite of the challenges created by needing to work at an international level, there is a reasonable expectation that legal strictures, both laws and contracts, will constrain in these businesses to a tolerable level. In contrast, *rogue* spammers actively seek to avoid accountability, to subvert barriers to their traffic, and to acquire unwitting and unwilling participation of machines owned by others. Independent of the legal details, the best social model to use for analyzing this latter group is crime. Often the activities do not violate particular laws, but what is most important is that the style of a spammer's conduct is the same as that of a criminal.

Unfortunately, the technical and operational world of spamming has also developed in scale and sophistication. Spamming used to entail one sender and one sending machine. Its performance was limited by the capacity of that machine and the bandwidth of its Internet connection. Today, rogue spammers control vast armies of compromised systems, called *zombies*, as shown in Figure 3. Zombies are owned by legitimate users who are unaware that their system has been compromised and is being used for spamming.

Figure 3: Rogue Spammer Control Network



The community of rogue spammers is remarkably well organized; it has become an extensive, underground economy. Some participants specialize in developing methods for breaking through filters. Others take over machines and turn them into zombies. Others sell the use of a zombie collection for periods of spamming. The estimated number of zombie systems is in the many tens of millions. After spam delivery, recipients often “click” to a transaction Web page. Web hosting is provided at multiple levels, in order to obscure the server side of the process, further reducing accountability.

Typically, spammers have the classic goal of selling products. However, they also can have political or religious motivations or even blatantly criminal intent, such as extortion. The ability to send very large number of messages to a specific destination gives spammers a tool that can be used to threaten an organization with a denial of service attack on their network.

Practical Efforts at Spam Control

It is tempting to believe that spam is an easy problem to solve, but history teaches us to be cautious. A web page located at <http://craphound.com/spamsolutions.txt> takes an irreverent approach in challenging simplistic proposals, by providing a checklist for the common weaknesses. In spite of its apparent whimsy, the checklist is surprisingly useful for screening proposals quickly.

The most common mechanism for spam control is a localized mechanism, the “filter”^[14], named for its conditionally permitting mail to flow through it. Filters typically are used within the recipient’s network (or Administrative Management Domain, as described later in this article.) However they may be placed anywhere along the path, notably including the MSA. Filters at the reception side cannot reduce Internet spam traffic. At the outbound side, they can. Filters have choices in the way they treat suspect messages. They can:

- Add a special annotation to the message
- Divert it into special storage
- Reject it back to its Handling Notification (RFC 2821 **MailFrom**) address or to the Client SMTP during the transfer session
- Simply delete it
- Accept it slowly, with “traffic shaping,” to control the rate of SMTP transmission

The difficult question is: What are the criteria that a filter should use? The difficult answer is: Many. This need to support a wide, and changing, variety of decision criteria has caused filtering engines to evolve into extensible platforms for spam detection and handling modules. As the mixture and complexity of filtering algorithms become more sophisticated, the overhead they entail has grown substantially larger.

It is convenient to divide techniques into three, basic classes of criteria, although each is complex:

- *Content analysis*, such as Bayesian statistics tracking of vocabulary and content hashing, to detect bulk duplication
- *Responsible Agent assessment*, either for permission (whitelist) or rejection (blacklist)
- *Traffic analysis*, such as rates at which messages come from the same author address or IP Host Address

Content analysis is always a matter of partial success (and partial failure.) It is usually statistical and depends upon a database of training messages, to establish vocabulary norms. Spammers are constantly developing techniques for bypassing the current analysis technologies. Further, different recipients on the same e-mail service can have wildly different statistical patterns of acceptable content. This makes fine-grained filtering by their service provider problematic.

It is clear that these tools for evaluating individual messages, or aggregate traffic flow, can have significant transient utility. However they cannot be effective, long-term tools, even with continuing enhancement. Notably they have little or no effect at reducing spam at its source. These post-hoc analysis tools have two inherent deficiencies, both of which are coupled to their using heuristics, rather than reliable, accurate and objective rules. The first is one of “false positives” in which legitimate mail is incorrectly labeled as spam. As an example, this could mean that an essential business transaction is not delivered, instead being classed as junk mail. Perhaps the most insidious example of this problem occurs when spammers send mail that purports to be from a well-known, legitimate business. This is called *phishing* and results in making *all* mail with the address suspect, so that legitimate postings of essential mail are not delivered.

The second problem with using heuristics is in the nature of an “arms race” between spammers and anti-spammers who must each constantly adapt techniques, consume more resources, yet never win. It does not help that those fighting spam have been losing the war, since spammers have tended to be more aggressive, more innovative and better organized...

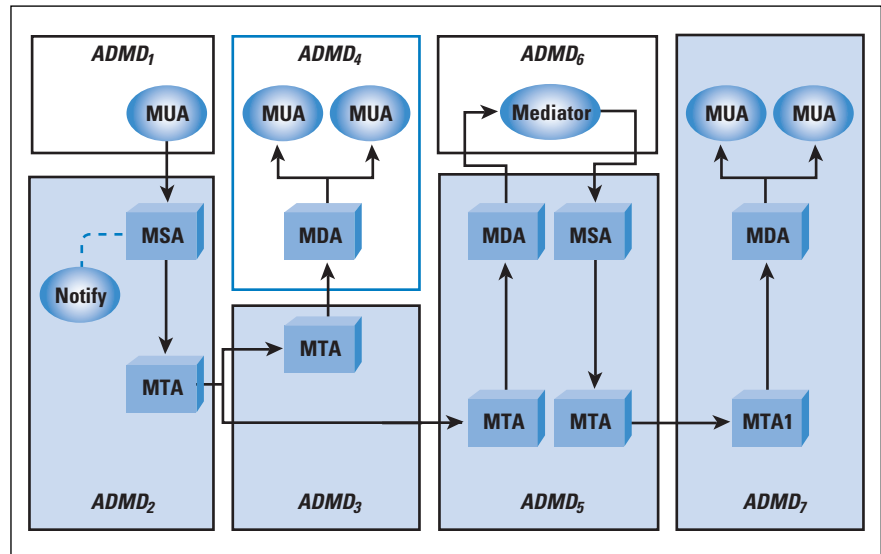
A different line of effort is based on the social assessment that the sender of an e-mail should be held accountable for it. The goal is to identify such an agent and then evaluate the agent's acceptability. This approach requires three enhancements to Internet mail:

- A clear sense of the boundaries between independent operational authorities
- A means of verifying an accountable identity that is associated with the message
- A means of formulating and sharing assessment information about accountable identities

Although e-mail operators often refer to *boundary* MTAs that face the open Internet, there is no accepted term for a region of e-mail components under unified authority. This article suggests a term derived from the OSI X.400 e-mail effort: *Administrative Management Domain* (ADMD) to mark these trust boundaries. They distinguish a collection of operational components subject to the same administrative policies, as discussed in [13].

An example of ADMDs is shown in Figure 4, and is derived from the scenario shown in Figure 2.

Figure 4: Independent Administrative Management Domains (ADMD)



The implied complexity of responsibilities and interactions is striking, even for this relatively modest case. For simplicity, think of the ADMDs labeled at the top of the Figure as representing users or value-added services, whereas the ADMDs labeled at the bottom could be a variety of classic Internet service (access) providers. The “boundary” agents are the ones with lines connecting over to another ADMD.

The increased diversity among Internet participants and ADMDs results in abuses such as spam. Proactive efforts to deal with these abuses require that we make changes in the nature of the trust between ADMDs and the way that that trust is enforced.

Accountability

Agent assessment seeks to hold an entity (agent) accountable for problematic e-mail. Who is a responsible agent for the content or for injecting the message into the MTS, and are they assessed as trusted or problematic?

There are two broad classes of accountable entities:

- *Content agents* comprise authors (RFC 2822 **From**) and those who are responsible for posting individual messages, as specified in the RFC 2822 **Sender** field. If the content agent is validated for a message, then the content probably reflects their intent. That is, it is unlikely that some other entity changed the content. Because the Notification Handler address (RFC 2821 **MailFrom**) appears in the SMTP protocol but is associated with the posting agent, it is often considered useful for analysis. Unfortunately the address often has no obvious relationship to the From field author or the Sender field posting agent, so its use for filtering can be problematic. However spammers often specify false Handling Notices addresses, in order to direct the mass of failed deliveries elsewhere. Consequently, it can be useful to validate the **MailFrom** address.
- *Operations agents* provide MTA or basic Internet access services. They are often held accountable for the impact of the bulk traffic their systems generate. Although they do not create the content, it is possible for them to enforce strict rules on their customers and to detect patterns of violations among them. Recommended practices for operators are beginning to obtain some consensus, such as with [15]. More are needed.

Assessment of agents can be proactive or reactive:

- *Accreditation* is the proactive registration by a sender, who aligns with a registry that extracts quality assurance commitments; any trust of the sender is therefore inherited from trust of the accreditation agency.
- *Reputation* refers to reactive evaluation of a sender's prior postings; for these, independent third parties evaluate the sender's history.

The functions that are combined, to establish useful accountability, comprise:

Identification: An identity label provides a unique reference to an entity.

Authentication: Validates the use of the identity label.

Authorization: Determines that the user associated with the identity is authorized to perform a particular function.

Assessment: Obtains an analysis of the trustworthiness or "quality" of the agency that is providing the authorization, or of the validated entity itself.

Unfortunately, many identities are involved in e-mail creation or transmission, as shown in Table 1.

Table 1: Roles for Internet Mail Identities

Type	Provided by	Identity of
MTA IP Host Address	Network-level service	SMTP client
EHLO Domain Name	RFC 2821 SMTP command	SMTP client
MTA Provider's IP Network Address	Network-level service	Site of SMTP client
Mail-From Mail Address	RFC 2821 SMTP command	Handling notices
From Mail Address	RFC 2822 header field	Author
Sender Mail Address	RFC 2822 header field	Posting agent
Received Domain Name	RFC 2822 header field	Relaying MTA site

Relative to an SMTP Server that is being asked to accept a message, the SMTP Client is an agent of the operator of the previous hop. Since the e-mail operator might be different from the operator of the IP access network that is hosting the e-mail service, it might entail a different identity. This highlights an interesting aspect of Table 1: Most of the identities associated with e-mail handling can be called “the sender.” Consequently, that term has become nearly meaningless, in anti-spam discussions.

Because identity listings are made explicitly in a database, they are capable of producing almost no false positives, although there might be many identities not listed and a listing might be inaccurate. Still, there are significant challenges with the use of identity-based filtering:

- Which identity should be used and how does it relate to spamming behaviors? Note that Table 1 listed quite a few choices. In addition an author can create bad content, but the identity listed in the RFC 2822 **From** field of that content might not be the actual author, even if that field is validated. The message might have originated on a compromised machine and used the identity associated with it, unbeknown to the owner of the machine. Also the operator of the mail-sending network might have nothing to do with creating content, but it might be reasonable to hold the operator accountable for aggregate traffic problems.
- How is the identity validated (authenticated)? What entity is doing the validation? How does it relate to the identity being validated? And why is it trusted? Can the validation mechanism, itself, be tricked?
- How is an identity determined to be a spammer or non-spammer? What entity is vouching for the quality of that identity and why is the vouching entity trusted?

Authentication Standards

Accountability requires having an accurate, reliable identity of the agent that is to be accountable. Authenticating an identity is, therefore, a prerequisite for assessment efforts. However it does not, by itself, ensure a positive assessment. Spammers can register and authenticate their identities, too.

Early anti-spam identity schemes use the IP Address of the client SMTP MTA that is sending directly to the server running the filter. The Address is provided by the underlying network service, and therefore has been trusted. However, spammers are becoming proficient at stealing IP Address space, such as by advertising routes that use allocated-but-unused blocks of IP Addresses! Also an IP Address changes as the host changes its attachment to the Internet, and it is affiliated with operators, not authors. This makes the IP Address obscure and unreliable, when attempting to assess e-mail.

A more recent focus is on the use of Domain Names, for references that are more stable and align better with the authority boundaries of Administrative Management Domains. Broadly there are two lines of effort at using Domain Names for validating messages being relayed. One associates the identity with the systems that handle the message along its path. These “path registration” schemes include Sender Policy Framework, Sender-ID, and Certified Server Validation. The other schemes tie a Domain Name identity to the message object. These include Domain-Keys Identified Mail, and Bounce-Address Tag Validation.

The *Sender Policy Framework* (SPF)^[16] has evolved over time, attempting to encompass multiple identities. It primarily uses the Domain Name in the RFC 2821 **MailFrom** command. It queries the *Domain Name System* (DNS) with that name and determines whether the IP address of the previous-hop MTA is registered under that name. Since any SMTP server along the transit path may choose to perform this query, SPF requires that the Domain Name contain a registration for every MTA along every delivery path for a message. (A common simplification for this model is to use it only between boundary MTAs, but this considerable constraint is not specified in SPF. Rather, its use is usually characterized as being more general.) Although the software overhead for SPF is quite small, the administrative overhead can become substantial, as the number of paths increase and as paths change. In addition, some sender SPF DNS configurations can trigger a very large number of queries per addressee. Lastly, the role of the RFC 2821 **MailFrom** command is to specify the Notification Handler address. This address might be entirely different from other origination information, making registration of all of the MTAs in the path problematic. SPF therefore has significant administrative problems with redirected traffic, such as when going through a third-party forwarding service.

Sender-ID (SID)^[17] uses a model similar to SPF, but it is based on the posting address Domain Name in the RFC 2822 **Sender** field (or RFC 2822 **From** field, if no **Sender** field is present.) Both SID and SPF sought IETF standardization in 2004 but the working group effort failed, due to lack of rough consensus convergence among participants and due to concerns over intellectual property claims.

Certified Server Validation (CSV)^[18] covers only the current client/server SMTP hop. The client specifies an operator's Domain Name in the RFC 2821 **EHLO** command. The server uses this name to query the DNS. It then validates the IP Address of the SMTP client and determines whether the Domain Name administrator has authorized the client to send mail. CSV also specifies a standard mechanism for querying an assessment service about the client's Domain Name.

DomainKeys Identified Mail (DKIM)^[19] specifies an accountable Domain Name that applies to a message during transit. It uses public key cryptography to digitally sign the message and provides guidance when the signing Domain Name differs from the Domain Name in the RFC 2822 **From** field.

DKIM Domain Name validation represents a significantly different goal from that of the strong authentication methods, such as [20, 21] which focus on long-term protection of message content. Also DKIM places its parametric information in a special RFC 2822 header field, rather than in the message body, so that it does not have any impact on recipient user agents that do not support DKIM. Although public key cryptography has relatively high computational cost, e-mail processing is usually i/o-bound, so that the real-world use of DKIM appears to have little impact on the aggregate message-handling capacity of a server.

Bounce Address Tag Validation (BATV)^[22] attacks the problem of mis-directed handling notices, such as bounces. It permits the creator of an RFC 2821 **MailFrom** bounce address to digitally sign it. When the bounce agent of that creator receives a message purporting to be a bounce, the agent can validate the address. Standardization of its format is needed so that e-mail intermediaries—such as some mailing list software—can determine the “core” of the mailbox portion. Since the creator of the signature semantics is the only consumer of the signature semantics, any signature algorithm can be used, including one based on symmetric keys. For convenience—and an existence proof—the BATV specification provides an example algorithm already in use.

Collaboration Support

Fighting spam must be a collaborative effort, which will benefit from using tools and standards that aid in exchanging information and performing coordination. To this end, standard methods of reporting spamming events, of characterizing particular spam, and of sending spam control data can be helpful. Some work in that direction is already underway.^[23] Fighting spam requires global operations collaboration; this will be aided by services to facilitate interactions between network administrators speaking different languages. It is also likely that there should be standards for the syntax and semantics of whitelists and blacklists.

Acknowledgement

The author wishes to express particular appreciation for the unusual amount of dialogue that took place with the reviewers of this article. It produced a substantially clearer and more concise article. It also highlighted the extraordinary diversity of views on the topic, in case one had had any doubt. In fact, the article by John Klensin which follows this one is a direct result of the dialog.

References

- [1] Hoffman, P. and D. Crocker, "Unsolicited Bulk Email: Mechanisms for Control," Internet Mail Consortium, UBE-SOL IMCR-008,
<http://www.imc.org/ube-sol.html>, revised May 4, 1998.
- [2] Postel, J. B., "Simple Mail Transfer Protocol," STD 10, RFC 821, August 1982.
- [3] Klensin, J., "Simple Mail Transfer Protocol," RFC 2821, April 2001.
- [4] Crocker, D.H., "Standard for the format of ARPA Internet text messages," STD 11, RFC 822, August 1982.
- [5] Resnick, P., "Internet Message Format," RFC 2822, April 2001.
- [6] Crocker, D., "Internet Mail Architecture," Internet Draft, **[draft-crocker-email-arch](#)**, April 2005.
- [7] Gellens, R. and J. C. Klensin, "Message Submission," RFC 2476, December 1998.
- [8] Myers, J. G. and M. T. Rose, "Post Office Protocol – Version 3," STD 53, RFC 1939, May 1996.
- [9] Crispin, M., "Internet Message Access Protocol – Version 4rev1," RFC 3501, March 2003.
- [10] Chandhok, R. and G. Wenger, "List-Id: A Structured Field and Namespace for the Identification of Mailing Lists," RFC 2919, March 2001.
- [11] Moore, K., "Simple Mail Transfer Protocol (SMTP) Service Extension for Delivery Status Notifications (DSNs)," RFC 3461, January 2003.
- [12] Freed, N. and N.S. Borenstein, "Multipurpose Internet Mail Extensions (MIME) Part One: Format of Internet Message Bodies," RFC 2045, November 1996.
- [13] Clark, D., Wroclawski, J., Sollins, K., and R. Braden, "Tussle in Cyberspace: Defining Tomorrow's Internet," ACM SIGCOMM, 2002.

- [14] Showalter, T., “Sieve: A Mail Filtering Language,” RFC 3028, January 2001.
- [15] Hutzler, C., Crocker, D., Resnick, P., Sanderson, R., and E. Allman, “Email Submission: Access and Accountability,” Internet Draft, **draft-hutzlerspamops-05**, October 2005.
- [16] Wong M., Schlitt M., “Sender Policy Framework (SPF) for Authorizing Use of Domains in EMAIL, version 1,” Internet Draft, **draft-schlitt-spf-classic-02**, June 2005.
- [17] Lyon J., Wong M., “Sender ID: Authenticating Email,” Internet Draft, **draft-lyon-senderid-core-01.txt**, May 2005.
- [18] Crocker D., Leslie J., Otis D., “Certified Server Validation (CSV),” Internet Draft, **draft-ietf-marid-csv-intro-02**, February 2005. Also see: <http://mipassoc.org/csv>
- [19] Allman E., Callas J., Delany M., Libbey M., Fenton J., Thomas M., “DomainKeys Identified Mail (DKIM),” Internet Draft, **draft-allman-dkim-base-00**, July 2005. Also see <http://mipassoc.org/dkim>
- [20] Ramsdell B. (ed.), “Secure/Multipurpose Internet Mail Extensions (S/MIME) Version 3.1 Message Specification,” RFC 3851, July 2004.
- [21] Elkins M., Del Torto D., Levien R., Roessler T., “MIME Security with OpenPGP,” RFC 3156, August 2001.
- [22] Levine J., Crocker D., Silberman S., Finch T., “Bounce Address Tag Validation (BATV),” Internet Draft, **draft-levine-mass-batv-00**, September 2004. Also see <http://mipassoc.org/batv>
- [23] Shafranovich, Y., “An Extensible Format for Email Feedback Reports,” Internet Draft, **draft-shafranovich-feedback-report-01.txt**, May 2005.

Ed.: This article is a revision of “Adapting Global Email for Controlling Spam,” in *Information Processing Society of Japan (IPSJ) Magazine—Special issue on Anti-Spam*, Japanese/English, Volume 46, No. 7, pp. 741–746, July 2005.

DAVE CROCKER is a principal with Brandenburg InternetWorking. He has authored or contributed to most Internet mail standards, and an assortment of e-mail products and businesses, as well as working on facsimile, security, ecommerce and EDI. He received the 2004 *IEEE Internet Award* for his work on e-mail. Dave is a contributor to the development efforts for DKIM, CSV and BATV, motivated by a strong desire to protect more than 30 years of professional investment that is being threatened by spamming. E-mail: dcrocker@bbiw.net

Taking Another Look at the Spam Problem

by John C. Klensin

The problem of unsolicited bulk e-mail on the Internet has been widely discussed, and many classes of solutions have been proposed. Dave Crocker's article discusses some of the background for the solutions generally, points to a semi-humorous list of ways in which proposed approaches fail, and compares several approaches based on source authentication. This article takes a somewhat contrarian view. It argues specifically that the traditional models for defining technological solutions and then letting the policy and legal communities work out the details of how to utilize them are seriously wrong in this particular case and that partially-effective methods of fighting spam actually cause more spam.

This article makes two main suggestions. First, attempts to design technological countermeasures to spam without a clear understanding of how far, and in what directions, the setters of social policy are willing to go are futile. The requirement is not just that there be social recognition that a problem exists. In order to design effective technological countermeasures with predictable and acceptable side-effects, we must first understand what measures society is willing to take—what laws it is willing to pass and enforce to make spam a criminal or civilly-punishable act—to set an appropriate context and set of boundary conditions. Without those conditions, design of technological countermeasures is likely to constitute poor engineering practice, not just futility. Second, deployment of spam counter-measures that are not completely effective largely shifts the burdens of spam from one recipient population to another while *increasing* the total amount of spam on the network.

His analysis and mine agree on several critical points. Solutions that discard important characteristics of today's e-mail environment permanently in order to make some short-term gains against spam are not acceptable. Approaches that require drastic and simultaneous changes to the ways in which e-mail works in order to function are not going anywhere. There is a difference between legitimate businesses who have decided, within the limits of existing legislation, to engage in mass, unsolicited, electronic mailings to promote their products and those bulk mailers who prefer to cover their tracks, hide linkages between sending addresses, hosts, and web sites (or create deceptive ones), and who use zombie mailers and other ways to avoid cost and detection. We also agree that spammers, or their tool suppliers, are creative, technically-knowledgeable, and able to react much more quickly than the spam-fighting community (especially the standards-based part of that community) to changes in operating conditions and countermeasures.

I suggest a further guideline to help us think about the problem: however small they might be on a per-message basis, there are costs associated with sending e-mail and costs associated with receiving it and eliminating undesirable content.

If an anti-spam “solution” is developed that permits the spammers to vastly increase the costs to the recipients without a proportionate increase in their own costs, that solution is not tenable. A serious effort to predict the impact of a proposed solution to spam, including costs to the end user and load on the network as the spammers adapt to it, should be a critical component of such efforts. But, while equivalent analyses of measures, likely responses, and countermeasures are standard with any (other) technique designed to enhance network security, they have been largely absent when new technological approaches to spam are proposed.

This is a different aspect of the so-called “arms race” problem. In a classic arms race, no one can really win, as Dave points out. But, more important, when such races stop, it is only because one party simply stops, is forced out of the game by external pressures, or becomes exhausted economically. As long as there are no economic constraints, every escalation is met with a counter-escalation, which is met with a counter-counter-escalation, and so on. It is this positive feedback cycle that characterizes a true arms race. The battle against spam demonstrates a particularly unfortunate variation on that pattern in which the incremental economic costs of trying to deploy new spam abatement measures appear to be much more severe than the costs to the spammers of the most obvious counter-measure to improved spam abatement procedures, simply sending out more traffic. This is discussed further and in context below.

Social Problems and Technological Solutions

In the technical and protocol design community, our normal model is to develop technology and then use it to inform the policy, social, and legal parts of the society who then need to sort things out on their side. One of the classic arguments for this approach, which does not seem relevant to the spam situation, is that the potential use or misuse of a technology will not, and should not, constrain its development. For spam, the situation appears to be exactly reversed: we need to understand what is feasible and plausible from social, political, legal, and regulatory standpoints in order to define the engineering solution space. If we do not know what behaviors society is willing to make illegal or subject to effective civil action and whether it is willing to enforce those laws or equivalent positions, we cannot adequately define the engineering solution space. That results, in turn, in a high risk of solving the wrong problem or an irrelevant one. Of course, recent history has shown a variety of irrelevant and costly solutions to spam proposed, and sometimes deployed.

The solution to spam is identical to the solution to most other significant social problems: society must determine that it is a problem, create effective rules prohibiting the problem, and then enforce those rules aggressively and consistently. Technical solutions that make it easier to identify spam and its sources can then be immensely useful, but they are only useful if designed to be effective within the framework set by those rules.

If, by contrast, societies are, in practice, unwilling to take effective social or legal action against spam and those who benefit from it, then this article suggests that anti-spam measures will tend to make the overall situation worse.

The question of spam beneficiaries provides a particularly good illustration of this point. So far, most legal systems in the world have taken the position that the act of spamming is the offense (if there is any offense at all). Operating a domain or web site to which the spam recipient is directed to buy a product or obtain another benefit is rarely considered a problem by either law enforcement or by the relevant ISP. While establishing cause and effect—that the spam was authorized or encouraged by the web site owner—can be quite difficult, there has, appropriately, been little examination of tools to detect or identify beneficiaries because doing so seems pointless. On the other hand, on the same theory that it is more useful to try to arrest the drug importer than the street dealer, a different set of laws about beneficiaries and spam-authorizers—those who, in at least some cases, pay the spammers to spam—might dramatically change the landscape.

Reducing Spam by the Percentages

A new technique or group of techniques that claims to be beneficial can have either positive or negative value with regard to the amount of spam that gets through, either overall or to the mailbox or a particular sample user. A technique can also result in significant increases in the amount of network bandwidth or server resources consumed if it is neutral or better with regard to the end user mailbox. As long as the spammers can increase the number of messages they send out, almost arbitrarily and at low or zero marginal cost, the percentage of spam that is filtered out is ultimately irrelevant. The key measurement is how the amount of spam that gets through to some exemplar user (or a statistical aggregate of them) changes. That change pattern can be net either positive or negative. Suppose a technique is introduced that causes an initial small incremental reduction in the amount of spam delivered. The patterns of the last several years suggest that the spammers will respond by making a large increase in the amount of traffic they send out. Since the costs of doing so are very low, it would arguably be irrational for them to do anything else. If the increased volume is enough larger than the amount of spam the new technique was able to stop, there is a net loss to the Internet overall: the small improvement may represent a percentage decrease in the amount of spam that gets through, but the amount seen by the representative user increases and the percentage claims are largely irrelevant.

Unless whatever methods that are used in an attempt to reduce the amount of spam actually stop it at, or very near, the point of origin, the net effect on users is to shift the amount of spam received from those who have deployed the latest and most effective countermeasures to those who have not yet done so. The total amount of spam-related traffic on the network just continues to rise. And, since most countermeasures have costs—either in processing time or in software licensing fees—the cost burdens on end users also continue to rise.

This would seem to argue for methods that cut off spam traffic close to the source, but attempts to design such methods have been fairly unsuccessful, sometimes because of another policy problem: the spammers argue that some people like receiving unsolicited bulk commercial e-mail so that cutting off bulk traffic near the point of origin prevents legitimate and desired traffic from transiting the network. Source-oriented techniques include not only technical approaches but efforts—by law or social pressure—to hold ISPs and mail providers responsible for all traffic emanating from their networks, thereby encouraging them to refuse to have spammers as customers, to aggressively enforce terms and conditions of service, and so on. The strongest advocates of the “blacklist” variation of those techniques continue to claim that they are very effective although some others in the community are not completely convinced.

The House-Burglar Analogy

In the absence of a coordinated approach that is oriented toward legal or social enforcement, most anti-spam techniques appear to induce more spam on the network. They do this by making simply sending much more traffic out the most rational behavior for a spammer who is faced with an abatement technique to adopt. They may enable shifting the burden of dealing with that spam from one person to another—in the same way that aggressive locks and alarm systems on one house slightly increases the relative burglary risk to the less-protected neighbor—but, as Dave’s article points out, we have no realistic plan for making it too expensive for the spammers to simply increase output.

Deterrents to burglary work moderately well because they increase the costs (in time, sophistication of the required tools, and so on) to the burglar. Equally important, they increase the risks of being caught and punished. In the present spam environment in most countries, we have no effective mechanism to increase costs and, at least statistically, the odds of being effectively punished even if caught are insignificant.

Shifting Burdens and Creating Preferred Classes of E-mail

The argument Dave presents for authenticated mail is ultimately that it can get expedited handling while non-authenticated mail is put aside for other methods of spam detection. That approach could be immensely effective at expediting receipt of some mail by the recipients who apply the needed checks, at least until the spammers begin authenticating their mail in a way that tricks the trust-establishment techniques. Prioritization of some messages and content will be effective as long as the fraction of such messages remains relatively small relative to the total number of messages received. As the percentage rises, one probably ends up either trusting all mail from a particular source, regardless of the author, or with a situation quite analogous to “whitelists,” although one that is much harder to trick than the original. Either is subject to attacks and scaling problems.

There is also the risk of abuse by providers who conclude that mail that cannot be authenticated well enough that their users can prioritize it should simply be rejected and who then define the conditions for adequate authentication in terms of a small circle of cooperating mail providers. Even if the types of authentication outlined in Dave's article are used only as intended, the costs to recipients will rise, perhaps rapidly, over time as percentages of messages bearing authentication information rises and sender authentication and authorization become just one more tool to distinguish probably-desired messages from probably-undesired ones.

Maybe there is not Enough Spam Yet

One of the depressing consequences of the reasoning discussed previously is that perhaps we have yet to see sufficient spam for governments and regulatory bodies to take the spam problem seriously—seriously enough to deploy effective laws and enforcement mechanisms. If spam-fighting methods shift the burdens of receiving spam away from those who have the resources to protect themselves they may simply place the spam impacts on others who have fewer resources. That pattern may, in turn, also reduce pressure on governments to take effective action and to do so in a way that would make the design constraints for effective technological approaches clear. If a collection of anti-spam methods have the effect of simultaneously increasing the amount of total spam on the network and of decreasing pressures on societies and governments to take effective action, are they really ones we want to deploy?

Conclusions

This article presents a rather grim view of the future if we continue on our present course. If we fail to examine the actual actions that societies and their governments are willing to take to deal with spam and spammers and to treat those actions and their limitations as design constraints on the technical and engineering approaches, we are likely to continue to see an ever-increasing amount of spam on the network. Spammers will not only adopt technical countermeasures to new techniques but they will also take advantage of their ability to simply increase message volumes (at almost no cost) to counter the effects of those techniques on the percentage of spam that is delivered. It may be time to finally deal with the spam problem as the difficult social issue that it is, rather than permitting societies and governments to continue to believe that a technological “silver bullet” is right around the corner and that no real social or political action, or commitment of law enforcement resources, is needed.

JOHN KLENSIN is an independent consultant based in Cambridge, Massachusetts. He has been involved in the design, development, and deployment of ARPANET and Internet applications, and occasionally lower-layer technologies, since the late 1960s and early 1970s. He has also been intermittently involved with Internet administrative and policy issues since the early 1980s. His current work primarily focuses on internationalization of the Internet on both technical and policy dimensions. E-mail: klensin@jck.com

Caveats in Testing Routing Protocol Convergence

by Russ White, Cisco Systems

In general, the main problems we find when testing routing protocols lie in generating accurate (or rather, realistic) data, as well as understanding the limitations of tests geared towards measuring routing protocol performance. Three areas of specific interest are covered in this article: defining convergence, taking realistic measurements, and creating realistic data.

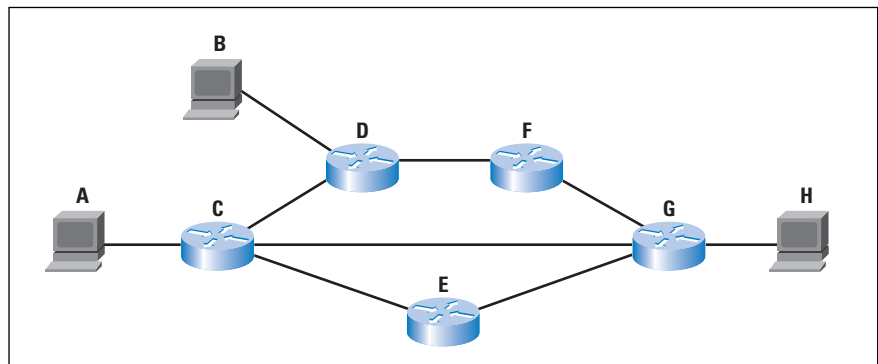
Defining Convergence

The first problem we face when trying to test routing is to define *convergence*. It seems like a simple question, but it's not, because there are so many different ways to measure convergence:

- How long does it take to begin forwarding traffic once a topology change has occurred?
- How long does it take for every router in the network to adjust to a topology change that has occurred?
- How long does it take for the forwarding information on a specific router to be updated once a topology change has occurred?
- How long does it take for the routing protocol to adjust to a topology change?

Each of these questions is actually completely different, as a short examination of the network in Figure 1, below, shows.

Figure 1: Test Network



Assume A is the traffic source for a test, and H is the sink, or the convergence measurement point. To measure the convergence time of this network, you send a stream of traffic from A to H; when the traffic stabilizes, the C to G link is taken down, and the length of the gap in traffic at H is measured. In this environment, we assume the path fails off of the C to G link, and onto the path through E.

This test assumes the traffic between B and H, or between A and B, will not be impacted by the link between C and G failing, but we do not know this will always be the case. In fact, it's possible that D and F will end up forming a *microloop* until they receive all the information needed to converge without the C to G link.

This microloop could last longer than C requires to recompute a path to H, so while the traffic from A to H may be successfully delivered, the network may not be in a fully converged state. The topic of microloop formation and avoidance is beyond the scope of this article.

In this small network, the time it takes for A to continue forwarding traffic to H may not be the same as the time it takes for the entire network to stabilize after the topology change. How long it takes for A to be able to reach H, and how long it takes for all the routers in the network to adjust to the topology change are two different questions. In this case, the concept of convergence is unclear, with several possible meanings; to properly build and understand the results of the test, we need to better understand the question being asked.

You could alter the test so only A, C, E, G, and H are in the network. This would provide a “clean” test of just the failover capabilities of the routing protocol being tested, as it’s implemented on the specific routers in the network, across the specific link types connecting the routers, in the simple failover situation. While the limited topology does limit the number of outputs being measured in the test, it also limits the closeness of the tested network to a real network design. The test can provide some very specific data points, but, once the test topology is simplified, it cannot provide a true picture of convergence in a larger, more complex topology.

Another option is to refine the test procedure so the traffic between B and H is tested as well as the traffic between A and H. Measuring traffic flow from every possible connected end point to every other possible connected end point on the network provides a number called *goodput*, which is the relation between the traffic injected into the network versus the traffic the network delivers across all paths.

Although this type of testing does provide more data in a more complex topology, it also has its drawbacks. For instance, if you are trying to compare two different implementations of a single protocol, or compare two different routing protocols, this test not only counts the amount of time required for the routing protocol to converge, it also tests the amount of time required to note the topology change, the time required to install the newly computed routes into the local routing table, and the time required to pass the changes from the routing table to the local forwarding tables. This might—or might not—be a good thing.

Isolating just the routing protocol can provide information about the performance of a specific implementation of the protocol in specific network designs, and under certain conditions. Including platform and media-specific issues—such as the installation of information into a local table—may cloud the picture. For instance, if the routing protocol can converge in milliseconds, but it takes seconds to determine that the link between C and G has failed, any changes in routing protocol convergence time will be lost in the much larger link failure detection time, reducing the value of the test.

In short, numerous tradeoffs are involved in designing a test to measure routing protocol convergence; you need to begin with the right questions, and understand the tradeoffs in the various tests you could, or might, run. There's no "simple" way to run a single test that will give you all the information you need to know to understand all possible implementations of a routing protocol on all possible platforms.

In the same way, it's important to keep these types of limiting factors in mind when reading, or using, test results provided by outside companies. It's fairly easy to look at a specific test for one measure, such as the number of neighbors a specific implementation of the *Border Gateway Protocol* (BGP) can support in specific conditions, and attempt to generalize those test results to much larger and varied real world networks. Quite often, the mapping isn't all that simple.

Taking Realistic Measurements

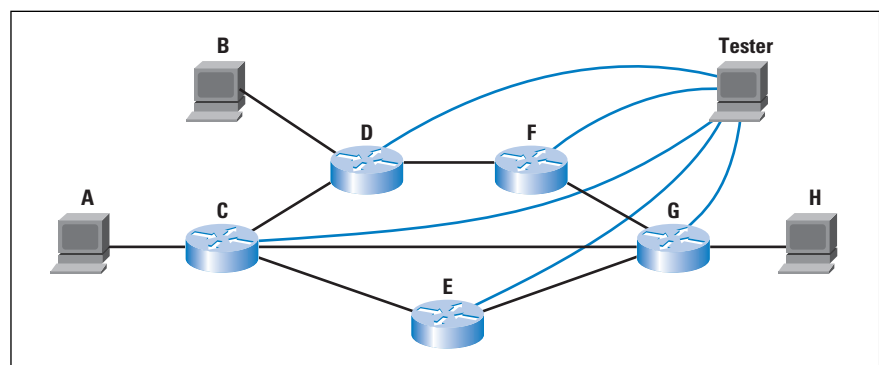
Assume you determine you want to test for protocol convergence by checking the routing tables at each router in the network in Figure 1, rather than trying to measure convergence by measuring traffic flow through the network. How would you go about doing this? There are two general types, or classes, of tests, that you could consider:

- *Black Box*: Treat the device as a black box, only using outside signals and controls, and never any output provided from the device itself.
- *White Box*: Use available output provided from the device itself, possibly with tests using signals outside the device, to determine when specific events on the device occur.

Obviously, black box testing is much more difficult, maybe impossible in some conditions, but, at the same time, can provide more "objective" measures of a devices' performance. Examples of black box tests for the *Open Shortest Path First* (OSPF) protocol are outlined in RFC 4061, RFC 4062, and RFC 4063. White box testing typically depends on *debug* and *show* commands to provide timestamped information about when specific events occur, such as when the routing protocol has received information about the topology change, when the routing protocol has finished computing the best path to each destination, and other events.

For simplicity, the network is reconfigured with a test measurement device, as shown in Figure 2, below.

Figure 2: Reconfigured Test Network



Some mechanism is used to determine when the routing protocol on each router has computed the correct routes; the network is connected, and allowed to converge. The link between C and G is taken down, and the time between the link failure and the correct routes being computed on C, D, E, F, and G is taken as the total convergence time in the network. This appears to be a straight forward test; what sorts of problems can we run in to here?

There are two possible mechanisms for determining when each device has correctly computed the routes after the C to G link fails:

- Some sort of “continuous output,” such as a *debug*, can be configured on each router, and the results collected and analyzed.
- The Tester can poll each device, using *show* commands, or some black box testing technique, to determine when device has recalculated the routes correctly.

Let’s examine each of these techniques separately.

Gathering Results from Continuous Router Output

The first, and simplest, mechanism is to gather the results from each router through debugging information provided by the protocol implementation which is generally used for troubleshooting and monitoring the routing protocol. There are three primary issues related to using this information you need to be aware of:

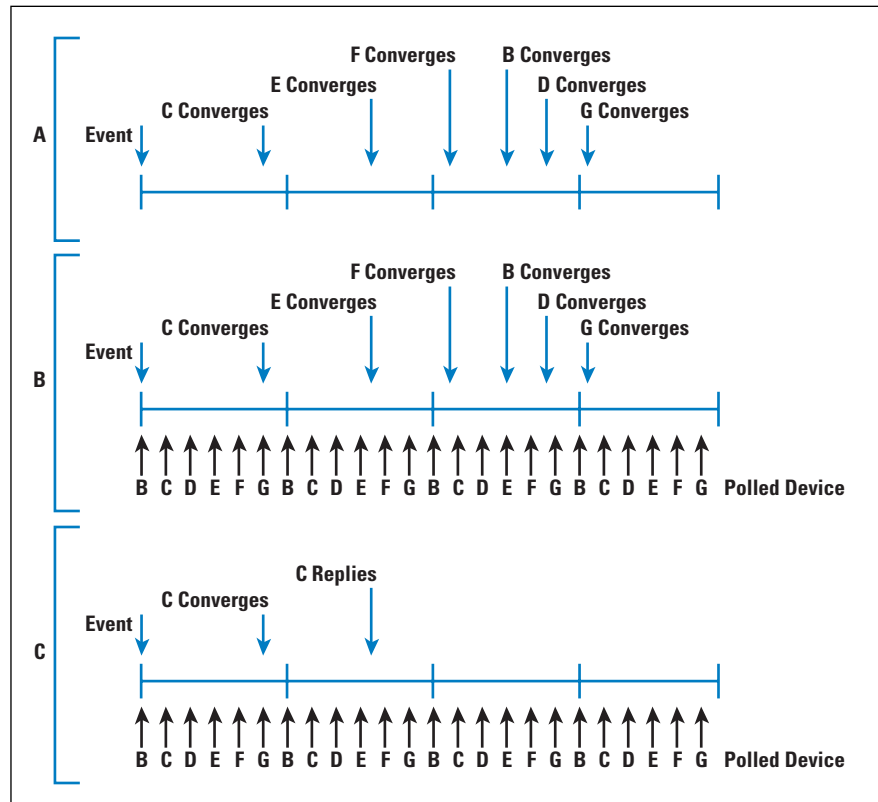
- The continuous stream of information provided by the device being tested can actually impact the test results, primarily because of the processor cycles required to record and display this information. In some situations, the additional cost is negligible, and in others, it’s simply not important (for instance, if the test is designed to show the differential between two situations, rather than provide absolute convergence times).
- If the timestamps injected by the devices being tested in the network are relied on, then the time clocks of every device must be synchronized. This synchronization must generally be within about 1/10th or less of the total variation in the test time for the results to be meaningful. In other words, if the timeclocks on all the devices are synchronized within one second of each other, and the results of the test are expressed in milliseconds, the actual test results are going to be lost in variations in the synchronization of the timeclocks.
- If the devices feed their information to the Tester, and the timestamp on the Tester is used to compare the event times within the network, the timestamps can be skewed by the packet processing requirements of the devices, as well as queuing delays in the Tester. Most routers prioritize routing traffic over switched traffic, and switched traffic over management traffic. There could be significant lags between an event occurring, and the router actually building a packet noting the occurrence of that event. Again, this is a matter of time differentials; if the test results are expressed in milliseconds, queuing delays alone can bury the results in noise.

We need to be careful when using *debug* or other continuous output to measure network convergence times in any given test, then. Quite often, we need to compare the granularity of the test results with the measurement technique used, and consider how much noise the measurement technique is actually likely to inject into the testing environment, compared to the test results granularity.

Polling Devices

Another common technique is to run some sort of process on the Tester which polls each device, either using some black box or white box measurement, to determine when each device finishes recalculating routes after the topology change has occurred. This type of test is also constrained by various factors that might not be obvious when you are designing a test, or examining the results of a test that uses it. Assume events in the network occur as Figure 3 illustrates.

Figure 3: Poll Testing Scenario



In Figure 3A, we assume that the Tester is able to poll every device in the network at the same time, once a second. The test shows the network converged at 4 seconds after the event, although the last router to converge, G, does so just after the 3 second mark. There can be a variation of the entire polling interval in the actual results without the test showing any difference in the convergence time of the network, implying that the polling interval must be much faster than the expected (measured) test results for the results to be meaningful. We normally suggest that the polling interval be about 10 times faster than the expected measurement rate, or that the Tester should poll every 1/10th of a second in this test, if the results are to be measured in seconds.

However, in real test environments, a test device cannot actually poll every device in the network at the same time. Instead, the Tester will poll one device periodically, rotating through the polled devices, so the longest time between any specific device being polled is the polling rate. We can call this rotating polling *serialization*, and the time it takes to rotate through all the devices the *serialization delay*. Here, we've spread out the polls across the total one second polling time, to illustrate, in Figure 3B. Three anomalies show up in this illustration:

- The total time for the network to converge is still just over three seconds, while the recorded test time is still in the four second range. This is similar to the problem we noted when we assumed the Tester was polling all the devices in the network at the same time.
- It appears, from our test results, that E and F have converged at about the same moment. In reality, their convergence is separated by almost one second. In some extreme cases, the devices may actually converge in the opposite order from the order they appear to converge.
- If the convergence order of D and G were to be reversed, the network would appear to converge almost a half a second faster, although the actual convergence time would remain constant. This could cause a widely diverging set of test results over multiple runs in what is, actually, a fairly consistent network convergence time.

Adding the serialization delay of polling isn't enough, however, to understand polling in real test environments. We also need to remember that each device which is polled must also answer each one of the polls, thereby introducing another variable amount of delay into the test results. For instance, in Figure 3C, C is polled once before and once after it converges. If we take the time that C answers as its convergence time, then we are also including processing time on C, which is variable, into C's total convergence time. However, if we take the polling time as C's convergence time, it's possible that the poll was received before C converged, and was processed, and answered, after C converged, skewing the results in the opposite direction.

Unfortunately, there are no simple answers to these problems. Instead, when you are designing a test, or examining the results of a test, the mechanism used to determine convergence, the rate at which that mechanism is used, and the reported final results, should be taken together, and considered closely. A test which reports results in milliseconds, but polls a large number of devices from a single test device, should be examined closely for serialization delay errors.

Use Real-Life Configuration Parameters and Prefix Attributes

Finally, we need to consider what is probably one of the most widely disregarded concerns in testing routing protocol implementations: building accurate and repeatable data sets to feed into the test. Let's examine a common test, to help in understanding this problem.

A network engineer sets up a router connected to a router testing device using a SONET link. The router tester is then configured to feed one million routes, through BGP, to the router being tested. The test is run, and the amount of time it takes for the router to accept and install all of the routes into its local tables is measured. The router is disconnected (we'll call this first router A), and another router (B) is connected. The same test is performed. In the end, the network engineer proclaims A has a better BGP implementation than B, because A accepted and installed the routes fed to it faster than B.

This sort of test, and these results, should raise a lot of red flags for anyone who's ever tested routers before. Many questions here are not answered:

- Were both routers tuned to optimum parameters for this specific test? Most routers are installed in a number of different situations in various networks, and most will perform better if they are tuned to fit the role they are playing in the network. This is similar to tuning a server for database use, or web server use.
- BGP is very sensitive to the data transmitted from one router to another; BGP implementers are generally aware of this, and use differing models of BGP behavior in different networks to tune their implementations. Specifically, in the case of BGP:
 - What percentage of the prefixes transmitted were of specific prefix lengths? What percentage of the routes transmitted were /24s, /23s, and so on?
 - How many different attribute sets were represented in the routing information transmitted? What number of unique attribute sets were included in the routes? For each attribute set, what percentage of the table did that attribute set represent?

Each of these questions can, and should, be compared to real world measures in the network the router is going to be installed in. There are some instances where protocol implementers have tuned their implementation for use in an Internet *Point of Presence* (POP), for instance, and the implementation doesn't fare as well as a route reflector, or the other way around. For some vendors, this tuning could even be on a platform by platform basis, making the job of characterizing a specific implementation through a simple test, like that described above, very difficult.

Conclusion

Designing, executing, and evaluating the results of a test attempting to measure network convergence is much more complex than it appears on the surface. In any given test situation, we need to ask:

- What was the test designed to measure? Is it measuring the appropriate outputs, in the correct ways, to actually measure this?

- What is the granularity of the test results and the actual network events, compared with the measurement techniques used in the test? Will normal test results get lost in the noise introduced by the measurement techniques?
- What is the data set used to build the test? Does it accurately reflect the data the routing protocol implementation will be handling in a real network (or more specifically, the real network the router will be installed in).

When designing, or evaluating, test results, there's a strong tendency to be dogmatic about the results, to say some specific test proves, in some way, a specific vendor, platform, protocol, or implementation, is "better." When evaluating tests in the real world, however, we need to be cautious of such statements, and try to examine the entire environment, considering test results with skepticism, and try to understand their limits—as well as their results.

For Further Reading

- [1] V. Manral, R. White, A. Shaikh, "Benchmarking Basic OSPF Single Router Control Plane Convergence," RFC 4061, April 2005.
- [2] V. Manral, R. White, A. Shaikh, "OSPF Benchmarking Terminology and Concepts," RFC 4062, April 2005
- [3] V. Manral, R. White, A. Shaikh, "Considerations When Using Basic OSPF Convergence Benchmarks," RFC 4063, April 2005.

RUSS WHITE works for Cisco Systems in the Routing Protocols Deployment and Architecture (DNA) team in Research Triangle Park, North Carolina. He has worked in the Cisco Technical Assistance Center (TAC) and Escalation Team in the past, has coauthored several books on routing protocols, including *Advanced IP Network Design*, *ISIS for IP Networks*, and *Inside Cisco IOS Software Architecture*. He is currently in the process of publishing a book on BGP deployment, and is the co-chair of the Routing Protocols Security Working Group within the IETF. E-mail: riw@cisco.com

Book Review

Running IPv6 *Running IPv6*, by Iljitsch van Beijnum, ISBN 1-59059-527-0, Apress, 2005. <http://www.apress.com/>

I've read a lot of books about emerging standards that read like "How I spent my summer vacation at a Standards Body." *Running IPv6* is *not* one of those. While van Iljitsch van Beijnum has been an active part of the IPv6 standards community, he has clearly done the homework of making it all work together. Weighing in at a compact 265 pages, *Running IPv6* really gets right to the point. The reader is assumed to have a working knowledge of IPv4.

Organization

The book starts off with a fairly typical introduction that explains why the author believes IPv6 is necessary. I find such introductions tedious, because if you've already forked out US \$44.95 for the book, the chances are that you're already motivated enough. This is, however, the only tedious chapter in the book.

What follows is a well written and organized primer for network administrators that covers how to configure end hosts, how get address space allocated, set up tunnels, and configure routers and the *Domain Name System* (DNS). The author covers in detail Linux, Windows, MacOS, Cisco's IOS (as well as that of other routing vendors), and Bind. We next move on to applications, IPv6 internals, transition strategies, and transit services.

Throughout, van Beijnum provides practical tips and advice on some of the pitfalls he found so the reader can avoid them. I particularly liked one case of whether to use eui-64 for the lower 64 bits of the address, pointing out the conflict between reducing configuration information (a good thing) and reduced readability (a bad thing).

The book primarily highlights differences between IPv4 and IPv6. This is important because it helps competent IPv4 administrators build on their existing knowledge. I know the last thing I want read about is how routing works when routing itself hasn't changed between versions. And I enjoyed reading, for instance, how *Dynamic Host Configuration Protocol Version 6* (DHCPv6) and stateless address auto-configuration differ from DHCPv4. I did not need nor want a primer in DHCP, but I did want to know about prefix delegation, which is not present in DHCPv4.

The author wastes no time on fluffy protocol niceties. Who cares, for instance, *how* a flow identifier is selected? What's important is that firewalls of the future may take advantage of it to determine flow direction, a major advance. Packet formats and semantics are only provided as they are needed by engineers to determine whether each component is performing correctly. The book is perhaps, therefore, best commended for what it lacks.

Unfortunately it lacks some subject matter I would like to have seen. Although van Beijnum covers how some common user applications, such as *telnet*, *ftp*, Web browsers and servers, and media players can use IPv6, business applications folks will be disappointed as there is no discussion of Oracle, SAP, or the like. The same is true for network management applications. And this may be a key roadblock to deployment of IPv6, as no self-respecting IT manager would deploy a service that cannot be managed. Such an obvious absence begs the question of whether those applications are IPv6 capable. On the bright side, you can try just about everything mentioned in the book because just about every tool mentioned either comes with the operating system or is freely available on the Internet. This book is not just theory.

A Must Read

It therefore shouldn't surprise anyone that I consider *Running IPv6* a "must read" for network engineers who have not yet played with IPv6. Even though Network Management Systems and business applications aren't covered, necessary protocol internals, semantics, operations, and troubleshooting are covered, therefore giving the reader a good knowledge base.

—*Eliot Lear, Cisco Systems, Inc.*
lear@cisco.com

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the "networking classics." In some cases, we may be able to get a publisher to send you a book for review if you don't have access to it. Contact us at **ipj@cisco.com** for more information.

A Pragmatic Report on IPv4 Address Space Consumption

Ole,

Thanks for a great round up in IPJ Volume 8, No. 3 on IPv6. This really helps focus where the state of the discussion needs to be in terms of addressing IPv6 deployment. You might be interested to know that this edition of the IPJ received tremendous interest in the UK. Within 24 hours of it arriving on your website, it was being distributed widely by several mailing lists serving communities from the *Ministry of Defence* (MoD) to important communications industry membership organisations. I received it myself at least three times from different lists!

Over recent months, I've seen a continuing trend to try to sideline IPv6 as not relevant to a particular discussion. IPv6 is either too low level for applications providers to think about, or too far off, or doesn't support some essential infrastructure service today. Some communities feel they have more than adequate IPv4 addresses to meet their foreseeable needs. These factors continue to drive debate on "if ever" rather than on "when" and "how" to deploy. That is, if the debate happens at all. All those who are investing in future IP-related services and networks need to read this edition of the *Internet Protocol Journal* for a reality check.

Tony Hain's article provides a compelling addition to the work you've already published by Geoff Huston on the analysis of IP address allocation, and is important food for thought that I think justifies increasing the urgency with which IPv6 support is treated. The discussion you hosted between Tony, Geoff with John Klensin and Fred Baker I think dealt very clearly with why the debate needs to be focused on the *how* and the *when* rather than on the *if*.

In the UK, we are seeing some significant investments made to enable IP level infrastructure with the intent of delivering profoundly new services into the twenty-first century, but none of these major investments appears to have included a vision for IPv6. So I think the point that was made concerning the current failure in making like-for-like investment decisions between v4 and v6 is hugely important for Chief Information Officers and Chief Financial Officers to take to their boards, or we will continue to find people investing for the past, rather than as they apparently believe, their future.

—Christian de Larrinaga
cde1@firsthand.net

Ole,

The analysis undertaken by Tony Hain and debated by some recognised experts makes it abundantly clear that the deployment of IPv6 is an immediate natural growth path to sustainability and global mass-market penetration of the Internet, beyond its worldwide current rate of less than 15%.

Tony has presented his study in the recent *IPv6 Forum Summits* (Seoul, Taipei, San Jose and Canberra) and obviously took a lot of people by surprise as previous studies maintained the suspense that the deployment of IPv6 should be an incremental transition and not an imminent and real migration. It was therefore decided to responsibly and morally act on this and renew a global Call to Action to set 2008 as a milestone of inevitable smooth transition in a softer form as a Y2K or Yv4 (The Year when IPv4 addresses will become hard to get) and get engineers to plan for it.

A global worldwide press release was published October 11, 2005 and can be read on the web site of the IPv6 Forum:

<http://www.ipv6forum.org>

The IPv6 Forum would like to recognise the work of *The Internet Protocol Journal* in watching diligently this space for the past couple of years and for initiating and orchestrating the constructive and consensual debate included at the end of the study, a contribution we trust is of great significance to the global good of the Internet.

—Latif Ladid, IPv6 Forum President
latif.ladid@village.uunet.lu

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 2005 Cisco Systems Inc.
All rights reserved.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA

The Internet Protocol Journal

March 2006

Volume 9, Number 1

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
Autonomous System Numbers	2
Working with IP Addresses	24
Letter to the Editor	35
Fragments	37

Autonomous Systems Numbers (ASNs) play an important role in the routing architecture of the Internet. An *Autonomous System* (AS) is, according to RFC 4271, "... a set of routers under a single technical administration, using an *interior gateway protocol* (IGP) and common metrics to determine how to route packets within the AS, and using an inter-AS routing protocol to determine how to route packets to other ASs." AS numbers are—like IP addresses—a finite resource, and predictions exist for when the AS number pool will be depleted. In our first article, Geoff Huston explains how ASNs work, and introduces us to the 4-byte ASN scheme that will allow for future growth beyond the currently predicted depletion date.

Our second article looks at another aspect of Internet routing and addressing—the IPv4 number space itself. Designers and operators of internets are often required to perform various address calculations in order to properly configure their networks. Russ White takes us through several exercises and introduces some "tricks of the trade" to make such calculations easier.

Our articles on spam in the last issue of IPJ prompted some feedback from our readers, and promises of more articles from other authors. This problem space clearly has more than a single solution. We look forward to bringing you more coverage of this topic in future editions.

The second issue of the *IETF Journal*, published by the Internet Society, is now available. Some people have asked me if I think of this new journal as a "competitor" to IPJ. I am happy to say that the *IETF Journal* is very much complementary to IPJ and covers important news from the IETF that we hope our readers will find interesting. You can access the *IETF Journal* by visiting: <http://ietfjournal.isoc.org>

The *IPJ Reader Survey* will soon close. We are grateful to the many readers who took the time to tell us about their reading habits, ideas for future articles, and other suggestions. Of course, we always welcome your feedback on any aspect of IPJ. Just drop us a line via e-mail to: ipj@cisco.com

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Exploring Autonomous System Numbers

by Geoff Huston, APNIC

So what are *Autonomous System Numbers* (ASNs), and what role do they play in the technology of the Internet? This article explores the role of ASNs as a critical element of the Internet routing architecture. We will first explore how the AS number space is structured, examine how ASNs are used in the interdomain routing environment and then look at the consumption rate of these numbers, and finally examine our options when we get to the point of likely ASN pool exhaustion. However, in order to put this into context, a brief overview of Internet routing architecture follows.

Internet Routing Architecture

Internet routing architecture is structured as a two-level hierarchy. The environment is first partitioned into *domains* with each domain using an internal routing environment. These network domains use an interior routing protocol (commonly referred to as an *Interior Gateway Protocol* [IGP]), which maintains a complete mapping set for the current internal topology of the domain, together with the set of “best paths” between any two points within the network domain. Although this approach of having a routing protocol automatically maintaining a comprehensive view of the current topology can be made to work within even quite large routing domains, such an approach does not scale to the size of the entire Internet. Fine-grained topology information is useful only in “local” situations, and is best omitted when forming a larger view of the network. Commonly used interior routing protocols include *Open Shortest Path First* (OSPF), *Intermediate System-to-Intermediate System* (IS-IS), and *Enhanced Interior Gateway Routing Protocol* (EIGRP).

The second level in the routing hierarchy is the *interdomain* routing domain. The interdomain routing environment describes how domains interconnect, but avoids the task of maintaining transit paths within each domain. In the interdomain space, a routing path to an address is described as a sequence of domains that must be transited to reach the domain that originates that particular address prefix. Today this interdomain space is maintained using Version 4 of the *Border Gateway Protocol* (BGPv4).

Each routing domain is a single administrative domain, operated within a uniform set of routing policies, and is operated independently from any other domain. The domain is in effect an autonomous unit in the overall routing architecture, and is termed an *Autonomous System* (AS). Each of these ASs is uniquely identified using an *Autonomous System Number* (ASN).

What Is an Autonomous System?

One of the best definitions of an Autonomous System can be found in an IETF document, RFC 4271^[4] that describes BGPv4:

“The classic definition of an Autonomous System is a set of routers under a single technical administration, using an *interior gateway protocol* (IGP) and common metrics to determine how to route packets within the AS, and using an inter-AS routing protocol to determine how to route packets to other ASs. Since this classic definition was developed, it has become common for a single AS to use several IGPs and sometimes several sets of metrics within an AS. The use of the term Autonomous System here stresses the fact that, even when multiple IGPs and metrics are used, the administration of an AS appears to other ASs to have a single coherent interior routing plan and presents a consistent picture of what destinations are reachable through it.”

The AS Number Pool

ASNs are drawn from a 16-bit number field, allowing for 65,536 possible values.

AS 0 is reserved, and may be used to identify nonrouted networks. The largest value—AS 65,535—is also reserved. The block of ASNs from 64,512 through 65,534 is designated for private use. ASN 23,456 is reserved for use in ASN pool transition. The remainder of the values, from 1 through to 64,511 (less 23,456), are available for use in Internet routing. The number space is unstructured, because there are no internal fields in the number structure, nor is there any aggregation or summarization capability for ASNs.

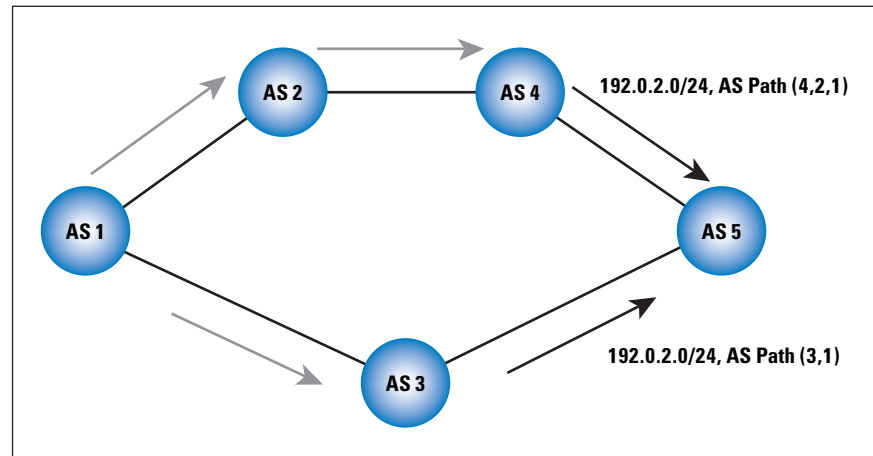
How AS Numbers Are Used in BGP

The interdomain routing space is constructed using two components: *address prefixes* and *AS numbers*, which are used as domain identifiers. Every prefix has an originating domain, known as the *Origin AS* from which reachability for the prefix is propagated across the interdomain space.

As the routing advertisement is propagated across the interdomain space, each prefix accumulates an associated “AS path.” When an address prefix advertisement transits a domain, the domain effectively “signs” the prefix advertisement by prepending its ASN to the AS path associated with the address prefix. At any point in the network the AS path describes a sequence of connected domains that forms a path from the current point to the originating domain. This setup is shown in Figure 1, where AS1 originates an advertisement for the address prefix **192.0.2.0/24**. At AS5, the AS receives two BGP advertisements for this prefix. One has the AS path (4, 2, 1), and the other has the AS path (3, 1).

The left-most number in the AS path list is the ASN of the adjacent AS from which the address prefix advertisement was received. The sequence of numbers indicates the sequence of ASs through which this update was propagated. The right-most, or final ASN, is the AS number of the AS that originated the address prefix advertisement, or *Origin AS*.

Figure 1: AS Path Generation in BGP



The AS path serves two purposes in interdomain routing: that of a *path length* metric and a *loop detection* mechanism.

The AS path is used as a path metric in the BGP path selection algorithm. When a domain receives two different BGP advertisements for the same address prefix, the default BGP selection process is that of selection of the advertisement of the minimal-length AS path, with each AS in the path counting as a single unit of “cost.” In the case of the example network in Figure 1, AS5 prefers to use the path through AS3 to reach the originating AS1, in preference to the longer path of AS4 and then AS2.

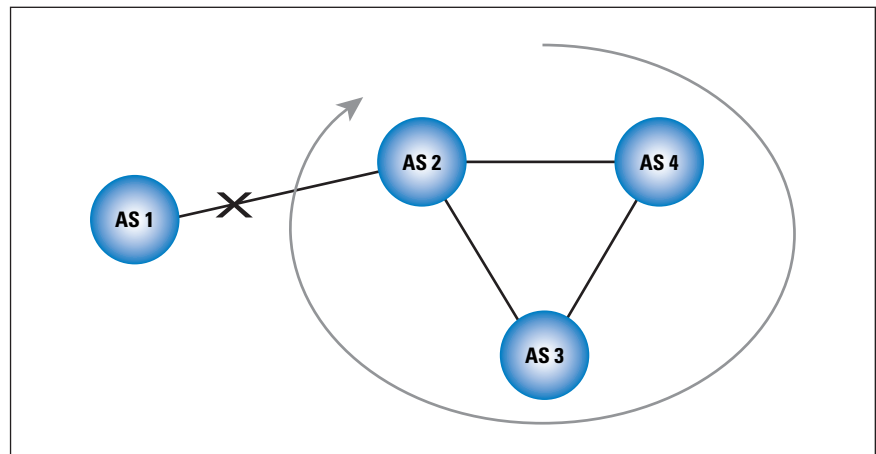
Although enumerating the AS path vector within the routing protocol is one way of passing the path cost through the routing domain, it may appear that the best path selection function could just as easily be supported by carrying a simple path cost metric of a domain transit counter, similar to that used by other distance vector routing protocols, such as *Routing Information Protocol Version 2* (RIPv2). However, the problem with distance vector protocols is the “count-to-infinity” dilemma.

To illustrate the need for explicit AS path enumeration in BGP, consider what happens when the AS path vector is replaced by a simple path cost metric. In the configuration shown in Figure 2, AS1 originates a routing advertisement toward AS2. AS2, AS3, and AS4 are interconnected in a simple loop configuration. When AS2 receives AS1’s advertisement with a path cost of 1, it passes the advertisement on to both AS3 and AS4, with a path cost of 2. Both AS3 and AS4 select as their best path this advertisement from AS2 with a path metric 2, corresponding to the AS path (2, 1).

Now if the connection between AS1 and AS2 is broken, then AS2 no longer sees AS1, and withdraws its best path to the prefix through AS1. AS2 then stops advertising a path to AS3 and AS4. But AS3 is already advertising a path to AS4, with a metric of 3, corresponding to the AS path (3, 2, 1). Upon the withdrawal of the advertisement from AS2, AS4 then selects this as its next best path, with a path cost of 3. AS4 then advertises this prefix to AS2 with a path cost of 4, corresponding to the AS path (4, 3, 2, 1).

At this point, without the explicit AS path in the advertisement, AS2 cannot deduce that this advertisement is, in fact, a loop. Accordingly, AS2 accepts this path with a metric of 4 as its best path. AS2 then advertises this to AS3 with a metric of 5, corresponding to the AS path (2, 4, 3, 2, 1). AS3 updates its best path to AS1 with this new metric and then sends an update to AS4, and so on. This process continues around the loop until the path cost metric reaches some defined maximal value. The higher the maximal value for the path cost metric, the longer the time taken to detect the loop condition. The smaller the maximal path cost metric, the smaller the span of network that the protocol can encompass. Setting the maximal path cost parameter requires some considerable care, and the operation of the protocol can be extremely slow to converge in terms of loop detection.

Figure 2: Loop Formation in Distance Vector Protocols



This form of loop can be averted by replacing the path cost counter with a fully enumerated AS sequence. Continuing the example in Figure 2, when AS2 withdraws its route to AS3 and AS4, AS4 still selects the other route it has heard, but this time the selected prefix has the path (3, 2, 1). When AS4 attempts to pass this advertisement to AS2, AS2 sees its own value in the associated AS path and rejects the advertisement. At the same time AS3 withdraws its advertisement to AS4, and at that point the prefix is dropped from the entire routing system. In this way the AS path acts as an efficient routing loop detector.

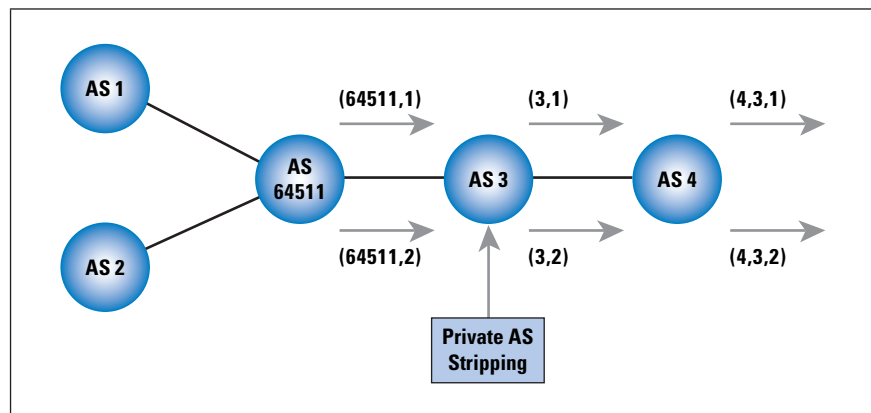
The use of ASNs and AS path vectors in BGP provides an effective solution to this classic problem of loop detection, as well as providing a simple and effective path-selection process.

Who Needs an AS Number?

Not every network needs to have its own ASN. The guiding principle is that ASNs are used to express distinct interdomain routing policies, and not every network has the requirement to express its own unique set of routing policies.

In the case where a network has a single upstream connection, the routing policies of the network are precisely the same as those of its upstream service provider, and there would normally be no need for the network to use a distinct ASN. Even if the network domain uses BGP for its upstream connection, the originating domain can use a private ASN (from the number range 64,512 – 65,534) to support the BGP session to the upstream network. The upstream network strips off the private ASN when it readvertises the prefix, and the upstream network appears to the rest of the Internet as the originating AS. Even if the AS has “downstream” networks it can still use a private AS, even when the downstream ASs are using public ASNs. The stripping of the private AS removes only the instances of the private AS from the AS path, and not the public ASNs (Figure 3).

Figure 3: Use of Private AS Numbers



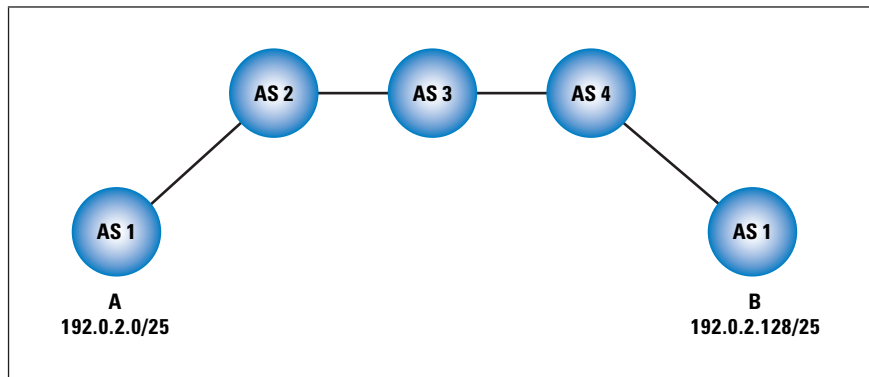
In the case where a network has two or more upstream transit connections, it is more likely that the network will use its own unique ASN. It is not always the case that a distinct ASN is required here, and the distinguishing factor is that of the network wanting to express particular routing policies. Where the network has no particular preference as to which of the upstream services should be used for incoming traffic, the network can also use a private ASN for each of its routing sessions. In such a case the external routing view would be that the prefix appears to be originated from multiple ASs.

In the case where there are multiple paths to reach the network, and where these paths need to be distinguished in the routing system by different AS paths that have the same originating AS (that is, there is a need to express a routing policy), then the network needs to use a unique ASN within the interdomain routing system.

Can an ASN Be Split Across Separated Subdomains?

There are many cases of dispersed networks that exist in multiple locations. If these locations are all administered by a single entity, it may be desirable to use a single ASN across all these domains. This scenario is possible, but considerable care needs to be exercised when designing the routing configuration. Figure 4 shows two distinct subdomains of AS1, and they are not interconnected internally.

Figure 4: Split AS



AS1 (A) advertises the prefix **192.0.2.0/25** to AS2, and this advertisement is propagated to AS2, AS3, and AS4. When AS4 passes this advertisement to the other segment of AS1 (B), this router rejects the advertisement because the associated AS path (4, 3, 2, 1) indicates that the route has already passed through AS1. Similarly, the first segment of AS1 (A) rejects the advertisement of **192.0.2.128/25** from AS2, because its path (4, 3, 2, 1) also indicates that a loop has formed. To restore complete connectivity between the distinct parts of AS1, AS1 needs to configure static routes at its edges. If AS1 (A) configures a static route to **192.0.2.128/25** pointing toward AS2, and AS1 (B) similarly configures a route to **192.0.2.0/25** through AS4, then the configuration enables full connectivity.

In more complex configurations where each of the segments of the network is multiply connected, the static route configuration becomes more complex. However, with very careful configuration, a single ASN can be distributed across multiple distinct networks.

AS Path Prepending and Path Poisoning

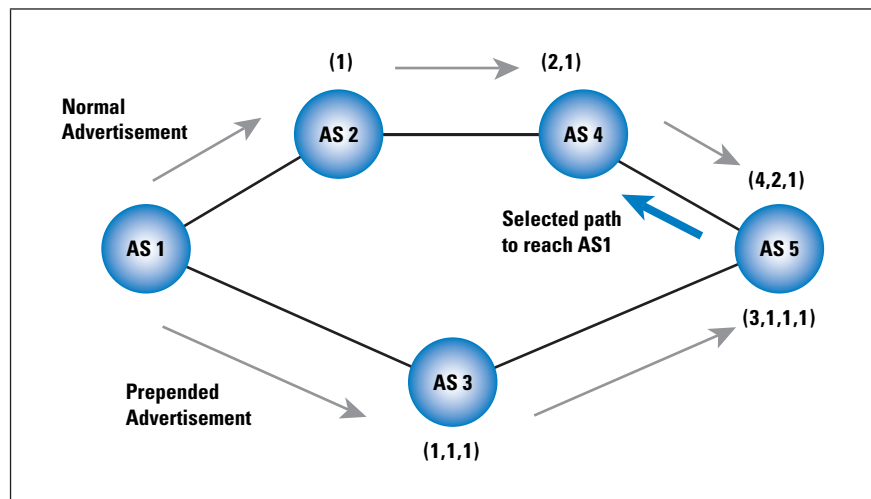
The basic mechanism of path preference in BGP is that of the *AS path length*. Where there are two advertised paths to reach a particular address prefix, the default selection algorithm in BGP is to prefer the advertisement with the *shorter* AS path length.

A multihomed domain may wish to have other domains prefer one particular path over another to reach it. This may be because the local domain wishes to optimize its traffic costs between the multiple upstream providers, balance the traffic load across multiple paths, or set up various forms of primary and backup relationships across the multiple provider upstream paths.

Although such policy preferences are often set up using BGP *communities*, BGP community signaling requires the cooperation of multiple parties in consistent interpretation of the community values. A more coarse form of expressing such policy preferences can be achieved through *AS path prepending*, a technique of deliberately extending the AS path length of a prefix advertisement by adding additional ASNs into the AS path of an advertised prefix. Normally the form of AS path prepending uses the local ASN to perform the prepending.

In the example in Figure 5, AS1 wants to express the policy to prefer incoming traffic through AS2, and use the link to AS3 only as a backup. To achieve this with AS path prepending, AS1 prepends itself twice in the AP path of the advertisement passed to AS3, in order to artificially lengthen the AS3 transit path. AS5 would have normally used the shorted AS path through AS3 to reach AS1. As a result of AS1 artificially lengthening its path to AS3, AS5 now selects the transit path through AS4 and AS2 to reach AS1.

Figure 5: AS Path Prepending

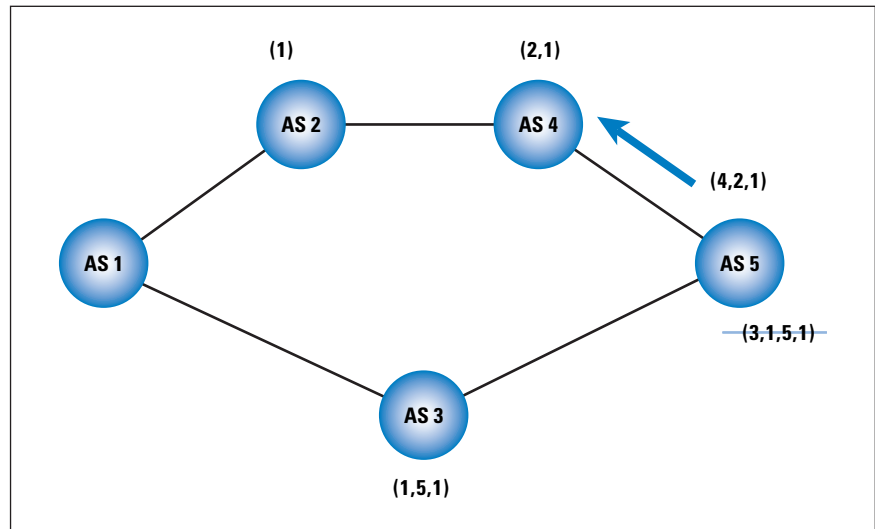


Of course AS path prepending is a very imprecise technique, and can often produce surprising results in real-world situations. A more deterministic method of traffic engineering uses additional signals attached to address prefix advertisements, through BGP communities.

A more subtle, and more controversial, prepending technique is that of so-called *AS path poisoning*, where an AS uses some other value to prepend in the AS path. In Figure 6, AS1 wants to express the policy that under no circumstances should AS5 use the transit through AS3 to reach AS1. In this case AS1 could use AS5 as the prepending value in its advertisement to AS3.

When AS5 receives this advertisement, the presence of its own ASN in the AS path means that it will not accept this advertisement, and prefers the transit path through AS4 and AS2. The difference between these two examples is that in the case where the connection between AS1 and AS2 is broken, none of AS2, AS4, or AS5 can possibly reach AS1 when this AS path poisoning technique is being used.

Figure 6: AS Path Prepending with AS Path Poisoning



AS Number Consumption

In this section we will look at the rate of consumption of ASNs, and estimate when they may be fully consumed. Of the 64,510 available AS numbers, as of January 2006 we have already allocated some 40,000, or well over half of the number pool. Two immediate questions arise—how long do we have before the number pool is completely exhausted, and what are our options for an expanded number pool that can encompass a larger interdomain routing environment?

The Factors for AS Number Consumption

Before looking at these two questions in further detail, it would be useful to understand the factors that affect AS number consumption.

From one perspective it is counterintuitive to assume that the Internet will evolve from tens of thousands of distinct routing domains to one of hundreds of thousands or even millions of distinct routing domains. It may appear that there is a reasonable level of correlation between the number of active *Internet Service Providers* (ISPs) in the Internet and the number of advertised ASNs. If forecasting a future demand for hundreds of thousands or even millions of ASNs, it would appear that we are forecasting continued fragmentation of the service provider industry with large numbers of small enterprises that, collectively, compose the Internet. This scenario does not appear to be likely.

The ISP industry is one with an underlying factor of economies of scale. Larger ISPs generally have access to more efficient use of resources and are more capable of sustaining a market share at competitive prices, with reasonable operating margins because of these economies of scale. Smaller providers tend to service niche markets, and in general are highly susceptible to pricing pressures in the competitive supply market. The overall result is strong pressure for continued aggregation in the service provider market, tending to aggregate to a smaller number of larger providers.

If the number of ASs in use is roughly commensurate to the number of service providers, then this view of the market dynamics would lead to a view that the service provider population is either in a state of equilibrium where the entrance of new niche-oriented players is much the same as the rate at which smaller players are aggregated into larger providers, or one of relatively small growth based on the larger dynamics of continued expansion of the Internet on a global basis.

In practice this has not been the case, and we see a continuous rate of consumption of new ASNs. This rate appears to be some 3,500 ASNs per year, and this consumption rate appears to have been steady since 2002 (see Figure 7). Accordingly, it appears that some additional factors affect AS number consumption rates.

One of these factors is the practice of *multihoming* at the edge of the network. Many end-site networks have business-critical needs for assured Internet connectivity, and a common way to achieve this connectivity is by using the services of two or more upstream providers. In such situations the end site may want to express different routing policies to each upstream provider, and it does so by using its own ASN and expressing these routing policies using BGP to each of its upstreams.

AS numbers are also used in other contexts. In *Multiprotocol Label Switching* (MPLS) Layer 3 networks, one form of generating the *Route Distinguisher* value for a VPN client network is through the use of concatenating the VPN host's AS number with a serial number. To what extent this semiprivate use of AS numbers in a VPN context contributes to the consumption rate of ASNs is difficult to assess, simply because the use of these numbers is not generally visible.

Even within the public Internet there are other contributory factors to AS number consumption. ISPs with diverse product portfolios may wish to express different routing policies for various product families, or express different routing policies in different regions of network coverage. Again this can be achieved through the use of distinct AS numbers of each routing policy set.

An associated contributory factor for AS Number consumption is that there is little incentive for AS Number return and recycling. With the current framework there is no direct cost to maintain an AS number allocation, and the overall characteristic of AS number allocation appears to be a “once and forever” allocation model. When AS numbers are no longer required, AS numbers generally do not return to the unallocated pool for subsequent reallocation.

Taken together, these factors lead to the conclusion that continued AS number consumption is based on a larger set of considerations than the dynamics of the service provider industry.

Accordingly, we can be a little more confident in making the assumption that the factors that have affected AS number consumption in the recent past will continue to be factors in the near-term future, leading to some further confidence in a predictive technique that uses recent consumption data to generate trends that can made predictive forecasts of future demands. We will apply this technique to AS number consumption data to make some forecasts of the time by which the current AS number pool will be exhausted.

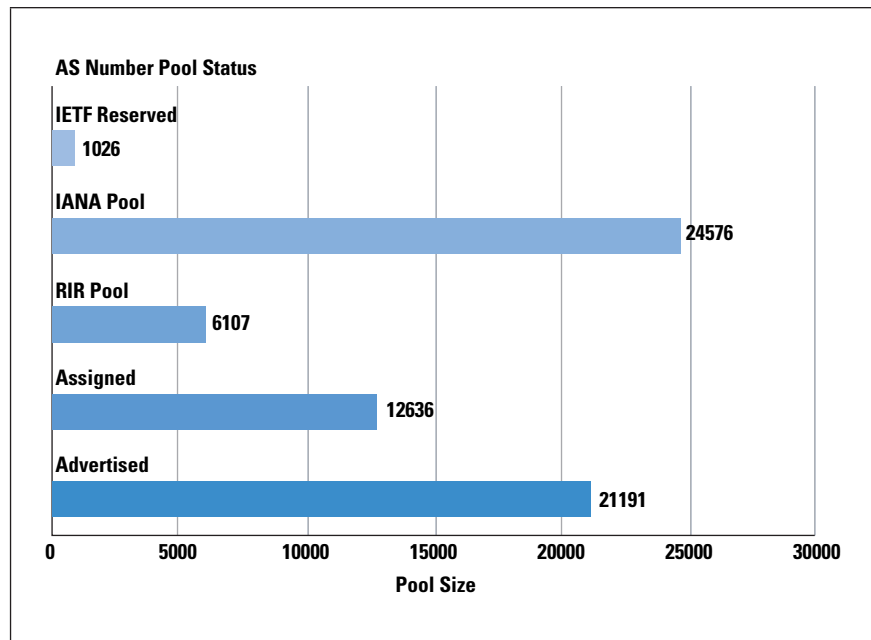
AS Number Pool Status

There are 65,536 AS numbers. As noted already, some 1,026 numbers are reserved and unable to be used in the public Internet, leaving 64,510 for use in the public Internet.

The pool of AS numbers is administered by the *Internet Assigned Numbers Authority* (IANA), and blocks of 1,024 numbers are allocated to the *Regional Internet Registries* (RIRs) periodically when the RIR’s pool drops below a threshold level.

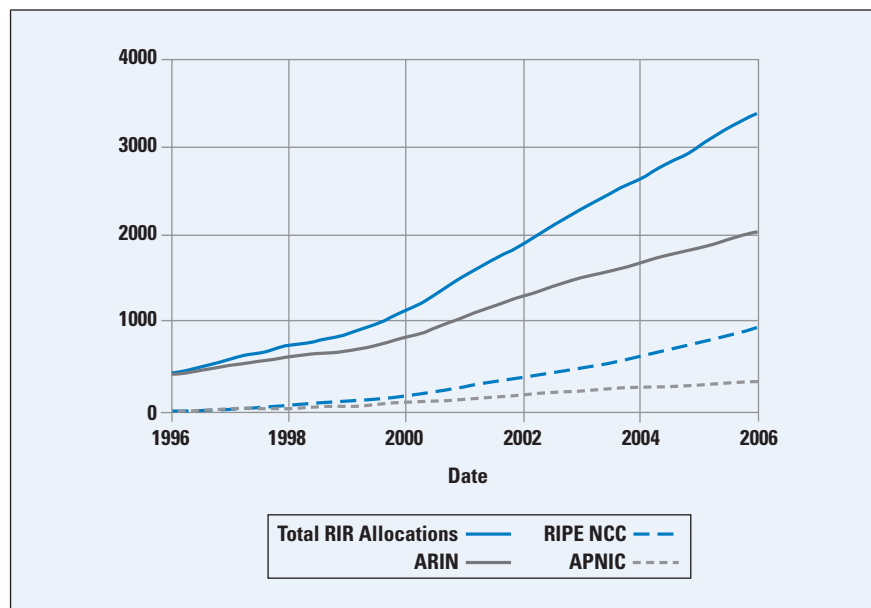
Of the 39,934 AS numbers that have been allocated by IANA by January 2006, there is a further classification of AS numbers. A working pool of numbers is held by the RIR for current assignment to ISPs. Of the assigned AS numbers, some are visibly used in the interdomain routing table of the public Internet, but others are not visible in the Internet. The breakdown of AS numbers into the RIR pool, assigned but not advertised, and assigned and advertised, as of January 2006, is shown in Figure 7. Of the 34,827 assigned AS numbers, some 21,191 are advertised; 12,636 have been allocated in the past, but are not currently advertised in the BGP routing table.

Figure 7: AS Number Status of
Advertised, Unadvertised,
and Unallocated Pools



The RIRs allocate ASNs to ISPs and end-user networks. A second time series can be generated, showing the cumulative sum of the RIR AS allocations (Figure 8). Not surprisingly, the time series shows the effects of the Internet boom across the period from 1999 through to late 2001 as a sharp upward trend in allocations. The subsequent market correction is also evident as a visible change in the AS allocation rate by early 2002.

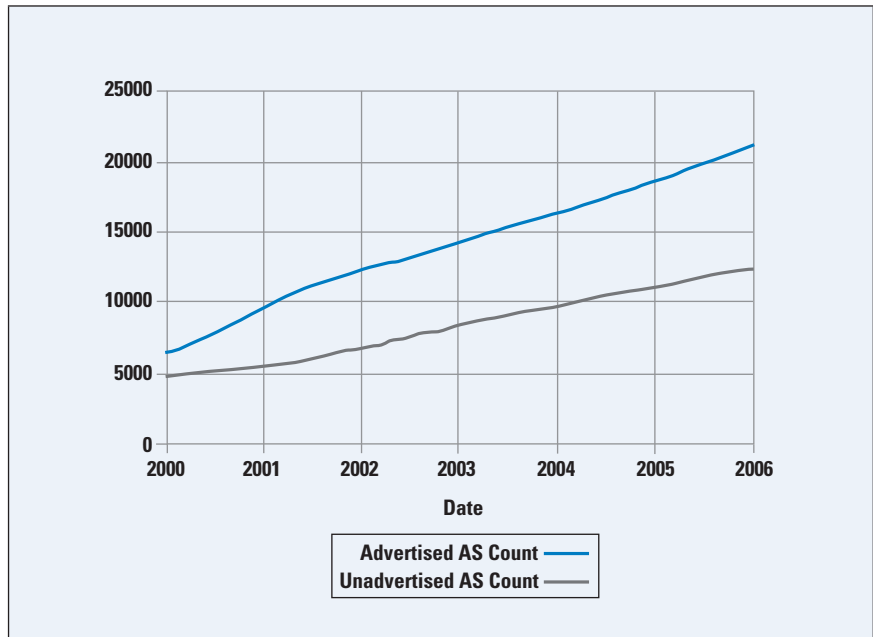
Figure 8: RIR Allocations



BGP AS Advertisements

In addition to allocation rates, a further source of ASN data is the interdomain routing table. The number of distinct ASs advertised in the interdomain routing space of the public Internet has been measured regularly since 1997. The time series of this count of advertised ASNs, and the complementary number of unadvertised ASNs, is shown in Figure 9.

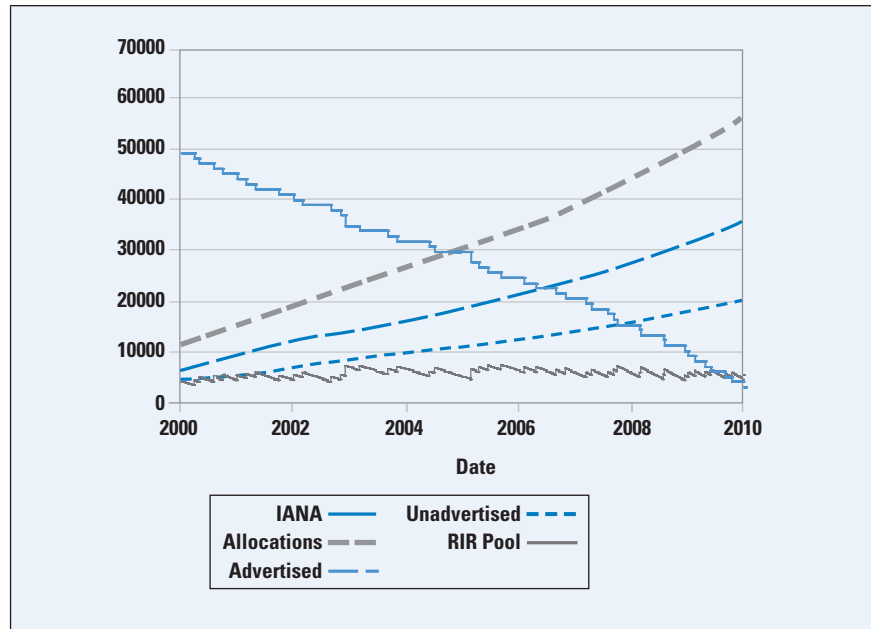
Figure 9: Advertised and Unadvertised AS Numbers



AS Number Consumption Projections

At this point, it is possible to make some projections on AS number consumption. The technique here is to use the past three years' consumption data (taking a starting point of January 2003) and derive an associated exponential function as a best fit to the 3-year data series in order to generate a trend function. This trend function is then projected forward in time to forecast the point in time when the resource reaches a certain threshold point. A considerable amount of detail is associated with this exercise, including the use of an exponential function as the best fit to the past 3 years' ASNs use rates (see <http://www.potaroo.net/tools/asns/>). However, for the purposes of this article it is appropriate to proceed to the outcome (Figure 10).

Figure 10: A Predictive Model of AS Number Consumption



From this model it appears that we are looking at steadily accelerating consumption of ASNs, and a projected date of late 2010 of exhaustion available AS numbers to allocate to ISPs.

The implication is that this model indicates that by late 2010 either the Internet should be using a new protocol for interdomain routing that does not rely on AS numbers at all, or, more likely, that the Internet should be using a version of BGP that supports the use of larger AS numbers that are drawn from a number pool significantly larger than 16 bits. The first option appears to be somewhat unrealistic, to say the least. And the second option, although simpler and very much the preferred path, is still going to take some time to deploy, particularly considering the growing size of the interdomain space of the Internet and the diversity of these component domains.

When contemplating a transition to a larger ASN pool, it should be remembered that every day there are more networks that will need to undertake a transition to a longer ASN field in their deployed instances of the BGP protocol.

The steps in this transition path appear to include:

- The completion of the relevant protocol standards for a larger ASN field in BGP
- The production of code in available implementations of BGP that support this protocol standard
- Various forms of testing this code, both in terms of its correct operation and interoperability and in terms of the correctness and viability of the relevant transition steps

- Developing the necessary infrastructural support system to manage the distribution of this new number pool
- A process of deployment of this protocol so that the deployment of larger ASNs can commence well before the point at which the existing AS number pool is exhausted

Even an aggressive schedule of transition across such a large and diverse network as the Internet will take many years to reach the final step. It also appears that a prudent course of action would see us reach that position not by 2010, but by 2008 at the latest, allowing us a margin of some 2 years (and some 10,000 remaining AS numbers) to complete the task.

32-Bit AS Numbers

In this part of the article we will look at the current proposal for a larger AS number pool. As of October 2005, the document defining this proposal is an IETF Internet Draft: **draft-ietf-idr-as4bytes-12.txt**. The proposed approach is to expand the size of the AS number pool space from 16 to 32 bits. In number terms this expands the number space from a pool of 65,536 numbers to 4,294,967,296 numbers. In terms of the current use of ASNs, the current scaling properties of the BGP routing protocol, and the use of ASs in the context of interdomain routing, a pool of some 4.3 billion numbers would easily encompass a network environment of significantly greater levels of domains, and interdomain interconnection density. Such a pool size would exceed some current guesses of the scaling capabilities of the BGP protocol by up to a further two orders of magnitude.

It is also proposed to preserve the first block of 65,536 32-bit ASNs to align with the allocations of the 16-bit numbers.

Let's use a new form of terminology here for 32-bit ASN values, where the first 65,536 ASNs are numbers that use the form "0.0" through "0.65535." The second set of 65,536 numbers would be written as 1.0 through to 1.65535, and so on. So here we will be using a number format of *<upper16 bits>.<lower 16 bits>*.

What is the inventory of concerns that need to specifically addressed in the transition to these 32-bit AS Numbers?

Obviously there is a need for some changes to the routing protocol, and an ordered interdomain transition is unrealistic to expect. More reasonable is an expectation of a piecemeal transition of domains, where individual domains transition their BGP platform to supporting 32-bit ASs in their own time. Domains that are currently using 16-bit ASs may have less reason to undergo an early transition to 32-bit AS support, whereas those domains that are assigned a nonmappable 32-bit ASN will find that they have to support 32-bit ASNs from the outset.

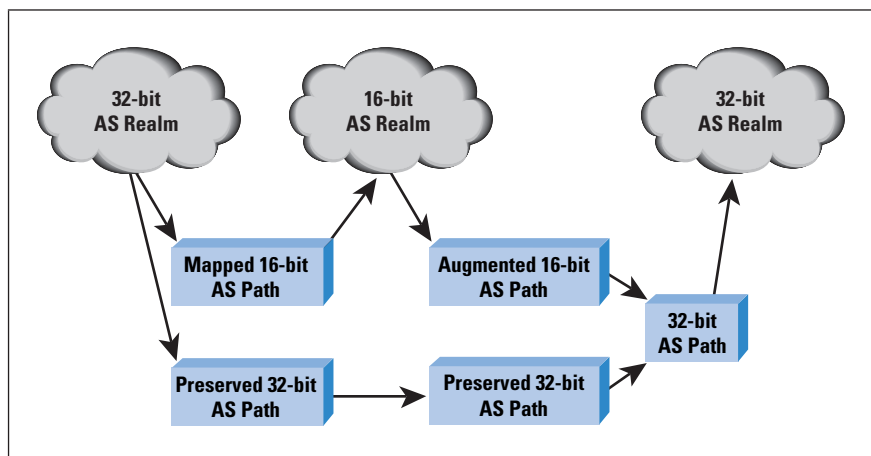
32-Bit Changes to BGP

BGP has two major parts within its protocol: opening up a BGP conversation with a peer BGP speaker, and then the transferring of protocol objects that describe reachability of address prefixes and associated attributes of these address prefixes. Both parts include AS Number components, and in considering changes to the current protocol, both parts of the protocol require some change. The message objects that need to be considered here are, therefore, the BGP OPEN message and the BGP UPDATE message.

The changes to BGP create a “NEW” BGP implementation that is capable of supporting a 32-bit ASN environment. The essential task of the changes is to define mechanisms that all NEW BGP speakers use to speak to each other and pass all ASN values in 32-bit fields. However, the Internet is way too large to set up a “flag day” at which point the entire collection of BGP speakers will undertake a switch from “OLD” BGP to NEW BGP. Accordingly, it is also necessary to define protocol interactions in NEW BGP where the transition in the Internet will be gradual and essentially uncoordinated. NEW BGP speakers will have to set up sessions with OLD BGP speakers, and of course OLD BGP speakers will also be peering with other OLD BGP speakers. The information associated with 32-bit AS paths must be passed across sections of the network that normally support only 16-bit AS paths. In other words, 32-bit AS information needs to be passed to OLD BGP speakers and between OLD BGP speakers.

The general approach adopted for transition is preserve AS path length information across the OLD and NEW BGP boundaries, while recognizing that some 32-bit AS information cannot be cleanly mapped into a 16-bit AS path. In order to preserve 32-bit information—a necessary step to prevent loop formation for 32-bit ASs—the 32-bit information is preserved across OLD transit paths and restored upon reentry into NEW BGP realms (Figure 11).

Figure 11: 16-Bit and 32-Bit AS Realms



Opening a BGP Session

The proposed approach is to initiate a NEW BGP session in a mode that is compatible with the OLD BGP protocol, and also inform the remote peer of its capability to conduct a NEW BGP conversation if the remote peer is also a NEW BGP speaker. NEW BGP speakers who open a peer session with an OLD BGP peer will ignore the NEW capability and operate their BGP peer session in OLD mode. A NEW BGP peer will respond positively to the NEW capability, and that BGP session can then operate in NEW mode.

The BGP OPEN message includes a fixed-length 16-bit *My_AS* field as well as potentially containing a capability query as part of the *Optional Parameters* section. In order to ensure that NEW and OLD speakers can communicate, this 16-bit *My_AS* field needs to be preserved in NEW BGP even when the Optional Parameters section includes the capability to undertake a NEW peering session. This may appear contradictory in the first instance, because the OPEN message then contains both a 16-bit ASN and a 32-bit *AS Capabilities Query*. The mechanism proposed for the OPEN message varies according to whether the NEW speaker is using a mappable ASN drawn from the original pool (that is, with a *My_AS* number in the range 0.0 through 0.65535), or it is using a number drawn from a higher-numbered 32-bit number block. In the first case the OPEN message would use the 16-bit mapped value in the *My_AS* field (dropping out the zero-valued high-order 16 bits of the AS value), whereas in the second case the BGP speaker would use for *My_AS* a special 16-bit value that is reserved for this purpose (AS 23456). In both cases the Optional Parameter section would include a capability code to indicate that the local BGP speaker can support 32-bit ASNs (Capability Code 65).

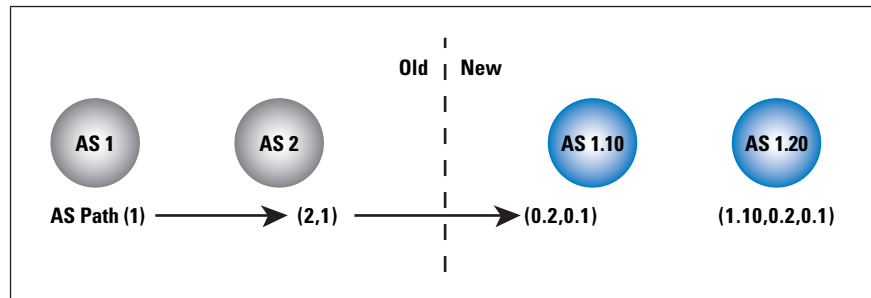
The side effect is that from the perspective of OLD BGP domains AS 23456 may appear to be connected to the interdomain network in many different locations. From the OLD BGP realm this does not present a protocol problem, although, as always, there is the potential here that this repeated use of AS 23456 as a 32-bit AS substitution token may create a somewhat confusing BGP view of the Internet from the perspective of the OLD BGP world.

The capability exchange uses a protocol described in RFC 3392. The NEW BGP speaker adds an optional capability field to the OPEN message. The 32-bit AS capability code 65 carries as its capability value the local 32-bit local ASN value. For a NEW peer this capability value is to be interpreted as the actual AS of the remote side, on the basis that the *My_AS* field in the body of the OPEN is either a truncation of the local 32-bit AS value (in the case of mappable 32-bit AS values), or the special value of AS 23456.

The BGP UPDATE Message

For a NEW BGP session (32-bit peering with 32 bits) the changes to the protocol are the use of 32-bit ASNs in the AS_PATH attribute of UPDATE messages. All 16-bit AS values are padded with a zero high-order 16 bits. If the AGGREGATOR attribute is used, it is similarly carried as a 32-bit value. So in the 32-bit peering, all 16-bit information is carried in mapped 32-bit ASNs (Figure 12).

Figure 12: OLD to NEW BGP
AS Path Mapping



In this way AS path length is preserved without change when translating 16-bit AS information into the 32-bit domain.

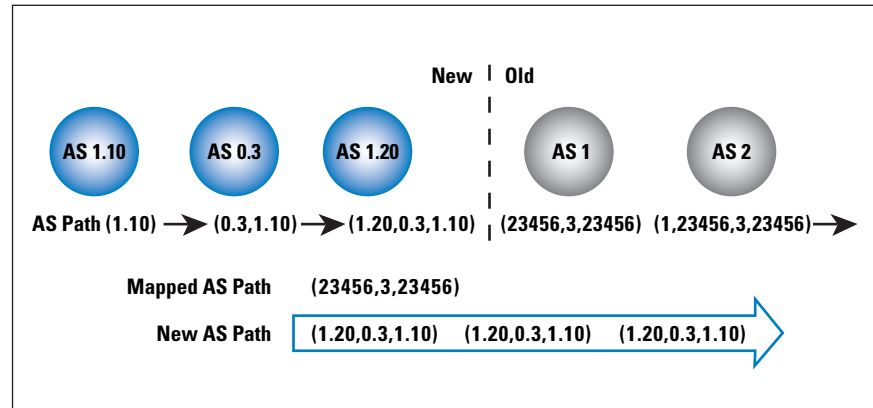
The next case is where an OLD BGP peers with a NEW BGP. We have already seen the simple case where the information is coming from a 16-bit path and there is no additional 32-bit information, and in this case the 16-bit values are simply mapped into 32-bit values, by padding the ASN values with 16 zero high-order bits. What about the reverse case where 32-bit information is being passed back into the 16-bit world?

This case has two parts: first creating an equivalent 16-bit AS path and second, packing up the 32-bit AS path information in such a way that it transits across the 16-bit domain in such a manner that it can be reassembled in any subsequent transition into a 32-bit domain. In the first case, the equivalent path information is constructed by stripping the high-order 16 bits off the AS value, as long as this part is all zeros. Where this is not possible—and the AS path contains one or more ASNs with non-zero high-order bits—then the transition ASN, 23456, is substituted in the place of each such ASN in the AS path. In this way the AS path length metric is preserved, and the prevention of count-to-infinity loops in the 16-bit domain is avoided.

The second part to this case is packaging up the 32-bit path into the OLD BGP session in such a way that it can be unpacked at any subsequent boundary back into a 32-bit routing realm. Here the proposal calls for new transitive community attributes to be carried in OLD BGP routing realms. These attributes are defined as transitive attributes, and should be passed through the OLD BGP peering sessions without alteration. It should be noted that this is not a protocol change as such, but it does require the explicit configuration support within OLD BGP implementations of this attribute as a transitive community.

The proposed mechanism is an extended community attribute called “NEW_AS_PATH.” When a NEW BGP speaker is speaking to an OLD BGP, the NEW BGP prepends its own AS value to the AS_PATH and copies this information into the NEW_AS_PATH attribute. It then translates the 32-bit AS path into a 16-bit equivalent AS path. The translation is straightforward, in that where the 32-bit AS has all zeros in the high-order 16 bits, the translation truncates the AS value to a 16-bit value, and where the high-order 16 bits are nonzero, the translation substitutes the reserved 16-bit value AS 23456 in its place (Figure 13).

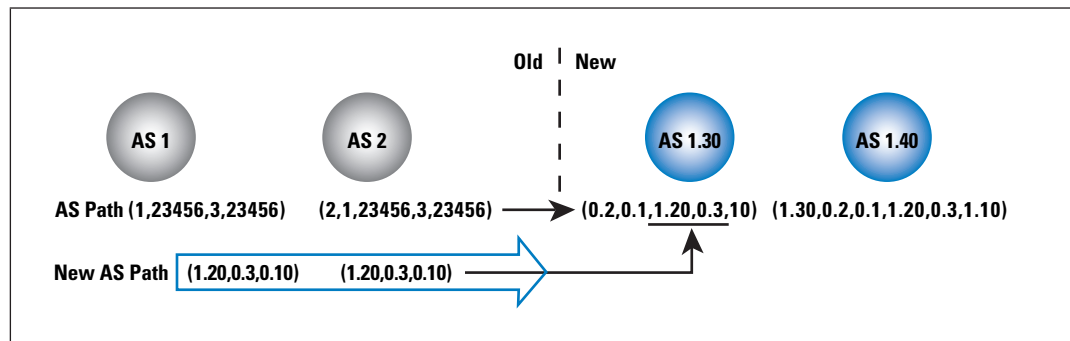
Figure 13: NEW to OLD BGP
AS Path Mapping



The transit across the OLD BGP domains leaves the NEW_AS_PATH untouched, and prepends 16-bit AS values to the AS_PATH. In other words, OLD BGP behaves as it always has. The NEW_AS_PATH is passed through the OLD realm as an opaque bit block.

The next transition is one from the OLD to the NEW domain when a NEW_AS_PATH attribute is present. In this case the NEW BGP speaker takes the AS path as presented by the OLD BGP speaker and converts the 16-bit AS values to 32-bit AS values by adding 16 bits of zero padding to each entry, as before. However, in this case the NEW BGP speaker then overwrites the trailing entries with the values specified by the NEW_AS_PATH attribute. The effective result is that the 32-bit AS path that entered the 16-bit sequence is prepended with the equivalent of the 16-bit transit AS sequence. The NEW_AS_PATH attribute is then removed from the BGP Update, leaving an intact 32-bit path as the AS_PATH attribute. This scenario ensures that the resultant BGP environment can accurately detect loops in both the NEW 32-bit and OLD 16-bit realms (Figure 14).

Figure 14: OLD to NEW BGP AS Path Mapping



What if there was a routing loop that traversed a mixed sequence of NEW and OLD routing realms? The restoration of the original 32-bit AS path at the OLD-to-NEW transition ensures that the potential loop is discarded, because a 16-bit AS sees its own AS in the 16-bit AS_PATH attribute, and a 32-bit AS also sees its own value in the 32-bit AS_PATH. The transition mapping ensures that the potential routing loop is detected by BGP.

The ability to perform AS path prepending is also unaltered in this mixed NEW and OLD BGP environment. The AS simply prepends its local AS value to the AS_PATH as usual. In the case of prepending on a NEW-to-OLD boundary, the prepended AS path is mapped into the NEW_AS_PATH attribute as described previously.

Earlier in this article we noted the less common use of AS path poisoning, where the prepending uses a different ASN value in order to ensure that the particular advertisement is not learned by a remote AS. For NEW BGP speakers there is no change to this capability. For OLD BGP speakers the AS path poisoning can be directed only toward 16-bit ASs, because the OLD BGP speaker has no knowledge of the structure or content of the NEW_AS_PATH attribute.

Another part of the BGP protocol that uses ASNs is the AGGREGATOR attribute. This attribute is attached to an update message when an AS combines two or more prefixes into a single aggregate prefix (a practice that is often referred to as “proxy aggregation”). The ASN of the aggregating AS is attached to the aggregate prefix advertisement as an AGGREGATOR attribute. The same ASN translation technique applies to AGGREGATOR attribute when an advertisement is passed across a transition point. In a NEW-to-OLD transition the AGGREGATOR may be a mappable ASN, in which case the value is truncated to 16 bits and no further action is required. Otherwise the 32-bit AGGREGATOR value is rewritten into a NEW_AGGREGATOR attribute and the transition 16-bit value, AS 2356, is placed into the AGGREGATOR attribute. On an OLD-to-NEW transition the NEW_AGGREGATOR attribute is copied back into the AGGREGATOR attribute, if defined; otherwise the AGGREGATOR is padded out with leading zeros.

Transition

Transition in this scheme is relatively straightforward. NEW BGP speakers can be deployed within the network in a piecemeal fashion without any major concerns, and no changes are required for OLD BGP speakers. The size of BGP UPDATE messages is slightly longer because of the extended length of the AS PATH attribute in NEW BGP and the NEW_AS_PATH attribute that has been added in the OLD BGP environment, but it should not prove to be a major factor.

BGP loop prevention appears to be adequately addressed in all commonly encountered situations, and there appears to be no other significant transition considerations from the perspective of BGP platforms.

This scenario implies a relatively straightforward transition, in that OLD BGP speakers do not have to migrate to NEW BGP capability just because 32-bit ASNs are deployed elsewhere in the network. As long as they transmit the NEW-AS_PATH update across their domain without attempting to alter it in any way, then the 32-bit routing realm will be able to perform loop detection and shortest AS path selection in a manner that is entirely consistent with the 16-bit routing realm. Deployment of NEW BGP code is required only when the local AS is numbered from the nonmappable 32-bit ASN space.

Alternatives to AS Numbers

It is certainly a challenging task to contemplate an environment in which a 32-bit ASN space is exhausted, but one would suppose that the same consideration was in the minds of the original BGP protocol designers when they opted to use 16-bit ASNs. Of course a 32-bit number pool is not double the pool size of a 16-bit number pool—it is 65,536 times larger. That does appear to lead one to believe that this time it will be a far more challenging task to exhaust this expanded number pool.

This approach of simply extending the number space appears to offer a path of minimal disruption and minimal change in terms of operational configuration, storage, message size, and processing overheads for BGP. Nothing much has changed here except the range of the number space, and some ancillary considerations relating to transitional arrangements.

Of course, other labeling spaces remain possibilities, and a shift to a different labeling scheme could well use the same transitional approach. There is no significance in the ASN apart from its uniqueness, and any other form of name space would function equally well in terms of its role in BGP. One could use strings such as domain names, URIs, fixed-length hashes of public keys, the public keys themselves, or even IPv6 addresses as distinguishing AS identifiers.

There is no direct requirement for summarization of ASN ranges within the protocol use, no requirement within the protocol to continue to use number identifiers, and no direct requirement to stick with values that are encoded in a fixed-length field.

However, such approaches would add to the size of BGP UPDATE messages, increase the storage requirements, and, perhaps marginally, increase processing overheads for BGP. The more complex the identity space the more complex the basic task of BGP configuration and the higher the possibility of mistakes. “Borrowing” AS identifiers from another name space, such as domain names, or derived URIs, has the associated concern that the uniqueness of the space is derived from the inherent stability and uniqueness of the name space upon which the identifiers are derived. It is definitely possible that at times this trust is misplaced.

Numbers are often the simplest of identifiers. This approach represents minimal change to the installed base of BGP speakers, and there is no requirement for an existing routing domain using a 16-bit ASN and OLD BGP to make any changes to its routing environment at all. The transition appears to offer flexibility, orderly transition, and minimal disruptions to existing operational practices.

Conclusion

We are certainly running out of available 16-bit ASNs, and an industry of the size of the Internet is no longer as agile as it may have been in the past to make the necessary adjustments to alleviate this situation. At present we need to have a considerable period of advance warning of change in something as fundamental as the interdomain routing space in order to be able to integrate changes into various operational cycles of testing and transitional deployment prior to integration into production environments. The first steps that need to be taken are the completion of the technical specification of this approach in the form of an Internet standard and the production and distribution of BGP implementations that support 32-bit ASNs from the existing BGP implementation suppliers. It would be preferable to get this transition process under way in the near future, while there is still time to complete the transition well before we exhaust the current 16-bit ASN space.

For Further Reading

- [1] “BGP Support for Four-octet AS Number Space,” E. Chen, Q. Vohra, work in progress, (**draft-ietf-idr-as4bytes-12.txt**), November 2005.
The 32-bit AS description and the associated transition considerations. This work is expected to be completed shortly, and published as an RFC as a proposed standard document.
- [2] “The AS Number Report”, G. Huston, (updated on a daily basis) **<http://www.potaroo.net/tools/asns>**
A longer description of the numerical analysis used in the prediction of AS Number exhaustion.
- [3] “ASN Missing in Action”, H. Uijterwaal, R. Wilhelm, Document RIPE-353, (**<http://www.ripe.net/docs/ripe-353.html>**), October 2005.
Another analysis of AS Number consumption has been performed by Henk Uijterwaal and Rene Wilhelm, using the RIR AS number allocation rate as the base for the predictive exercise.
- [4] “A Border Gateway Protocol 4 (BGP 4),” Y. Rekhter, Ed. T. Li, Ed. S. Hares, Ed., RFC 4271, January 2006.
- [5] “Capabilities Advertisement with BGP-4,” R. Chandra, J. Scudder, RFC 3392, November 2002.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for almost two decades, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector, and has served time with Telstra, where he was the Chief Scientist in the company’s Internet area. Geoff is currently the Internet Research Scientist at the Asia Pacific Network Information Centre (APNIC). He has been a member of the Internet Architecture Board, and currently co-chairs two Working Groups in the IETF. He is author of a number of Internet-related books. E-mail: **gih@apnic.net**

Working with IP Addresses

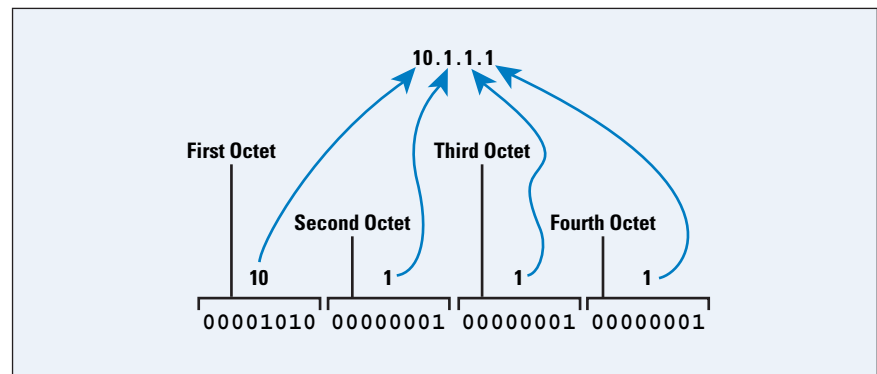
by Russ White, Cisco Systems

IP addresses, both IPv4 and IPv6, appear to be complicated when you first encounter them, but in reality they are simple constructions, and using a few basic rules will allow you to find the important information for any situation very quickly—and with minimal math. In this article, we review some of the basics of IPv4 address layout, and then consider a technique to make working with IPv4 addresses easier. Although this is not the “conventional” method you might have been taught to work with in IP address space, you will find it is very easy and fast. We conclude with a discussion of applying those techniques to the IPv6 address space.

Basic Addressing

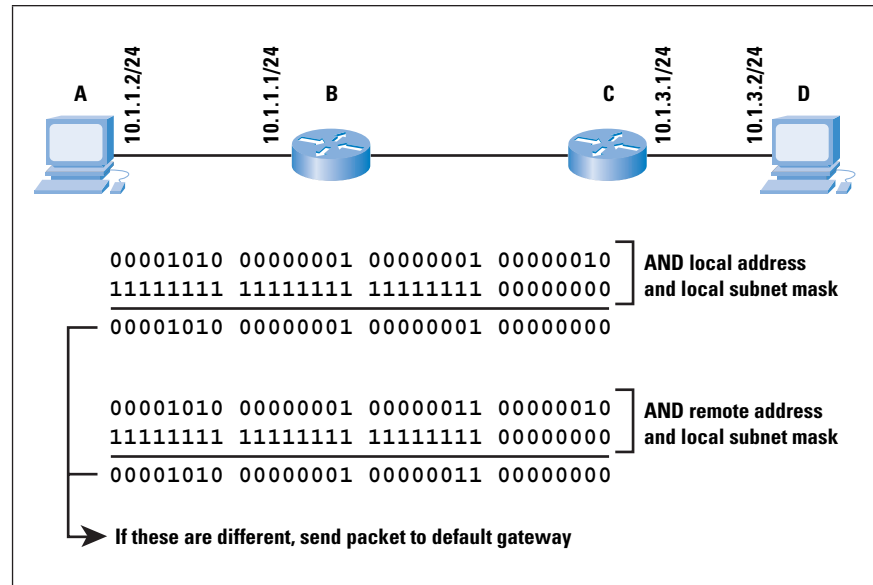
IPv4 addresses are essentially 32-bit binary numbers; computer systems and routers do not see any sorts of divisions within the IPv4 address space. To make IPv4 addresses more human-readable, however, we break them up into four sections divided by dots, or periods, commonly called “octets.” An octet is a set of eight binary digits, sometimes also called a “byte.” We do not use byte here, because the real definition of a byte can vary from computer to computer, whereas an octet remains the same length in all situations. Figure 1 illustrates the IPv4 address structure.

Figure 1: IPv4 Address Structure



Because each octet represents a binary (base 2) number between 0 and 2^8 , each octet will be between 0 and 255. This part of IPv4 addresses is simple—but what about subnet masks? To understand a subnet mask, we need to understand how a device actually uses subnet masks to determine where to send a specific packet, as Figure 2 illustrates.

Figure 2: Subnet Masks



If host A, which has the local IP address **10.1.1.2** with a subnet mask of **255.255.255.0**, wants to send a packet to **10.1.3.2**, how does it know whether D is connected to the same network (broadcast domain) or not? If D is connected to the same network, then A should look for D's local Layer 2 address to transmit the packet to. If D is not connected to the same network, then A needs to send any packets destined to D to A's local default gateway.

To discover whether D is connected or not, A takes its local address and performs a logical AND between this and the subnet mask. A then takes the destination (remote) address and performs the same logical AND (using its local subnet mask). If the two resulting numbers, called the *network address* or *prefix*, match, then the destination must be on the local segment, and A can simply look up the destination in the *Address Resolution Protocol* (ARP) cache, and send the packet locally. If the two numbers do not match, then A needs to send the packet to its default gateway.

Note: ARP is a protocol used to discover the mappings between the IP addresses of devices attached to the same network as the local device and the Layer 2 address of devices attached to the same network as the local device. Essentially, a device sends an ARP broadcast containing the IP address of some other device it believes to be connected, and the device with the specified IP address replies with its Layer 2 address, providing a mapping between these two addresses.

If a subnet mask is a "dotted decimal" version of the binary subnet mask, then what is the prefix length? The prefix length is just a shorthand way of expressing the subnet mask. The prefix length is the number of bits set in the subnet mask; for instance, if the subnet mask is **255.255.255.0**, there are 24 1's in the binary version of the subnet mask, so the prefix length is 24 bits. Figure 3 illustrates network masks and prefix lengths.

Figure 3: Prefix Lengths

Binary Mask	Prefix Length	Subnet Mask
11111111 00000000 00000000 00000000	/8	255.0.0.0
11111111 10000000 00000000 00000000	/9	255.128.0.0
11111111 11000000 00000000 00000000	/10	255.192.0.0
11111111 11100000 00000000 00000000	/11	255.224.0.0
11111111 11110000 00000000 00000000	/12	255.240.0.0
11111111 11111000 00000000 00000000	/13	255.248.0.0
11111111 11111100 00000000 00000000	/14	255.252.0.0
11111111 11111110 00000000 00000000	/15	255.254.0.0
11111111 11111111 00000000 00000000	/16	255.255.0.0
11111111 11111111 10000000 00000000	/17	255.255.128.0
11111111 11111111 11000000 00000000	/18	255.255.192.0
11111111 11111111 11100000 00000000	/19	255.255.224.0
11111111 11111111 11110000 00000000	/20	255.255.240.0
11111111 11111111 11111000 00000000	/21	255.255.248.0
11111111 11111111 11111100 00000000	/22	255.255.252.0
11111111 11111111 11111110 00000000	/23	255.255.254.0
11111111 11111111 11111111 00000000	/24	255.255.255.0
11111111 11111111 11111111 10000000	/25	255.255.255.128
11111111 11111111 11111111 11000000	/26	255.255.255.192
11111111 11111111 11111111 11100000	/27	255.255.255.224
11111111 11111111 11111111 11110000	/28	255.255.255.240
11111111 11111111 11111111 11111000	/29	255.255.255.248
11111111 11111111 11111111 11111100	/30	255.255.255.252
11111111 11111111 11111111 11111110	/31	255.255.255.254
11111111 11111111 11111111 11111111	/32	255.255.255.255

Working with IPv4 Addresses

Now that we understand how an IPv4 address is formed and what the subnet length and prefix length are, how do we work with them? The most basic questions we face when working with an IP address follow:

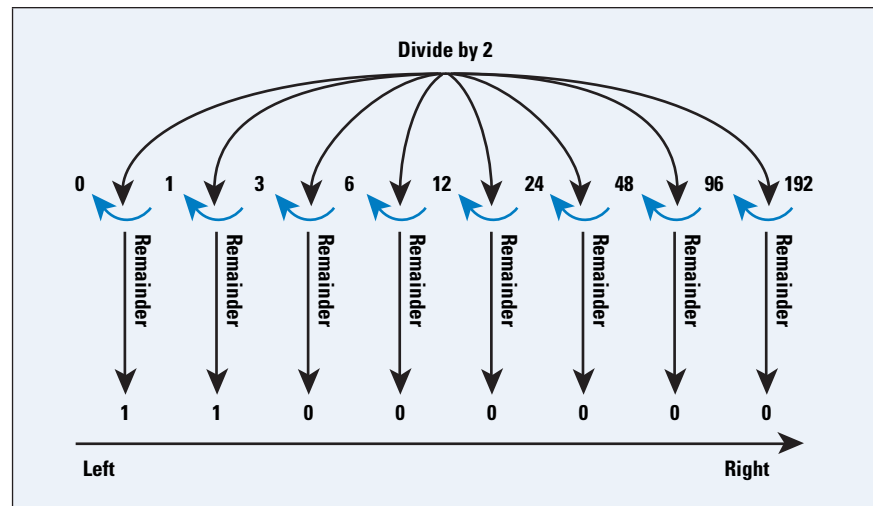
- What is the network address of the prefix?
- What is the host address?

There are two ways to find the answers to these questions: the hard way and the easy way. We cover the hard way first, and then show you the easy way.

The Hard Way

The hard way to determine the prefix and host addresses is to convert the address into binary, perform logical AND and NOR operations on the address and the subnet mask, and then convert the resulting numbers back to decimal. Figure 4 illustrates the process of converting a single octet of the IPv4 address into binary; the number converted in this case is 192.

Figure 4: Binary Conversion



The process is simple, but tedious; divide the octet value by 2, take the remainder off, and then divide by 2 again, until you reach 0. The remainders, reversed in direction, are the binary numbers representing the value of the octet. Performing this process for all four octets, we have the binary IP address, and can use logical AND and NOR operations to find the prefix (network address) and the host address, as Figure 5 shows for the address **192.168.100.80/26**.

Figure 5: Address Calculation

Network		11000000	10101000	01100100	01010000
		192	168	100	80
		11111111	11111111	11111111	11000000
		8	+8	+8	+2 == 26
AND		11000000	10101000	01100100	01000000
		192	168	100	64
Host		11000000	10101000	01100100	01010000
		192	168	100	80
		11111111	11111111	11111111	11000000
		8	+8	+8	+2 == 26
NOR		00000000	00000000	00000000	00010000
		0	0	0	16

The Easy Way

All this conversion from binary to decimal and from decimal to binary is tedious— is there an easier way? Yes. First, we start with the observation that we work only with the numbers within one octet at a time, no matter what the prefix length is. We can assume all the octets before this *working octet* are part of the network address, and octets after this *working octet* are part of the host address.

The first thing we need to do, then, is to find out which octet is our *working octet*. This task is actually quite simple: just divide the prefix length by 8, discard the remainder, and add 1. The following table provides some examples.

Address	Hard Math	Working Octet
192.158.100.80/26	$(26 \div 8) + 1 = 4$	4
10.1.1.48/23	$(23 \div 8) + 1 = 3$	3
172.31.80.10/22	$(22 \div 8) + 1 = 3$	3

*Note: Another way to look at this task is that you will ignore the octets indicated by the division. For instance, for **192.168.100.80/26**, the result of dividing 26 by 8 is 3, so you will ignore the first three octets of the IP address, and work only with the fourth octet. This process has the same result.*

When we know the working octet, what do we do with it? Well, we could simply use the procedure outlined, convert the single octet to binary, perform AND and NOR operations on it with the right bits from the subnet mask, and then put it all back together to find the network and host addresses—but there is an easier way to find the network and host parts of the working octet. Start by doing the same math, only this time we want to work with the remainder rather than the result.

192.168.100.80/26

$26 \div 8 = 3$ with a remainder of 2

Take the remainder, and use the following table to find the corresponding *jump* within the octet; this number is the distance, in decimal form, between the network addresses within the octet.

1	2	3	4	5	6	7	8
128	64	32	16	8	4	2	1

In this chart, the first line represents the prefix length *within this octet*, the second line represents the prefix value when this bit is set to 1, the number of hosts in the subnet for this prefix length, and the *jump* between network addresses with the specified prefix length.

The number 2 corresponds to 64, so the jump is 64—there is a network at 0, 64, 128, 192, and 224 in this octet. Now all we need to do is figure out which one of those networks this address is in. This task is fairly simple: just take the largest network number that fits into the number in the working octet. In this case, the largest number that fits into 80 is 64, so our network address is **192.168.100.64/26**.

Now, what about the host address? That is easy when we have the network address—just subtract the network address from the IP address, and you have the host address within the network: $80 - 64 = 16$. This process takes a little practice, but it is not hard when you become accustomed to the steps.

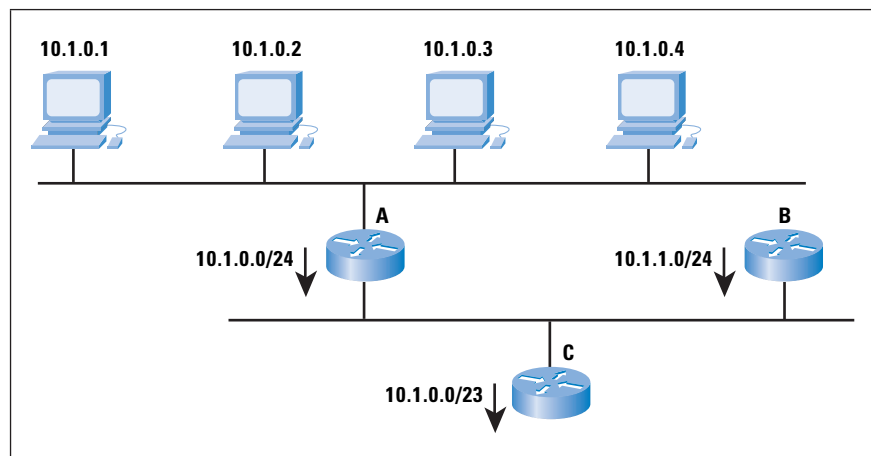
Address	Hard Math	Working Octet	Remainder	Jump	Network	Host
192.158.100.80/26	$26 \div 8 = 3$	$3 + 1 = 4$	2	64	192.168.100.64/26	$80 - 64 = 16$
10.1.1.48/23	$23 \div 8 = 2$	$2 + 1 = 3$	7	2	10.1.0.0/23	$1 - 0 = 1.48$
172.31.80.10/22	$22 \div 8 = 2$	$2 + 1 = 3$	6	4	172.31.80.0/22	$80 - 80 = 0.10$

In the second and third examples, you see that the working octet is actually the third, rather than the fourth, octet. To find the host address in these examples, you simply find the host address in the third octet, and then “tack on” the fourth octet as part of the host address as well, because part of the third octet—and all of the fourth octet—are actually part of the host address.

Summarization and Subnets

Subnets and supernets are probably the hardest part of IP addressing for most people to understand and handle quickly, but they are both based on a very simple concept—*aggregation*. Figure 6 shows how aggregation works.

Figure 6: Address Aggregation



The figure shows four hosts with the addresses **10.1.0.1**, **10.1.0.2**, **10.1.0.3**, and **10.1.0.4**. Router A advertises **10.1.1.0/24**, meaning: “Any host within the address range **10.1.0.0** through **10.1.0.255** is reachable through me.” Note that not all the hosts within this range exist, and that is okay—if a host within that range of addresses is reachable, it is reachable through Router A. In IP, the address that A is advertising is called a *network address*, and you can conveniently think of it as an address for the wire the hosts and router are attached to, rather than a specific device.

For many people, the confusing part comes next. Router B is also advertising **10.1.1.0/24**, which is another network address. Router C can combine—or aggregate—these two advertisements into a single advertisement. Although we have just removed the correspondence between the wire and the network address, we have not changed the fundamental meaning of the advertisement itself. In other words, Router C is saying: “Any host within the range of addresses from **10.1.0.0** through **10.1.1.255** is reachable through me.” There is no wire with this address space, but devices beyond Router C do not know this, so it does not matter.

To better handle aggregated address space, we define two new terms, *subnets* and *supernets*. A subnet is a network that is contained entirely within another network; a supernet is a network that entirely contains another network. For instance, **10.1.0.0/24** and **10.1.1.0/24** are both subnets of **10.1.0.0/23**, whereas **10.1.0.0/23** is a supernet of **10.1.0.0/24** and **10.1.1.0/24**.

Now we consider a binary representation of these three addresses, and try to make more sense out of the concept of aggregation from an addressing perspective; Figure 7 illustrates.

Figure 7: Aggregation Details

00001010	00000001	00000000	00000000	10.1.0.0/24
11111111	11111111	11111111	00000000	
00001010	00000001	00000001	00000000	10.1.1.0/24
11111111	11111111	11111111	00000000	
00001010	00000001	00000000	00000000	10.1.0.0/23
11111111	11111111	11111110	00000000	

Changing Bit

By looking at the binary form of **10.1.0.0/24** and **10.1.1.0/24**, we can see that only the 24th bit in the network address changes. If we change the prefix length to 23, we have effectively “masked out” this single bit, making the **10.1.0.0/23** address cover the same address range as the **10.1.0.0/24** and **10.1.1.0/24** addresses combined.

The Hardest Subnetting Problem

The hardest subnetting problem most people face is that of trying to decide what the smallest subnet is that will provide a given number of hosts on a specific segment, and yet not waste any address space. The way this sort of problem is normally phrased is something like the following:

You have 5 subnets with the following numbers of hosts on them: 58, 14, 29, 49, and 3, and you are given the address space **10.1.1.0/24**. Determine how you could divide the address space given into subnets so these hosts fit into it.

This appears to be a very difficult problem to solve, but the chart we used previously to find the jump within a single octet actually makes this task quite easy. First, we run through the steps, and then we solve the example problem to see how it actually works.

- Order the networks from the largest to the smallest.
- Find the smallest number in the chart that fits the number of the largest number of hosts + 2 (*you cannot, except on point-to-point links, use the address with all 0's or all 1's in the host address; for point-to-point links, you can use a /31, which has no broadcast addresses*).
- Continue through each space needed until you either run out of space or you finish.

This process seems pretty simple, but does it work? Let's try it with our example.

- Reorder the numbers 58, 14, 29, 49, 3 to 58, 49, 29, 14, 3.
- Start with 58.
 - The smallest number larger than (58 + 2) is 64, and 64 is 2 bits.
 - There are 24 bits of prefix length in the address space given; add 2 for 26.
 - The first network is **10.1.1.0/26**.
 - The next network is **10.1.1.0 + 64**, so we start the next “round” at **10.1.1.64**.
- The next block is 49 hosts.
 - The smallest number larger than (49 + 2) is 64, and 64 is 2 bits.
 - There are 24 bits of prefix length in the address space given; add 2 for 26.
 - We start this block at **10.1.1.64**, so the network is **10.1.1.64/26**.
 - The next network is **10.1.1.64 + 64**, so we start the next “round” at **10.1.1.128**.
- The next block is 29 hosts.
 - The smallest number larger than (29 + 2) is 32, and 32 is 3 bits.
 - There are 24 bits of prefix length in the address space given; add 3 for 27.
 - We start this block at **10.1.1.128**, so the network is **10.1.1.128/27**.

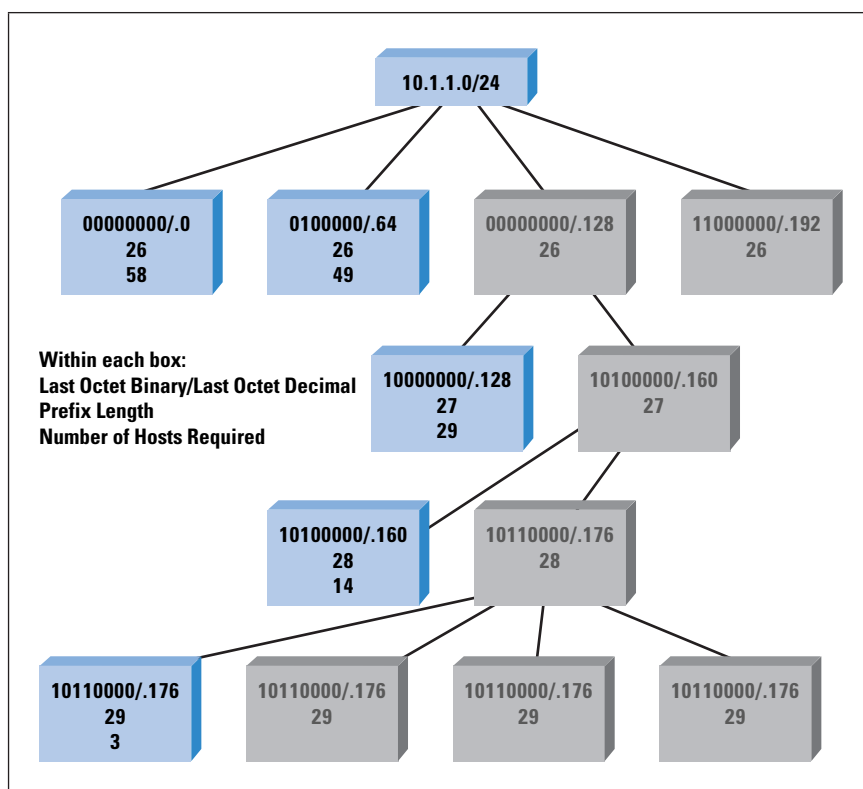
- The next network is **10.1.1.128 + 32**, so we start the next “round” at **10.1.1.160**.
- The next block is 14 hosts.
 - The smallest number larger than $(14 + 2)$ is 16, and 16 is 4 bits (actually equal, but it still works).
 - There are 24 bits of prefix length in the address space given; add 4 for 28.
 - We start this block at **10.1.1.160**, so the network is **10.1.1.160/28**.
 - The next network is **10.1.1.160 + 16**, so we start the next “round” at **10.1.1.176**.

The last block is 3 hosts.

- The smallest number larger than $(3 + 2)$ is 8, and 8 is 5 bits.
- There are 24 bits of prefix length in the address space given; add 5 for 29.
- We start this block at **10.1.1.176**, so the network is **10.1.1.176/29**.
- This is the last block of hosts, so we are finished.

It is a simple matter of iterating from the largest to the smallest block, and using the simple chart we used before to determine how large of a *jump* we need to cover the host addresses we need to fit onto the subnet. Figure 8 illustrates the resulting hierarchy of subnets.

Figure 8: Subnet Chart



In this illustration:

- The first line in each box contains the final octet of the network address in binary and decimal forms.
- The second line in each box contains the prefix length.
- The third line indicates the number of hosts the original problem required on that subnet.
- Gray boxes indicate blocks of address space that are unused at that level.

Working with IPv6 Addresses

IPv6 addresses appear to be much more difficult to work with—but they really are not. Although they are larger, they are still made up of the same fundamental components, and hosts and routers still use the addresses the same way. All we really need to do is realize that each pair of hexadecimal numbers in the IPv6 address is actually an octet of binary address space. The chart, the mechanisms used to find the network and host addresses, and the concepts of super and subnets remain the same.

For example, suppose we have the IPv6 address **2002:FF10:9876:DD0A:9090:4896:AC56:0E01/63** and we want to know what the network number is (host numbers are less useful in IPv6 networks, because they are often the MAC address of the system itself).

- $63 \div 8 = 7$, remainder 7.
- The working octet is the 8th, which is 0A.
- Remainder 7 on the chart says the jump is 2, so the networks are **00, 02, 04, 06, 08, 0A, 0C, and 0E**.
- The network is **2002:FF10:9876:DD0A::/63**.

The numbers are longer, but the principle is the same, as long as you remember that every *pair* of digits in the IPv6 address is a single octet.

Summary

IP addresses appear to be very complex on first approach, but their inbuilt structure actually provides easy ways to divide the problems into pieces and approach one piece of the problem at a time—the same way we design and build networks on a large scale. If you learn to use some simple techniques and understand how IP addresses are structured, they are relatively easy to work with.

For Further Reading

The following IETF *Requests for Comments* (RFCs) provide information on IP addressed and addressing structures:

- [1] V. Fuller, T. Li, J. Yu, K. Varadhan, “Supernetting: an Address Assignment and Aggregation Strategy,” RFC 1338, June 1992.
- [2] E. Gerich, “Guidelines for Management of IP Address Space,” RFC 1466, May 1993.
- [3] Y. Rekhter, T. Li, “An Architecture for IP Address Allocation with CIDR,” RFC 1518, September 1993.
- [4] V. Fuller, T. Li, J. Yu, K. Varadhan, “Classless Inter-Domain Routing (CIDR): an Address Assignment and Aggregation Strategy,” RFC 1519, September 1993.
- [5] Y. Rekhter, B. Moskowitz, D. Karrenberg, G. J. de Groot, E. Lear, “Address Allocation for Private Internets,” RFC 1918, February 1996.

RUSS WHITE works for Cisco Systems in the Routing Protocols Deployment and Architecture (DNA) team in Research Triangle Park, North Carolina. He has worked in the Cisco Technical Assistance Center (TAC) and Escalation Team in the past, has co-authored several books on routing protocols, including *Advanced IP Network Design*, *IS-IS for IP Networks*, and co-author of *Practical BGP*. He is the co-chair of the Routing Protocols Security Working Group within the IETF. E-mail: riw@cisco.com

Letter to the Editor

Dear Editor,

I read with interest the article by Dave Crocker in the December 2005 issue of IPJ (Volume 8, No. 4) titled “Challenges in Anti-Spam Efforts.” However, I was surprised not to find any mention of *graylisting*, an effective anti-spam technique. The technique is not a programmatic or analytical approach to the spam problem but rather relies on exploiting the general behavioral weakness of spam delivery (that spammers typically want to try an address just once for their delivery of spam contents). The technique provides a pragmatic solution to the contemporary bulk commercial e-mail problem to a large extent.

If you are planning for a sequel of this article, I would strongly advocate mentioning the technique for the general benefit of the community.

I administrate a national ISP of considerable size in Pakistan, and the extent to which graylisting has helped us in fighting against spam is amazing.

Successful spam-fighting techniques (especially those that are still far from being widely adopted and worked upon) of today make good candidates for future efforts. My enthusiasm for graylisting is chiefly a result of the benefits our company has derived from its use, but I also want to champion its use because I think it is not widely adopted among peer ISPs because of ignorance. Hence my enthusiastic advocacy of this unsung hero in the fight against spam.

Citations:

Graylisting, <http://en.wikipedia.org/wiki/Graylisting>

—Tee Emm, Supernet, Pakistan
tm@super.net.pk

The Author responds:

Dear Editor,

I appreciate Tee Emm’s concern that graylisting was not explicitly cited in my article.

I must use the cliché of “limited space” as my primary excuse for omitting graylisting. The tight constraints for a brief article required some difficult choices. As I mentioned at the end of the article, the people reviewing it before publication were particularly helpful (and vigorous). The question of what detail to include was a major focus. My decision was to have only a basic review of existing techniques, because the focus of the article was on future activities.

I believe the work on detection and reaction mechanisms against “bad actors” is reasonably mature. I also believe that the creation of a trust overlay for e-mail, to permit better handling of messages from “good actors,” is very early and in need of much more focus.

With that said, I think I can also generate a plausible claim that graylisting is a form of “traffic shaping,” which is cited in the article.

I primarily meant the traffic shaping reference to be about the technique of tracking aggregate (statistical) flow from a specific address. However, I think that graylisting constitutes a simple—albeit quite useful—mechanism that is designed to slow down the senders, to limit their impact. As Tee Emm notes, graylisting is based on the spammers’ pattern of giving up, after a single failure to send the message. That is the ultimate “shaping,” I think.

Certainly a summary of existing techniques is a worthy topic. It has become quite a rich topic, and matured to a level of qualifying as an area of administration and operations specialization.

As for a follow-up article, I do not have one planned, currently. If I do another one, I hope it will be about open mechanisms for achieving authentication and assessment (vetting) of good actors.

Perhaps next year.

For reference, I should note that there has been some public follow-up on the article, when CircleID reprinted a posting I made about it: http://www.circleid.com/posts/challenges_in_anti_spam_efforts/

Again, I appreciate Tee Emm’s interest and comment.

—Dave Crocker, *Brandenburg Internet Working*
dcrocker@bbiw.net

IETF @ 20

The *Internet Engineering Task Force* (IETF) and the *Internet Society* (ISOC) celebrate the 20th anniversary of the IETF, the world's leading Internet standards development body. The IETF is a large open international community of network designers, operators, vendors, and researchers concerned with the evolution of the Internet architecture and the smooth operation of the Internet. Its principal task today is the development and publication of technical specifications and standards for basic Internet protocols. It is open to any interested individual.

The first IETF meeting was held on the afternoon of January 16, 1986, in San Diego, California. As a community-driven activity the IETF went on to pioneer a unique, open process for standards development. Open to all, and based on principles such as “rough consensus and running code,” the IETF has enabled the development of standards that have supported every aspect of the Internet's phenomenal growth.

“The IETF is unique,” said Brian Carpenter, IETF Chair. “Unlike other standards bodies, there is very little in the way of formal hierarchy and there are no membership requirements or fees. The IETF welcomes broad participation by anyone interested in the future technical evolution and stability of the Internet—and IETF standards are available to all, without charge.”

“There is global recognition of the achievements of the IETF in its support of the development of Internet technology. As the demands on the Internet increase, the IETF clearly has a vital role to play in ensuring that Internet technologies continue to evolve in a coherent and coordinated manner,” said Leslie Daigle, chair of the *Internet Architecture Board* (IAB) which provides architectural oversight of IETF activities.”

“The success of the IETF has largely been due to a pragmatic, consensus-based approach to technology standards development,” noted Lynn St. Amour, President and CEO of ISOC. “Many of the principles of co-operation and collaboration that were developed in the IETF are now being successfully applied in other global forums. ISOC is proud to be associated with the IETF—we value its members' accomplishments over the last 20 years and look forward to celebrating these achievements over the course of 2006.”

ISOC has declared 2006 “The Year of the IETF” and will be running several activities during the year in celebration of the IETF's 20th anniversary. For more information, see: <http://ietf20.isoc.org>

The Internet Society is a not-for-profit membership organization founded in 1992 to provide leadership in Internet related standards, education, and policy. With offices in Washington, DC, and Geneva, Switzerland, it is dedicated to ensuring the open development, evolution and use of the Internet for the benefit of people throughout the world. ISOC is the organizational home of the IETF and other Internet-related bodies who together play a critical role in ensuring that the Internet develops in a stable and open manner. For over 13 years ISOC has run international network training programs for developing countries and these have played a vital role in setting up the Internet connections and networks in virtually every country connecting to the Internet during this time.

ISOC Welcomes WSIS Proposal

Delegates meeting at the *World Summit on the Information Society* (WSIS) in Tunis have affirmed their commitment to build on the governance mechanisms that have enabled the Internet's incredibly successful growth.

ISOC welcomes the recognition by WSIS of how the effectiveness of the existing arrangements for Internet governance has helped make the Internet the highly robust, dynamic and geographically diverse medium that it is today.

"We are delighted that there is now much broader recognition of the achievements of the organisations that support the Internet community," said Lynn St. Amour, President and CEO of the ISOC. "These organizations, along with their open, consensus-based processes clearly have a vital role to play in the further development of the Internet. It is also significant that the WSIS debate has moved beyond the details of technical administration and on to broader issues that require increased coordination by stakeholders in order to ensure the continued stability of the Internet."

The WSIS recommendation includes a proposal for a new forum for multi-stakeholder policy dialogue—the *Internet Governance Forum*. ISOC, together with partner organizations from the Internet community, has always worked to encourage full engagement in such dialogues by all those with an interest in the Internet's future. ISOC believes that the forum's success depends upon the fullest participation by all stakeholders. At the same time, ISOC is pleased to note that the proposed forum would have no oversight function and would have no involvement in the day-to-day operations of the Internet.

"ISOC will facilitate increased cooperation and information sharing amongst all parties interested in Internet governance and we look forward to playing an active role in the new forum as is expected of us by the global community," said Lynn St. Amour. "We very much hope that the Tunis summit will lead to some real and positive outcomes that will help bring the benefits of the Internet to people everywhere—especially to those who are yet to be connected."

ISOC, along with some of its partner organisations—the *Number Resource Organisation* (NRO), the IETF, *London Internet Exchange* (LINX), the *Internet Corporation for Assigned Names and Numbers* (ICANN) and the *Council of European National Top level domain Registries* (CENTR)—were present at the *ICT 4 All* exhibition held in conjunction with WSIS.

For more information about the organizations listed above visit:

<http://isoc.org>

<http://ietf.org>

<http://iab.org>

<http://www.intgovforum.org>

<http://www.linx.net>

<http://nro.org>

<http://www.centri.org>

<http://www.itu.int/wsis>

Upcoming Events

The *Internet Engineering Task Force* (IETF) will meet in Montreal, Canada, July 9–14, 2006. For more information, visit:

<http://ietf.org>

ACM's *SIGCOMM 2006* will be held in Pisa, Italy, September 11–15, 2006. For more information, visit:

<http://www.acm.org/sigs/sigcomm/sigcomm2006>

The *North American Network Operators Group* (NANOG) will meet in St. Louis, MO October 8–10, 2006. For more information, see: **<http://nanog.org>**

The *American Registry for Internet Numbers* (ARIN) will meet (jointly with NANOG) in St. Louis, October 11–13, 2006. For more information, see: **<http://arin.net>**

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, Professor, WIDE Project
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 2006 Cisco Systems Inc.
All rights reserved.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA
--

The Internet Protocol Journal

June 2006

Volume 9, Number 2

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
Gigabit TCP.....	2
Instant Messaging.....	27
Letters to the Editor.....	38
Corrections	43
Book Review.....	44
Fragments	47

In our June 2000 issue we wrote: “Two protocols used in the Internet are so important that they deserve special attention: the *Internet Protocol* (IP) from which this journal takes its name, and the *Transmission Control Protocol* (TCP). IP is fundamental to Internet addressing and routing, while TCP provides a reliable transport service that is used by most Internet applications, including interactive Telnet, file transfer, electronic mail, and Web page access via HTTP. Because of the critical importance of TCP to the operation of the Internet, it has received much attention in the research community over the years. As a result, numerous improvements to implementations of TCP have been developed and deployed.” We return to TCP in this issue with a look at its performance at gigabit speeds. Geoff Huston describes numerous research proposals related to TCP and discusses lessons learned by operators and researchers involved with this protocol.

My first encounter with the Internet (then called the ARPANET) took place in 1976 when I visited the *Norwegian Defence Research Establishment* (NDRE) at Kjeller, about 20 kilometers from Oslo, Norway. At NDRE, one of the researchers, named Pål, showed me a teletype terminal that was connected through the ARPANET to a host computer at SRI International in Menlo Park, California. After a few minutes, the teletype started printing messages from someone called “Geoff” on the other end of the line. Pål typed back, passing on questions from myself about the weather in California and so on. I later learned that the host computer was a PDP-10 model KA10 running the TENEX operating system. TENEX could “link” two terminals together so that anything typed on one terminal would appear on the other, and conversely. This primitive “chat” system is the forerunner of today’s *Instant Messaging* (IM) environment. David Strom gives an overview of the current state of IM solutions in our second article.

The article “Working with IP Addresses” in our last issue sparked several comments, some of which are included in our Letters to the Editor section. A few readers also noticed some errors in the article, so we have included the corrections in this issue. We very much appreciate your feedback. Please send your comments to: ipj@cisco.com

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Gigabit TCP

by Geoff Huston, APNIC

In looking back over some 30 years of experience with the Internet, the critical component of the Internet Protocol Suite that has been the foundation of its success as the technology of choice for the global communications system is the *Internet Protocol* (IP) itself, working an overlay protocol that can span almost any form of communications media. But I would also like to nominate another contender for a critical role within IP, namely the reliable transport protocol that sits on top of IP, the *Transmission Control Protocol* (TCP), and its evolution over time. In support of this nomination is the fact that the end-to-end rate-adaptive control algorithm that was adopted by TCP represented a truly radical shift from the reliable gateway-to-gateway virtual circuit flow control systems used by other protocols of similar vintage. It is also interesting to note that TCP is not designed to operate at any particular speed, but it attempts to operate at a speed that uses its fair share of all available network capacity along the network path. The fundamental property of the TCP flow control algorithm is that it attempts to be maximally efficient while also attempting to be maximally fair.

Previous articles on this topic, “TCP Performance”^[12] and “The Future for TCP”^[13] looked at the design assumptions behind TCP and its performance characteristics. The essential characteristic of TCP is that it attempts to establish a dynamic equilibrium with other concurrent sessions and opportunistically use all available network capacity. It achieves this by constantly altering its flow characteristics, continually probing the network to see if higher speeds are supportable, while also being prepared to immediately decrease the current sending rate in the face of received signals of network congestion.

In a world where network infrastructure capacity and complexity are related to network cost and delivered data is related to network revenue, TCP fits in well. The minimal assumptions that TCP makes about the capability of network components permit networks to be constructed using simple transmission capabilities and simple switching systems. “Simple” often is synonymous with cheap and scalable, and there is no exception here. TCP also attempts to maximize data delivery through adaptive end-to-end flow rate control and careful management of retransmission events. In other words, TCP is an enabler for cheaper networking for both the provider and consumer. For the consumer the offer of fast cheap communications has been a big motivation in the increase in demand for Internet-based services, and this—more than any other factor—has been the major enabling factor for the increased use of the Internet itself. “Cheap” is often enough in this world, and TCP certainly helps to make data communications efficient and therefore cheap.

Although TCP is highly effective in many networking environments, that does not mean it is highly effective in every environment. For example:

- In those wireless environments where there is significant wireless noise, TCP may confuse the outcome of radio-based signal corruption and the corresponding packet drop with the outcome of network congestion, and consequently the TCP session may back off its sending rate too early and back off for too long.
- TCP also backs off too early when the network routers have insufficient buffer space. This effect is more subtle, but it is related to the coarseness of the TCP algorithm and the consequent burstiness of TCP packet sequences. These bursts, which occur at up to twice the bottleneck capacity rate, are smoothed out by network buffers. Buffer exhaustion in the interior of the network causes packet drop, which causes the generation of a loss signal to the active TCP session, which, in turn, either halves its sending rate or—in the worst case—resets the session state and restarts with a single packet exchange. Particularly in wide-area networks, where the end-to-end delay-bandwidth product becomes a significant factor, TCP uses the network buffers to sustain a steady-state throughput that matches the available network capacity. Where the interior buffers are under-configured in memory it is not possible to even out the TCP bursts to continuously flow through the constrained point at the available data rate.
- TCP also asks its end hosts to have local capacity equal to the available network capacity on the forward and reverse paths. The reason is that TCP does not discard data until the remote end has reliably acknowledged it, so the sending host has to retain a copy of the data for the time it takes to send the data plus the time for the remote end to send the matching acknowledgement.

Even accounting for these limitations, it is true to say that TCP works amazingly well in most environments. Nevertheless, one area is proving to be quite a fundamental challenge to TCP as we know it, and that is the domain of wide-area, very-high-speed data transfer.

Very-High-Speed TCP

End host computers, even laptop computers these days, are typically equipped with Gigabit Ethernet interfaces, and have gigabytes of memory and internal data channels that can move gigabits of data per second between memory and the network interface. Current IP networks are constructed using multigigabit circuits and high-capacity switches and routers (assuming there is still a quantitative difference between these two forms of packet switching equipment). If the end hosts and the network both can support gigabit transmissions then a TCP session should be able to operate end to end at gigabits per second, and achieve the same efficiency at gigabit speeds as it does today at megabit speeds—right?

Well, no, not exactly!

This conclusion is not obvious, particularly when the TCP Land Speed Record is now at some 7Gbps across a distance that spans 30,000 km of network. What is going on?

Let's return to the basics of TCP to understand some of the variables with very-high-speed TCP. TCP operates in one of two states, that of *slow start* and *congestion avoidance*.

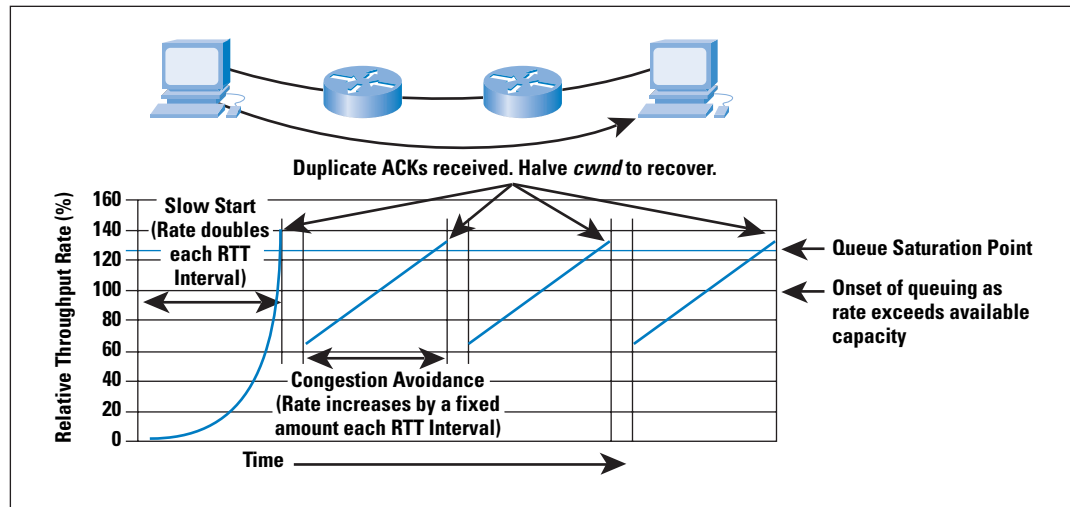
- *Slow start* mode is the initial mode of operation of TCP in any session, as well as its “reset” mode. In this mode, TCP sends two packets in response to each ACK packet that advances the sender's window. In approximate terms (*delayed* ACKs notwithstanding), this mode allows TCP to double its sending rate in each successive lossless *round-trip-time* (RTT) interval. The rate increase is exponential, effectively doubling each RTT interval, and the rate increase is bursty, effectively sending data into the network at twice the bottleneck capacity during this phase.

Sending data into the network at twice the bottleneck data speed is possible because of the “ACK clocking” property of TCP. Disregarding the complications of the TCP delayed ACK mechanism for a second, a TCP receiver generates a new ACK packet each time a packet arrives at the receiver. The sending rate of the ACKs is, in effect, the same as the receiving rate for the data packets. Assuming a one-way data transfer, so that ACK packets in the reverse direction are of minimal size, and assuming minimal jitter on the reverse path from the receiver back to the sender, the arrival rate of ACKs at the sender is comparable to the arrival rate of data packets at the receiver. In other words, the return ACK rate is comparable to the bottleneck capacity of the forward network path from sender to receiver. Sending two packets per received ACK is effectively sending packets into the network at twice the bottleneck capacity. At the bottleneck point the switching unit receives twice the amount of data than it can transmit to the output device over a period that corresponds to the delay-bandwidth product of the bottleneck link. Hence the comment that TCP is a *bursty* protocol, particularly at startup. For this reason TCP tends to operate more effectively across network switching elements that are generously endowed with memory, or have for each output port a buffer capacity roughly equal to the delay-bandwidth product of the link that is attached to that port.

- In the other operating mode, that of *congestion avoidance*, TCP sends an additional segment of data for each loss-free round-trip time interval. This increase is additive rather than exponential, increasing the sender's speed at the constant rate of one segment per RTT interval.

TCP undertakes a state transition upon the detection of packet loss. Small-scale packet loss (of the order of 1 or 2 packets per loss event) causes TCP to halve its sending rate and enter congestion avoidance mode, irrespective of whether it was in this mode already. Repetition of this cycle gives the classic sawtooth pattern of TCP behavior, and the related derivation of TCP performance as a function of packet loss rate. Longer sustained packet loss events cause TCP to stop using the current session parameters, recommence the congestion control session using the restart window size, and enter the slow start control mode once again. (See Figure 1).

Figure 1: TCP Behavior



But what happens when two systems are at opposite sides of a continent with a high-speed path between them? How long does it take for a single TCP session to get up to a data transfer rate of 10 Gbps? Can a single session operate at a sustained rate of 10 Gbps?

Let's look at a situation such as the network path from Brisbane, on the eastern side of the Australian continent, to Perth on the western side. The cable path is essentially along the southern coast of the continent, so the RTT delay is 70 ms, implying that there are 14.3 round-trip intervals per second. Let's also assume that the packet size being used is 1500 octets, or 12,000 bits, and the TCP initial window size is a single packet. And let's also assume that the bottleneck capacity of the host-to-host path between Brisbane and Perth is 10 Gbps.

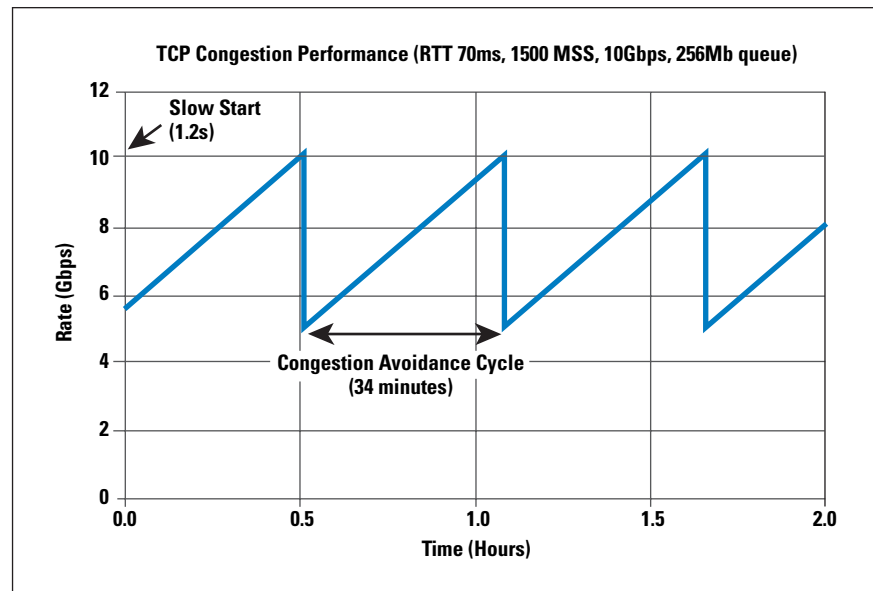
In a simple slow start model the sending speed doubles every 70 ms, so after 17 RTT intervals where the sending rate has doubled for each interval, or after some 1.2 seconds have elapsed, the transfer speed reaches 11.2 Gbps (assuming a theoretical host with sufficiently fast hardware components, sufficiently fast internal data paths, and adequate memory). At this stage let's assume that the sending rate exceeded the buffer capacity at the bottleneck point in the network path. Packet drop will occur, because the critical point buffers in the network path are now saturated.

At the point of reception of an ACK sequence that signals packet loss, the TCP sender's congestion window will halve, as will the TCP sending rate, and TCP will switch to congestion avoidance mode. In congestion avoidance mode the rate increase is 1 segment per RTT, equivalent to sending an additional 12 kilobits per RTT, or, given the session parameters as specified previously, equivalent to a rate increase of 171 kbps each RTT. So how long will it take TCP to recover and get back to a sending rate of 10 Gbps?

If this were a T1 circuit where the available path bandwidth is 1.544 Mbps, and congestion loss occurred at a sending rate of 2 Mbps (higher than the bottleneck transmission capacity due to the effect of queuing buffers within the network), then TCP would rate halve to 1 Mbps and then use congestion avoidance to increase the sending rate back to 2 Mbps. Within the selected parameters of a 70-ms RTT and 1500-byte segment size, this process involves using congestion avoidance to inflate the congestion window from 6 segments to 12. This process takes 0.42 seconds. So as long as the network can operate without packet loss for the session over an order of 1-second intervals, then TCP can comfortably operate at maximal speed in a megabit-per-second network.

What about our 10-Gbps connection? The first estimate is the amount of usable buffer space in the switching elements. Assuming a total of 256 MB of usable queue space on the network path prior to the onset of queue saturation, the TCP session operating in congestion avoidance mode will experience packet loss some 590 RTT intervals after reaching the peak transmission speed of 10 Gbps, or some further 41 seconds, at which point the TCP sending rate in congestion avoidance mode is 10.1 Gbps. For all practical purposes the TCP congestion avoidance mode causes the sawtooth oscillation of this ideal TCP session between 5.0 Gbps and 10.1 Gbps. A single iteration of this sawtooth cycle takes 2062 seconds, or 34 minutes and 22 seconds. The implication here is that the network has to be stable in terms of no packet loss along the path for time scales of the order of tens of minutes (or some billions of packets), and corresponding transmission bit error rates that are less than 10^{-14} . It also implies massive data sets to be transferred, because the amount of data passed in just one TCP congestion avoidance cycle is 1.95 terabytes (1.95×10^{12} bytes). It is also the case that the TCP session cannot make full use of the available network bandwidth, because the average data transfer rate is 7.55 Gbps under these conditions, not 10 Gbps. (See Figure 2).

Figure 2: TCP Behavior at High Speed



Clearly something is unexpected with this scenario, because it certainly looks like it is a difficult and lengthy task to fill a long-haul, high-capacity cable with data, and TCP is not behaving as expected. Although experimenting with the boundaries of TCP is in itself an interesting area of research, some practical problems here could well benefit from this type of high-speed transport.

A commonly quoted example, and certainly one of the more impressive ones is the Large Hadron Collider at CERN:

“The CERN Particle Physics lab in Geneva, Switzerland, successfully transmitted a data stream averaging 600Mbytes per second for 10 days to seven countries in Europe and the US. It was a crucial test of the computing infrastructure for the Large Hadron Collider being built at CERN. The LHC will be the most data intensive physics instrument ever built, generating 1500 Megabytes every second for a decade or more.”

—*New Scientist*, 30 April 2005

TCP and the Land Speed Record

The TCP Land Speed Record was originally an informal effort to achieve record-breaking TCP transfer speeds across IP networks. The late 1980s and early 1990s saw some noted milestones, particularly with Van Jacobson’s efforts in achieving sustained 10-Mbps and 45-Mbps TCP transfer speeds.

This activity has been incorporated into the Internet2 program, with the introduction of some formal rules about what constitutes a TCP Land Speed effort. In particular, the rules now have times, distances, and TCP constraints, and they call for the use of operational networks. Updates to the record have been posted frequently in recent years, and as of May 2006 the IPv4 single stream record is a TCP session operating at 7.21 Gbps for 30 minutes over 30,000 km of fibre path.

It is certainly possible to have TCP perform for sustained intervals at very high speed, as the land speed records for TCP show, but something else is happening here, and a set of preconditions need to be met before attempting to set a new record:

- First, it is good—indeed essential—to have the network path all to yourself. Any form of packet drop is a major problem here, so the best way to ensure no packets are lost is to keep the network path all to yourself.
- Secondly, it is good—indeed essential—to have a fixed latency. If the objective of the exercise is to reach a steady-state data transmission, then any change in latency, particularly a reduction in latency, has the risk of a period of oversending, which in turn has a risk of packet loss. So keep the network as stable as possible.
- Thirdly, it is good—indeed essential—to have extremely low bit error rates from the underlying transmission media. Data corruption causes checksum failure, which causes packet drop.
- Lastly, it is essential to know in advance both the round-trip latency and the available bandwidth.

You can then multiply these two numbers together (RTT and bandwidth), divide by the packet size, round down, and be sure to configure the sending TCP session to have precisely this buffer size, and the receiver to have a slightly larger size. And then start up the session.

The intention here is for TCP to use slow start to the point where the sender runs out of buffer space, at which point it will continue to sit at this buffer speed for as long as the sender, receiver and network path all remain in a stable state. For the example configuration of a 10-Gbps system with 70 ms RTT, setting a buffer limit of 116,000 packets will cause the TCP session to operate at 9.94 Gbps. As long as the latency remains steady (no jitter), with no bit errors, and as long as there is no other cross traffic, in theory this sending rate can be sustained indefinitely, with a steady stream of data packets being matched by a steady stream of ACK packets.

Of course, this situation is artificially constrained. The real concerns here with the protocol are in the manner in which it shares a network path with other concurrent sessions as well as its ability to fill the available network capacity. In other words, what would be good to see is a high-speed, high-volume version of TCP that could coexist on a network with all other forms of traffic, and, perhaps more ambitiously, that this high-speed form of TCP could share the network fairly with other traffic sessions while at the same time making maximal use of the network. The problem with TCP in its current incarnation is that it takes way too long in its additive increase mode (congestion avoidance) to recover its sustainable operating speed when operating at high speed across transcontinental-size network paths. If we want very-high-speed TCP to be effective and efficient, then we need to look at changes to TCP for high-speed operation.

High-Speed TCP

There are two basic approaches to high-speed TCP: parallelism of existing TCP, or changes to TCP to allow faster acceleration rates in a single TCP stream.

Using parallel TCP streams as a means of increasing TCP performance is an approach that has existed for some time. The original HTTP specification, for example, allowed the use of parallel TCP sessions to download each component of a Webpage (although HTTP 1.1 reverted to a sequential download model because the overheads of session startup appeared to exceed the benefits of parallel TCP sessions in this case). Another approach to high-speed file transfer through parallelism is that of GRID FTP. The basic approach is to split up the communications payload into numerous discrete components, and send each of these components simultaneously. Each component of the transfer can be between the same two endpoints (such as GRID FTP), or can be spread across multiple endpoints (as with BitTorrent).

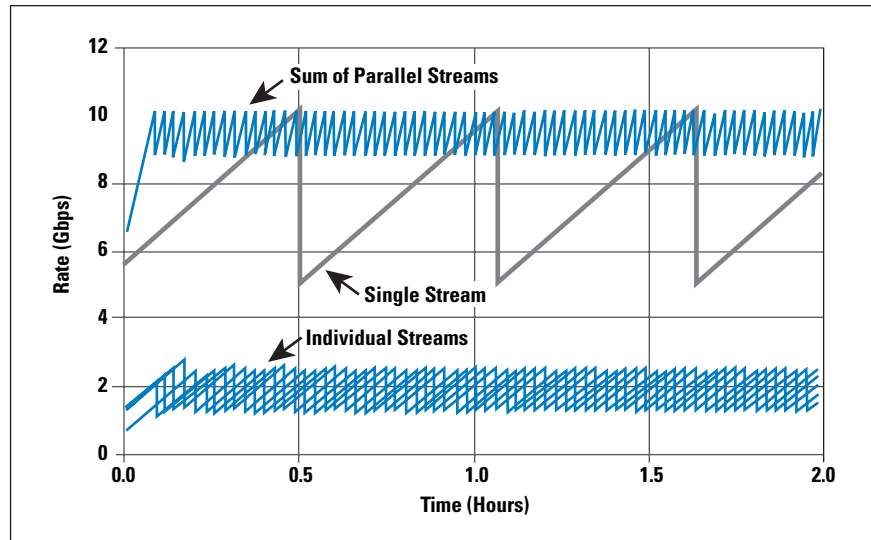
But for parallel TCP to operate correctly, we need to have already assembled all the data (or at a minimum know where all the data components are located). Where the data is being generated in real time (such as observatories or particle colliders) in massive quantities, there may be no choice but to treat the data set as a serial stream and use a high-speed transport protocol to dispatch it. In this case the task is to adjust the basic control algorithms for TCP to allow it to operate at high speed, but also to operate “fairly” on a mixed-traffic high-speed network.

Parallel TCP

Using parallelism as a key to higher speed is a common computing technique, and lies behind many supercomputer architectures. The same can apply to data transfer, where a data set is divided into numerous smaller chunks, and each component chunk is transmitted using its own TCP session. The underlying expectation here is that when using some number, N , of parallel TCP sessions, a single packet drop event will most probably cause the fastest of the N sessions to rate halve, because the fastest session will have more packets in flight in the network, and is therefore the most likely session to be impacted by a packet drop event. This session will then use congestion avoidance rate increase to recover, implying that the response to a single packet drop is reduction of the sending rate by at most $1/(2N)$. For example, using five parallel TCP sessions, the response to a single packet drop event is to reduce the total sending rate by $1/(2 \times 5)$, or $1/10$, as compared to the response from a single TCP session, where a single packet drop event would reduce the sending rate by $1/2$.

A simulated version of five parallel sessions in a 10 Gbps session is shown in Figure 3.

Figure 3: Parallel TCP Simulation:
Single vs Parallel Streams

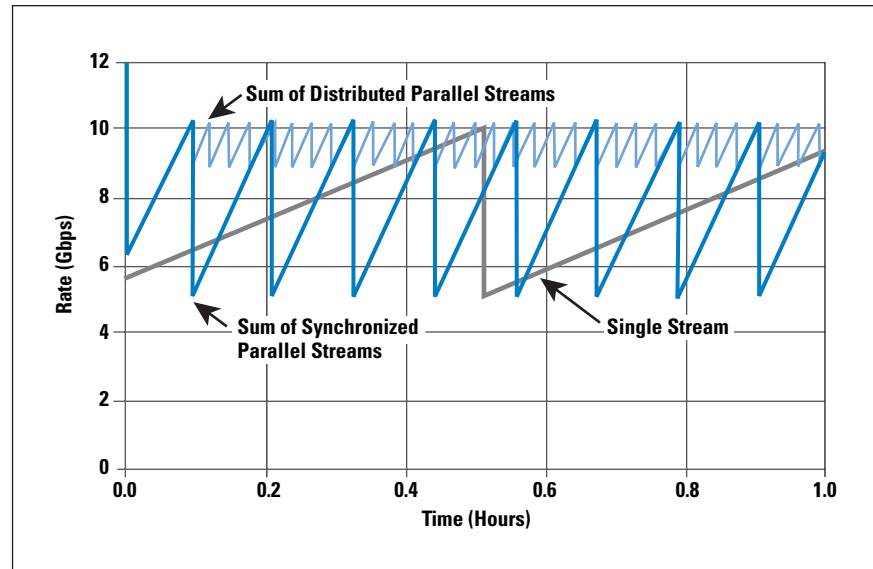


The essential characteristic of the aggregate flow is that under lossless conditions the data flow of N parallel sessions increases at a rate N times faster than a single session in congestion avoidance mode. Also the response to an isolated loss event is that of rate halving of a single flow, so that the total flow rate under ideal conditions is between R and $R \times (2N - 1)/2N$, or a long-term average throughput of $R \times (4N - 1)/4N$. For $N = 100$ our theoretical 10-Gbps connection could now operate at 9.9 Gbps.

Of course practice is different from theory, and a considerable amount of work has looked at the performance of parallel TCP under various conditions, in terms of both maximizing throughput and choosing the most efficient number of parallel active streams to use. Part of the problem is that although simple simulations, such as that used to generate Figure 4, tend to evenly distribute each of the parallel sessions to maximize the throughput, there is the more practical potential that the individual sessions self-synchronize. Because the parallel sessions have a similar range of window sizes, it is possible that at a given point in time a similar number of packets will be in the network path from each stream. If the packet drop event is a multiple packet drop event, such as a tail-drop queue, then it is entirely feasible that numerous parallel streams will experience packet loss simultaneously, and there is the consequential potential for the streams to fall into synchronization.

The two extremes, evenly distributed and tightly synchronized multiple streams, are indicated in Figure 4. The average throughput of parallel synchronized streams is the same as a single stream over extended periods in this simulation, and both are certainly far worse than an evenly distributed set of parallel streams.

Figure 4: Comparison of Parallel TCP:
Synchronized and Distributed
Streams



One way to address this problem is to reunite these parallel streams into a single controlled stream that exhibits the same characteristics as evenly spread parallel streams. This approach, MulTCP, is considered in the next section.

If all this analysis of parallel TCP streams sounds a little academic and unrelated to networking today, it is useful to note that many *Internet Service Providers* (ISPs) currently see *BitTorrent* traffic as their highest-volume application. BitTorrent is a peer-to-peer protocol that undertakes transfer of datasets using a highly parallel transfer technique. Under BitTorrent the original dataset is split into blocks, each of which can be downloaded in parallel. The subtle twist here is that the individual sessions do not have the same source points, and the host may take feeds from many different sources simultaneously, as well as offering itself as a feed point for the already downloaded blocks. This behavior exploits the peer-to-peer nature of these networks to a very high extent, potentially not only exploiting parallel TCP sessions for speed gains, but also exploiting diverse network paths and diverse data sources to avoid single path congestion. Considering its effectiveness in terms of maximizing transfer speeds for high-volume datasets and its relative success in truly exploiting the potential of peer-to-peer networks—and of course the dramatic acceptance of BitTorrent and its extensive use—BitTorrent probably merits closer examination, but perhaps that is for another time and an article of its own.

Very High Speed Serial TCP

The other general form of approach is to reexamine the current TCP control algorithm to see if there are parameter or algorithm changes that could allow TCP to undertake a better form of rate adaptation to these high-capacity, long-delay network paths. The aim here is to achieve a good congestion response algorithm that does not amplify transient congestion conditions into sustained disaster areas, while at the same time being able to support high-speed data transfers, thereby making effective use of all available network capacity.

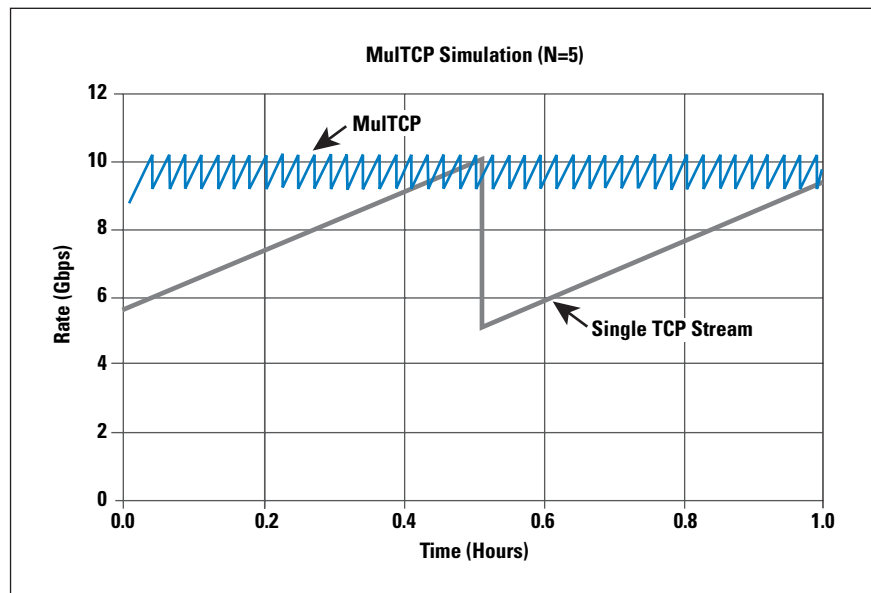
We also want TCP to behave sensibly in the face of other TCP sessions, so that it can share the network with other TCP sessions fairly.

MulTCP

The first of these approaches is *MulTCP*^[1], which is a single TCP stream that behaves in a manner equivalent to N parallel TCP sessions, where the virtual sessions are evenly distributed in order to achieve the optimal outcome in terms of throughput. The essential changes to TCP are in congestion avoidance mode and the reaction of packet loss. In congestion avoidance mode MulTCP increases its congestion window by N segments per RTT, rather than the default of a single segment. Upon packet loss, MulTCP reduces its window by $W/(2N)$, rather than the default of $W/2$. MulTCP uses a slightly different version of slow start, increasing its window by 3 segments per received ACK, rather than the default value of 2.

MulTCP represents a simple change to TCP that does not depart radically from the TCP congestion control algorithm. Of course when choosing an optimal value for N , some understanding of the network characteristics would help. If the value for N is too high, the MulTCP session has a tendency to claim an unfair amount of network capacity, but if the value is too low, it does not necessarily take full advantage of available network capacity. Figure 5 shows MulTCP compared to a simulation of an equivalent number of parallel TCP streams and a single TCP stream ($N = 5$ in this particular simulation).

Figure 5: MulTCP



Good as this is, there is the lingering impression that we can do better. It would be better not to have to configure the number of virtual parallel sessions; it would be better to support fair outcomes when competing with other concurrent TCP sessions over a range of bandwidths; and it would be better to have a wide range of scaling properties.

There is no shortage of options here for fine-tuning various aspects of TCP to meet some of these preferences, ranging from adaptations applied to the TCP rate control equation to approaches that view the loading onto the network as a power spectrum problem.

HighSpeed TCP

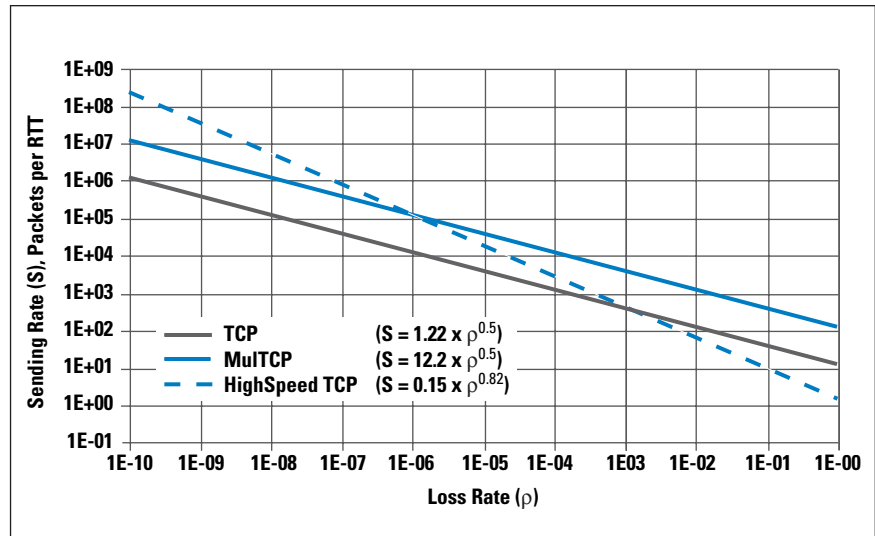
Another approach, described in [2], “HighSpeed TCP for Large Congestion Windows” looks at this from the perspective of the TCP rate equations, developed by Sally Floyd at ICIR.

When TCP operates in congestion avoidance mode at an average speed of W packets per RTT, then the number of packets per RTT varies between $(2/3)W$ and $(4/3)W$. Each cycle takes $(2/3)W$ RTT intervals, and the number of packets per cycle is therefore $(2/3)W^2$ packets. This result implies that the rate can be sustained at W packets per RTT as long as the packet loss rate is 1 packet loss per cycle, or a loss rate, ρ , where $\rho = 1/((2/3)W^2)$. Solving this equation for W gives the average packet rate per RTT of $W = \sqrt{(1.5)/(\rho)}$. The general rate function for TCP, R , is therefore: $R = (MSS/RTT) \times (\sqrt{(1.5)/(\rho)})$, where MSS is the TCP packet size.

Taking this same rate equation approach, what happens for N multiple streams? The ideal answer is that the parallel streams operate N times faster at the same loss rate, or, as a rate equation the number of packets per RTT, W_N , can be expressed as $W_N = N(\sqrt{(1.5)/(\rho)})$, and each TCP cycle is compressed to an interval of $(2/3) (W_N^2/N^2)$.

But perhaps the desired response is not to shift the TCP rate response by a fixed factor of N —as is the intent with MulTCP—but to adaptively increase the sending rate through increasing values of N as the loss rate falls. The proposition made by HighSpeed TCP is to use a TCP response function that preserves the fixed relationship between the logarithm of the sending rate and the logarithm of the packet loss rate, but alters the slope of the function, such that TCP increases its congestion avoidance increment as the packet loss rate falls. This relationship is shown in Figure 6 where the log of the sending rate is compared to the log of the packet loss rate. MulTCP preserves the same relationship between the log of the sending rate and the log of the packet loss rate, but alters the offset, whereas changing the value of the exponent of the packet loss rate causes a different slope in the rate equation.

Figure 6: TCP Response Functions



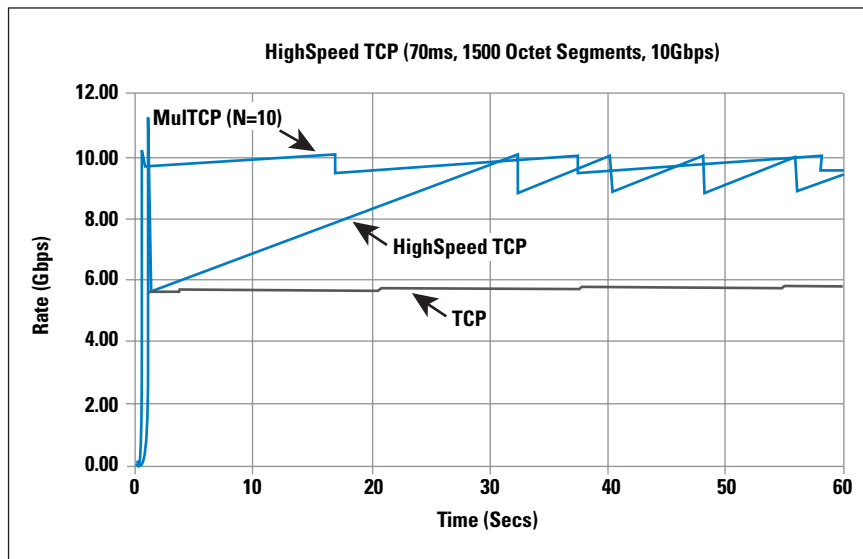
One way to look at the HighSpeed TCP proposal is that it operates in the same fashion as a turbocharger on an engine; the faster the engine is running, the higher the turbo-charged boost to the normal performance of the engine. Below a certain threshold value the TCP congestion avoidance function is unaltered, but when the packet loss rate falls below a certain threshold value then the higher speed congestion avoidance algorithm is invoked. The higher-speed rate equation proposed by HighSpeed TCP is based on achieving a transfer rate of 10 Gbps over a 100-ms latency path with a packet loss rate of 1 in 10 million packets. Working backward from these parameters gives us a rate equation for W , the number of packets per RTT interval of $W = 0.12/\rho^{0.835}$, approximately equivalent to a MulTCP session where the number of parallel sessions, N , is raised as the TCP rate increases.

This result can be translated into two critical parameters for a modified TCP: the number of segments to be added to the current window size for each lossless RTT time interval, and the number of segments to reduce the window size in response to a packet loss event. Conventional TCP uses values of 1 and $(1/2)W$, respectively. The HighSpeed TCP approach increases the congestion window by 1 segment for TCP transfer rates up to 10 Mbps, but then uses an increase of some 6 segments per RTT for 100 Mbps, 26 segments at 1 Gbps and 70 segments at 10 Gbps. In other words the faster the TCP rate that has already been achieved, then the greater the rate acceleration. Highspeed TCP also advocates a smaller multiplicative decrease in response to a single packet drop, so that at 10 Mbps the multiplier would be $1/2$, at 100 Mbps the multiplier is $1/3$, at 1 Gbps it is $1/5$, and at 10 Gbps it is set to $1/10$.

What does this process look like? Figure 7 shows a HighSpeed TCP simulation. What is not easy to discern is that during congestion avoidance HighSpeed TCP opens its sending window in increments of 53 through 64 segments each RTT interval, making the rate curve slightly upward during this window expansion phase.

HighSpeed TCP manages to recover from the initial rate halving from slow start in about 30 seconds, and operates at an 8-second cycle, as compared to the 38-minute cycle of a single TCP stream, or a 10-stream MultTCP session that operates at a 21-second cycle.

Figure 7: HighSpeed TCP Simulation



One other aspect of this work concerns the so-called slow start algorithm, which at these speeds is not really slow at all. The final RTT interval in our scenario has TCP attempting to send an additional 50 MB of data in just 70 ms, meaning an additional 33,333 packets are pushed into the network queues. Unless the network path is completely idle at this point, it is likely that hundreds—if not thousands—of these packets will be dropped in this step, pushing TCP back into a restart cycle. HighSpeed TCP has proposed a limited slow start to accompany HighSpeed TCP that limits the inflation of the sending window to a fixed upper rate per RTT to avoid this problem of slow start overwhelming the network and causing the TCP session to continually restart. Other changes for HighSpeed TCP are to extend the limit of three duplicate ACKs before retransmitting to a higher value, and a smoother recovery when a retransmitted packet is itself dropped.

Scalable TCP

Of course HighSpeed TCP is not the only offering in the high-performance TCP stakes.

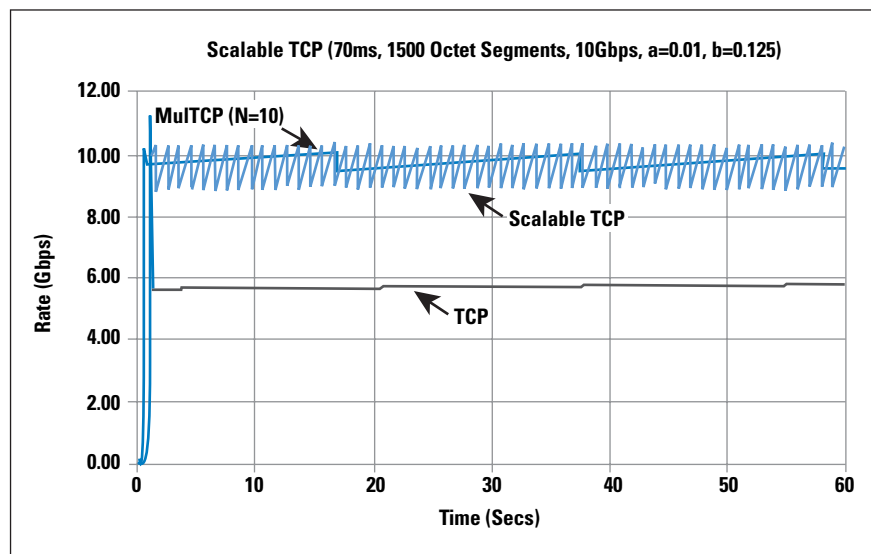
Scalable TCP^[3], developed by Tom Kelly at Cambridge University, attempts to break the relationship between TCP window management and the RTT time interval. It does this by noting that in “conventional” TCP, the response to each ACK in congestion avoidance mode is to inflate the sender’s congestion window size (*cwnd*) by $(1/cwnd)$, thereby ensuring that the window is inflated by 1 segment each RTT interval. Similarly the window halving on packet loss can be expressed as a reduction in size by $(cwnd/2)$. Scalable TCP replaces the additive function of the window size by the constant value *a*.

The multiplicative decrease is expressed as a fraction b , which is applied to the current congestion window size.

In Scalable TCP, for each ACK the congestion window is inflated by the constant value a , and upon packet loss the window is reduced by the fraction b . The relative performance of Scalable TCP as compared to conventional TCP and MulTCP is shown in Figure 8.

The essential characteristic of Scalable TCP is the use of a multiplicative increase in the congestion window, rather than a linear increase, effectively creating a higher frequency of oscillation of the TCP session, probing upward at a higher rate and more frequently than HighSpeed TCP or MulTCP. The frequency of oscillation of Scalable TCP is independent of the RTT interval, and the frequency can be expressed as $f = \log(1 - b) / \log(1 + a)$. In this respect, longer networks paths exhibit similar behavior to shorter paths at the bottleneck point. Scalable TCP also has a linear relationship between the log of the packet loss rate and the log of the sending rate, with a greater slope of HighSpeed TCP.

Figure 8: Scalable TCP



BIC and CUBIC

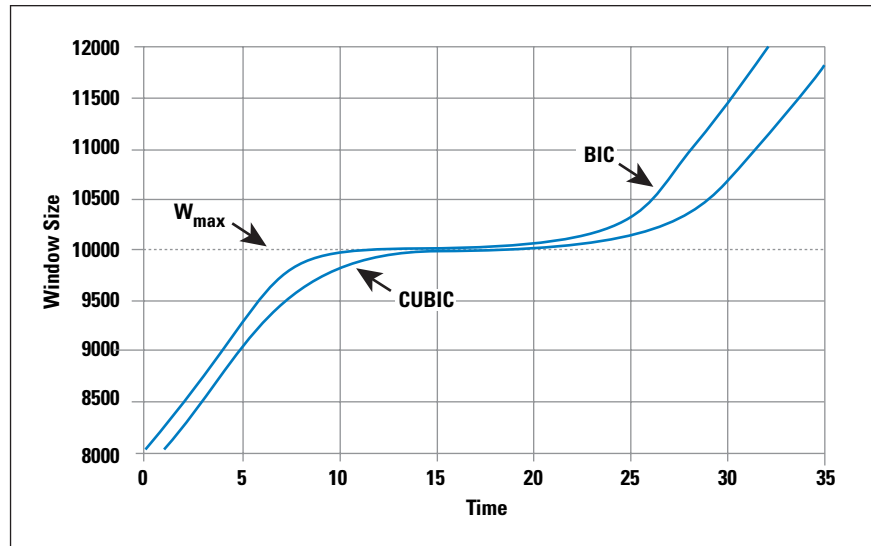
The common concern here is that TCP underperforms in those areas of application where there is a high bandwidth-delay product. The common problem observed here is that the additive window inflation algorithm used by TCP can be very inefficient in long-delay, high-speed environments. As can be seen in Figure 10, the ACK response for TCP is a congestion window inflation operation where the amount of inflation of the window is a function of the current window size and some additional scaling factor.

Binary Increase Congestion Control (BIC)^[4] takes a different view, by assuming that TCP is actively searching for a packet sending rate that is on the threshold of triggering packet loss, and uses a binary chop search algorithm to achieve this efficiently.

When BIC performs a window reduction in response to packet drop, it remembers the previous maximum window size, as well as the current window setting. With each lossless RTT interval BIC attempts to inflate the congestion window by one half of the difference between the current window size and the previous maximum window size. In this way BIC quickly attempts to recover from the previous window reduction, and, as BIC approaches the old maximum value, it slows down its window inflation rate, halving its rate of window inflation each RTT. This process is not quite as drastic as it may sound, because BIC also uses a maximum inflation constant to limit the amount of rate change in any single RTT interval. The resultant behaviour is a hybrid of a linear and a non-linear response, where the initial window inflation after a window reduction is a linear increase, but as the window approaches the previous point where packet loss occurred the rate of window increase slows down. BIC uses the complementary approach to window inflation when the current window size passes the previous loss point. Initially further window inflation is small, and the size of the window inflation value doubles for each RTT, up to a limit value, beyond which the window inflation is once more linear.

BIC can be too aggressive in low RTT networks and in slower speed situations, leading to a refinement of BIC, namely CUBIC^[5]. CUBIC uses a third-order polynomial function to govern the window inflation algorithm, rather than the exponential function used by BIC. The cubic function is a function of the elapsed time since the previous window reduction, rather than the implicit use by BIC of an RTT counter, so that CUBIC can produce fairer outcomes in a situation of multiple flows with different RTTs. CUBIC also limits the window adjustment in any single RTT interval to a maximum value, so the initial window adjustments after a reduction are linear. Here the new window size, W , is calculated as $W = C(t - K)^3 + W_{\max}$, where C is a constant scaling factor, t is the elapsed time since the last window reduction event, W_{\max} is the size of the window prior to the most recent reduction and K is a calculated value: $K = (W_{\max} \beta / C)^{1/3}$. This function is more stable when the window size approaches the previous window size W_{\max} . The use of a time interval rather than an RTT counter in the window size adjustment is intended to make CUBIC more sensitive to concurrent TCP sessions, particularly in short RTT environments.

Figure 9 shows the relative adjustments for BIC and CUBIC, using a single time base. The essential difference between the two algorithms is evident in that the CUBIC algorithm attempts to reduce the amount of change in the window size when near the value where packet drop was previously encountered.

Figure 9: Window Adjustment
for BIC and CUBIC

Westwood

The “steady state” mode of TCP operation is one that is characterized by the “sawtooth” pattern of rate oscillation. The additive increase is the means of exploring for viable sending rates while not causing transient congestion events by accelerating the sending rate too quickly. The multiplicative decrease is the means by which TCP reacts to a packet loss event that is interpreted as being symptomatic of network congestion along the sending path.

BIC and CUBIC concentrate on the rate increase function, attempting to provide for greater stability for TCP sessions as they converge to a long-term available sending rate. The other perspective is to examine the multiplicative decrease function, to see if there is further information that a TCP session can use to modify this rate decrease function.

The approach taken by Westwood^[6], and a subsequent refinement, Westwood+^[7], is to concentrate on the halving by TCP of its congestion window in response to packet loss (as signaled by three duplicate ACK packets). The conventional TCP algorithm of halving the congestion window can be refined by the observation that the stream of return ACK packets actually provides an indication of the current bottleneck capacity of the network path, as well as an ongoing refinement of the minimum RTT of the network path. The Westwood algorithm maintains a bandwidth estimate by tracking the TCP acknowledgement value and the inter-arrival time between ACK packets in order to estimate the current network path bottleneck bandwidth. This technique is similar to the “Packet Pair” approach, and that used in the TCP Vegas. In the case of the Westwood approach the bandwidth estimate is based on the receiving ACK rate, and is used to set the congestion window, rather than the TCP send window. The Westwood sender keeps track of the minimum RTT interval, as well as a bandwidth estimate based on the return ACK stream. In response to a packet loss event, Westwood does not halve the congestion window, but instead sets it to the bandwidth estimate times the minimum RTT value.

If the current RTT equals the minimum RTT, implying that there are no queue delays over the entire network path, then the sending rate is set to the bandwidth of the network path. If the current RTT is greater than the minimum RTT, the sending rate is set to a value that is lower than the bandwidth estimate, and allows for additive increase to once again probe for the threshold sending rate when packet loss occurs.

The major concern here is the potential variation in inter-ACK timing, and although Westwood uses every available data and ACK pairing to refine the current bandwidth estimate, the approach also uses a low pass filter to ensure that the bandwidth estimate remains relatively stable over time. The practical result here is that the receiver may be performing some form of ACK distortion, such as a delayed ACK response, and the network path contains jitter components in both the forward and reverse direction, so that ACK sequences can arrive back at the sender with a high variance of inter-ACK arrival times. Westwood+ further refines this technique to account for a false high reading of the bandwidth estimate due to ACK compression, using a minimum measurement interval of the greater of the RTT or 50 ms.

The intention here is to ensure that TCP does not over-correct when it reduces its congestion window, so that the problems relating to the slow inflation rate of the window are less critical for overall TCP performance. The critical part of this work lies in the filtering technique that takes a noisy sequence of measurement samples and applies an anti-aliasing filter followed by a low-pass discrete-time filter to the data stream in order to generate a reasonably accurate available bandwidth estimate. This estimate is coupled with the minimum RTT measurement to provide a lower bound for the TCP congestion window setting following detection of packet loss and subsequent fast retransmit repair of the data stream. If the packet loss is caused by network congestion the new setting will be lower than the threshold bandwidth (lower by the ratio $RTT_{min} / RTT_{current}$), so that the new sending rate will also allow the queued backlog of traffic along the path to clear. If the packet loss is caused by media corruption, the RTT value will be closer to the minimum RTT value, in which case the TCP session-rate backoff is smaller, allowing for a faster recovery of the previous data rate.

Although this approach has direct application in environments where the probability of bit-level corruption is intermittently high, such as often encountered with wireless systems, it also has some application to the long-delay, high-speed TCP environment. The rate backoff of TCP Westwood is one that is based on the $RTT_{min} / RTT_{current}$ ratio, rather than rate halving in conventional TCP, or a constant ratio, such as used in MulTCP, allowing the TCP session to oscillate its sending rate closer to the achievable bandwidth rather than performing a relatively high-impact rate backoff in response to every packet loss event.

H-TCP

The observation made by the proponents of H-TCP^[9] is that better TCP outcomes on high-speed networks is achieved by modifying TCP behavior to make the time interval between congestion events smaller. The signal that TCP has taken up its available bandwidth is a congestion event, and by increasing the frequency of these events TCP will track this resource metric with greater accuracy. To achieve this tracking, the H-TCP proponents argue that both the window increase and decrease functions may be altered, but in deciding whether to alter these functions, and in what way, they argue that a critical factor lies in the level of sensitivity to other concurrent network flows, and the ability to converge to stable resource allocations to various concurrent flows.

“While such modifications might appear straightforward, it has been shown that they often negatively impact the behaviour of networks of TCP flows. High-speed TCP and BIC-TCP can exhibit extremely slow convergence following network disturbances such as the start-up of new flows; Scalable-TCP is a multiplicative-increase multiplicative-decrease strategy and as such it is known that it may fail to converge to fairness in drop-tail networks.”

Work-in-progress: **draft-leith-tcp-htcp-01.txt**

H-TCP argues for minimal changes to the window control functions, observing that in terms of fairness a flow with a large congestion window should, in absolute terms, reduce the size of their window by a larger amount than smaller-sized flows, as a means of readily establishing a dynamic equilibrium between established TCP flows and new flows entering the same network path.

H-TCP proposes a timer-based response function to window inflation, where for an initial period, the existing value of one segment per RTT is maintained, but after this period the inflation function is a function of the time since the last congestion event, using an order-2 polynomial function where the window increment in each RTT interval, $\alpha = (\frac{1}{2}T^2 + 10T + 1)$, where T is the elapsed time since the last packet loss event. This equation is further modified by the current window reduction factor β where $\alpha' = 2 \times (1 - \beta) \times \alpha$.

The window reduction multiplicative factor, β , is based on the variance of the RTT interval, and β is set to RTT_{min} / RTT_{max} for the previous congestion interval, unless the RTT has a variance of more than 20 percent, in which case the value of $\frac{1}{2}$ is used.

H-TCP appears to represent a further step along the evolutionary path for TCP, taking the adaptive window inflation function of HighSpeed TCP, using an elapsed timer as a control parameter as was done in Scalable TCP, and using the RTT ratio as the basis for the moderation of the window reduction value from Westwood.

FAST

FAST^[10] is another approach to high-speed TCP. FAST is probably best viewed in context in terms of the per packet response of the various high speed TCP approaches, as indicated in the following Control and Response table:

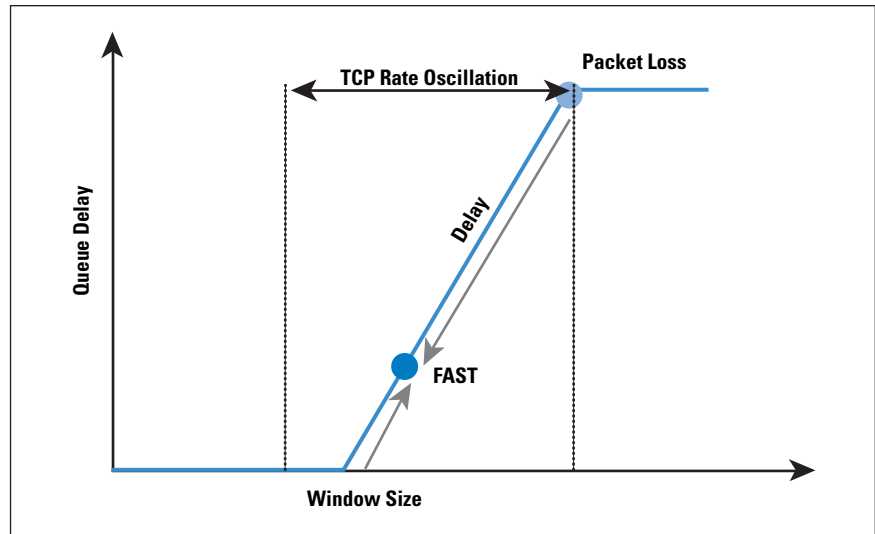
Type	Control Method	Trigger	Response
TCP	AIMD(1,0.5)	ACK response Loss response	$W = W + 1/W$ $W = W - W \times 0.5$
MulTCP	AIMD(N,1/2N)	ACK response Loss response	$W = W + N/W$ $W = W - W \times 1/2N$
HighSpeed TCP	AIMD(a(w), b(w))	ACK response Loss response	$W = W + a(W)/W$ $W = W - W \times b(W)$
Scalable TCP	MIMD(1/100, 1/8)	ACK response Loss response	$W = W + 1/100$ $W = W - W \times 1/8$
FAST	RTT Variation	RTT	$W = W \times (\text{base RTT}/\text{RTT}) + \alpha$

All these approaches share a common structure of window adjustment, where the sender's window is adjusted according to a control function and a flow gain. TCP, MulTCP, HighSpeed TCP, Scalable TCP, BIC, CUBIC, Westwood, and H-TCP all operate according to a congestion measure that is based on ACK clocking and a packet loss trigger. What is happening in these models is that a bottleneck point on the network path has reached a level of saturation such that the bottleneck queue is full and packet loss is occurring. It is noted that the build up of the queue prior to packet loss would have caused a deterioration of the RTT.

This fact leads to the observation made by FAST, that another form of congestion signalling is one that is based on RTT variance, or cumulative queuing delay variance. FAST is based on this latter form of congestion signalling.

FAST attempts to stabilize the packet flow at a rate that also stabilizes queue delay, by basing its window adjustment, and therefore its sending rate, such that the RTT interval is stabilized. The window response function is based on adjusting the window size by the proportionate amount that the current RTT varies from the average RTT measurement. If the current RTT is lower than the average, then window size is increased, and if the current RTT is higher then window size is decreased. The amount of window adjustment is based on the proportionate difference between the two values, leading to the observation that FAST exponentially converges to a base RTT flow state. By comparison, conventional TCP has no converged state, but instead oscillates between the rate at which packet loss occurs and some lower rate (Figure 10).

Figure 10: TCP Response Function vs. FAST



FAST maintains an exponential weighted average RTT measurement and adjusts its window in proportion to the amount by which the current RTT measurement differs from the weighted average RTT measurement. It is harder to provide a graph of a simulation of FAST as compared to the other TCP methods, and the more instructive material has been gathered from various experiments using FAST.

XCP — End-to-End and Network Signalling

It is possible to also call in the assistance of the routers on the path and call on them to mark packets with signaling information relating to current congestion levels. This approach was first explored with the concept of ECN, or *Explicit Congestion Notification*, and has been generalized into a transport flow control protocol, called XCP,^[11] where feedback relating to network load is based on explicit signals provided by routers relating to their relative sustainable load levels. Interestingly this digresses from the original design approach of TCP, where the TCP signaling is set up as effectively a heartbeat signal being exchanged by the end systems, and the TCP flow control process is based upon interpretation of the distortions of this heartbeat signal by the network.

XCP appears to be leading into a design approach where the network switching elements play an active role in end-to-end flow control, by effectively signalling to the end systems the current available capacity along the network path. This setup allows the end systems to respond rapidly to available capacity by increasing the packet rate to the point where the routers along the path signal that no further capacity is available, or to back off the sending rate when the routers along the path signal transient congestion conditions.

Whether such an approach of using explicit router-to-end host signals leads to more efficient very high-speed transport protocols remains to be determined, however.

Where Next?

The basic question here is whether we have reached some form of fundamental limitation of the TCP window-based congestion control protocol, or whether it is a case that the window-based control system remains robust at these speeds and distances, but that the manner of control signalling will evolve to adapt to an ever-widening range of speed extremes in this environment.

Rate-based pacing, as used in FAST can certainly help with the problem of the problem of guessing what are “safe” window inflation and reduction increments, and it is an open question as to whether it is even necessary to use a window inflation and deflation algorithm or whether it would be more effective to head in other directions, such as rate control, RTT stability control or adding additional network-generated information into the high-speed control loop. Explicit router-based signaling, such as described in XCP, allows for quite precise controls over the TCP session, although what is lost there is the adaptive ability to deploy the control system over any existing IP network.

However, across all these approaches, the basic TCP objectives remain the same: what we want is a transport protocol that can use the available network capacity as efficiently as possible—and as quickly as possible—minimizing the number of retransmissions and maximizing the effective data throughput.

We also want a protocol that can adapt to other users of the network, and attempt to fairly balance its use with competing claims for network resources.

The various approaches that have been studied to date all represent engineering compromises in one form or another. In attempting to optimize the instantaneous transfer rate the congestion control algorithm may not be responsive to other concurrent transport sessions along the same path. Or in attempting to optimize fairness with other concurrent sessions, the control algorithm may be unresponsive to available network path capacity. The control algorithm may be very unresponsive to dynamic changes in the RTT that may occur during the session because of routing changes in the network path. Which particular metrics of TCP performance are critical in a heterogeneous networking environment is a topic where we have yet to see a clear consensus emerging from the various research efforts.

However, we have learned a few things about TCP that form part of this consideration of where to take TCP in this very-high-speed world:

- The first lesson is that TCP has been so effective in terms of overall network efficiency and mutual fairness because everyone uses much the same form of TCP, with very similar response characteristics. If we all elected to use radically different control functions in each of our TCP implementations then it appears likely that we would have a poorly performing chaotic network subject to extended conditions of complete overload and inefficient network use.

- The second lesson is that a transport protocol does not need to solve media level or application problems. The most general form of transport protocol should not rely on characteristics of specific media, but should use specific responses from the lower layers of the protocol stack in order to function correctly as a transport system.
- The third lesson from TCP is that a transport protocol can become remarkably persistent and be used in contexts that were simply not considered in the original protocol design, so any design should be careful to allow generous margins of use conditions.
- The final lesson is one of fair robustness under competition. Does the protocol negotiate a fair share of the underlying network resource in the face of competing resource claims from concurrent transport flows?

Of all these lessons, the first appears to be the most valuable and probably the most difficult to put into practice. The Internet works as well as it does today largely because we all use the much same transport control protocol. If we want to consider some changes to this control protocol to support higher-speed flows over extended latency, then it would be perhaps reasonable to see if there is a single control structure and a single protocol that we can all use.

So deciding on a single approach for high-speed flows in the high-speed Internet is perhaps the most critical part of this entire agenda of activity. It is one thing to have a collection of differently controlled packet flows each operating at megabits-per-second flow rates on a multi-gigabit network, but it is quite a frightening prospect to have all kinds of different forms of flows each operating at gigabits per second on the same multigigabit network. If we cannot make some progress in reaching a common view of a single high-speed TCP control algorithm then it may indeed be the case that none of these approaches will operate efficiently in a highly diverse high-speed network environment.

Acknowledgment

I must acknowledge the patient efforts of Larry Dunn in reading through numerous iterations of this article, correcting the text and questioning some of my wilder assertions. Thanks Larry.

However, whatever errors may remain are, undoubtedly, all mine.

Further Reading

There is a wealth of reading on this topic, and here any decent search engine can assist. However if you are interested in this topic and want a starting reference that describes it in a very careful and structured manner, then I can recommend the following two sources as a good way to start exploring this topic to gain an overview of the current state of the art in this area:

- “HighSpeed TCP for Large Congestion Windows,” S. Floyd, RFC 3649, December 2003.

Floyd’s treatment of this topic is precise, encompassing, and wonderfully presented. If only all RFCs were of this quality.

- Proceedings of the Workshops on Protocols for Fast Long-Distance Networks.

These workshops have been held in:

2003: <http://datatag.web.cern.ch/datatag/pfldnet2003/>

2004: <http://www.didc.lbl.gov/PFLDnet2004/program.htm>

2005: <http://www.ens-lyon.fr/LIP/RESO/pfldnet2005/>

References

- [1] “Differentiated End-to-End Internet Services Using a Weighted Proportional Fair Sharing TCP,” J. Crowcroft and P. Oechslin, ACM SIGCOMM *Computer Communication Review*, Volume 28, No. 3, pp. 53–69, July 1998.
- [2] “HighSpeed TCP for Large Congestion Windows,” S. Floyd, RFC 3649, December 2003.
- [3] “Scalable TCP: Improving Performance in High-Speed Wide Area Networks,” T. Kelly, ACM SIGCOMM *Computer Communication Review*, Volume 33, No. 2, pp. 83–91, April 2003.
- [4] “Binary Increase Congestion Control (BIC) for Fast Long-Distance Networks,” L. Xu, K. Harfoush, and I. Rhee, *Proceedings of IEEE INFOCOMM 2004*, March 2004.
- [5] “CUBIC: A New TCP-Friendly High-Speed TCP Variant,” I. Rhee, L. Xu, <http://www.csc.ncsu.edu/faculty/rhee/export/bitcp/cubic-paper.pdf>, February 2005.
- [6] “TCP Westwood: Congestion Window Control Using Bandwidth Estimation,” M. Gerla, M. Y. Sanadidi, R. Wang, A. Zanella, C. Casetti, and S. Mascolo, *Proceedings of IEEE Globecom 2001*, Volume 3, pp. 1698–1702, November 2001.
- [7] “Linux 2.4 Implementation of Westwood+ TCP with Rate-Halving: A Performance Evaluation over the Internet,” A. Dell’Aera, L. A. Greco, and S. Mascolo, Tech. Rep. No. 08/03/S, Politecnico di Bari, http://deecal03.poliba.it/mascolo/tcp%20westwood/Tech_Rep_08_03_S.pdf
- [8] “End-to-end Internet packet dynamics,” V. Paxson, *Proceedings of ACM SIGCOMM 97*, pp. 139–152, 1997.

- [9] “H-TCP: TCP Congestion Control for High Bandwidth-Delay Product Paths,” D. Leith, R. Shorten, Work in Progress, June 2005. Internet Draft: **draft-leith-tcp-htcp-00.txt**
- [10] “FAST TCP: Motivation, Architecture, Algorithms, Performance,” C. Jin, X. Wei, and S. H. Low, *Proceedings of IEEE INFOCOM 2004*, March 2004.
- [11] “Congestion Control for High Bandwidth-Delay Product Networks,” D. Katabi, M. Handley, and C. Rohrs, ACM SIGCOMM *Computer Communication Review*, Volume 32, No. 4, pp. 89–102, October 2002.
- [12] “TCP Performance,” Geoff Huston, *The Internet Protocol Journal*, Volume 3, No. 2, June 2000.
- [13] “The Future for TCP,” Geoff Huston, *The Internet Protocol Journal*, Volume 3, No. 3, September 2000.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for almost two decades, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector, and has served time with Telstra, where he was the Chief Scientist in the company’s Internet area. Geoff is currently the Internet Research Scientist at the Asia Pacific Network Information Centre (APNIC). He has been a member of the Internet Architecture Board, and currently co-chairs three Working Groups in the IETF. He is author of several Internet-related books. E-mail: **gih@apnic.net**

How Instant Messaging Is Transforming the Enterprise Network

by David Strom

Instant Messaging (IM) has come of age and is close to becoming one of those protocols that offers something for everyone. Once the province of chatty teens looking to replace phone conversations with electronic ones, IM is now a corporate mainstay and part of a new breed of applications that are built around “presence detection,” the ability to determine when someone—or something—is online and available to communicate.

Indeed, IM is rapidly spreading across the corporate world and becoming an able replacement for overflowing voicemail and e-mail inboxes that are clogged with spam and buried in irrelevant and non-time-sensitive postings. If you must get through to a busy corporate executive, IM is becoming the fastest and most effective method of communicating. Move over BlackBerry.

IM offers several benefits today, having taken some lessons learned by other Internet protocols of the past. First, it has a solid user and developer base. Second, it has a relatively simple building-block structure like the best of Internet protocols, with well-defined clients and servers. Third, interoperability efforts are beginning to pay off among the leading independent and private IM systems. Fourth, open-source rules are making inroads in all the right places. Fifth, Microsoft is a friend (for once) of IM and helping matters—rather than playing its usual monopolist role in this space, the company is actually encouraging future developments and interoperability. Finally, a new collection of advanced applications is taking hold that will take advantage of the existing Internet and IM infrastructure and create some very sophisticated IM applications.

Let's examine more closely where IM originated, where it is going, and what the specific implications are for each of these developments and for networking professionals. As a warning, this article by its very nature takes some positions on products and vendors. These opinions are solely those of the author, and they represent nothing wider or more inclusive.

User Base

The IM servers are operated by either public network or private entities. The major difference between the two is that the public systems operate across the Internet and can be accessed by any users who download the appropriate client software and create their own identity. Message traffic is usually transmitted in plaintext and without any encryption whatsoever.

The private IM systems are usually maintained by a corporate IT department and operate behind firewalls; they offer message encryption, message retention, and archiving; prepopulated buddy lists that are integrated into the corporate authentication and directory servers; and better security and privacy that are specific to a particular set of corporate users. These private systems are not available to the public and are designed strictly for employee communications or communications among particular trading partners of the corporation.

The four most popular public IM systems are currently all in corporate hands: Microsoft, Yahoo, eBay/Skype, and AOL. Actually, we should make that five systems because AOL owns two separate networks, *AOL Instant Messenger* (AIM) and *I seek you* (ICQ). Introduced in November 1996, ICQ was actually the first general-purpose IM system combining presence or a list of contacts with the ability to send messages. Other popular systems include the open-source Jabber and Tencent QQ, the latter very popular in China. Estimates vary widely as to the total number of nonduplicated users—because many people have multiple accounts and use multiple systems—but it is safe to say that more than 150 million users are active across all these systems at any moment. The most recent estimates of active users are as follows:^[1]

IM System	Estimate of Active Users
AIM	53 million active users
ICQ	15 million active users
Skype	10 million active users
MSN Messenger	29 million active users
Yahoo Messenger	21 million active users
Jabber	13.5 million enterprise users
Tencent QQ	10 million active users

Why IM Is So Popular for Businesses

But these numbers are more about individuals using IM. They hide the real story over the past several years, the rise of IM as a solid enterprise communications tool. Corporate IM usage has skyrocketed the last several years, and one survey has found IM users in more than 50 percent of American corporations^[2]. As mentioned earlier, there are public and private IM systems. The vast majority of the private IM systems are for institutional use for communications inside a company or among several suppliers, customers, and other trading partners.

The largest players in the private IM space are Microsoft Office Live Communications Server and IBM/Lotus' Sametime, although Jabber Corporation (not to be confused with the Jabber Software Foundation) is also gaining a strong following. We will discuss more about the role of open source in a moment, but first let's examine the reasons why IM has become so popular among so many business users.

First, workers have become more mobile and more difficult to track down. As secretarial support disappears and voicemail becomes more the norm, you want to know when people are actually at their desk—or laptop—these days. Staffs are more far-flung, and the global village becomes a lot smaller when you use IM to “talk” to someone halfway across the planet and get an immediate response. Finding someone who is available requires more than just making a phone call or exchanging e-mail messages. IM automatically tells you who is available—and who is not—at any given hour of the day.

Second, e-mail is no longer the productivity tool it once was because pipes are clogged with spam, viruses, and phishing attacks. Getting a quick response—that is, within minutes—through e-mail now seems so quaint, so “last year.”

Third, IM enables better collaboration and a tighter sense of community. With IM, you can educate an entire team, give the team feedback in real time, develop relationships, and cement the team together. It is a nice antidote and countermeasure to connect all these home-based and remote workers.

Fourth, the next generation of IM is not just about text chats; it also offers solid integration with voice and video. Voice and video calling is now part of Microsoft, Yahoo, Apple, and AOL IM software as well as part of the Skype network, which pioneered the feature. These audio and video extensions are becoming more popular with the private Lotus and Microsoft systems as well.

Finally, the real-time features of IM and its ability to track someone down no matter where they are located are attractive to customers, partners, and suppliers that need a guaranteed method of communication. IM is becoming the critical technology ingredient for corporations that are looking for faster response times, tying their customers closer together, and enabling teleworkers to communicate across the globe.

Components

Following are some definitions and explanations for those unfamiliar with the world of IM. Every IM network is composed of clients, servers, and protocols to connect them.

Each IM client has three major pieces:

- A buddy list or roster of friends with whom you wish to communicate—The list is organized by groups that you specify, such as “friends,” “work colleagues,” “family,” and so forth. The list indicates who is online, who is available to talk to, and who is offline or blocked by the user from communicating. Users organize their buddies in different ways and have complete control over the categories, naming conventions, and the like.

- A separate window that shows the text chats in process—Users type in this window and view the responses of their correspondents.
- Any additional features for video and audio chats and for file transfers between users

The last item bears some further discussion. All major IM products are moving beyond their roots of simple text chats toward more integrated and sophisticated communications, including real-time voice and video calls. Indeed, the mixture of *Voice over IP* (VoIP) and IM is a potent and popular one, accounting for the rapid uptake in Skype's adoption around the world. To use Skype as an example (although Yahoo has begun offering similar phone calling features in its IM client, and the others are soon to follow), users can make phone calls to the land-line phone numbers for a few pennies per minute—even calls to numbers in other countries. This is part of its attraction, along with voice mailboxes that are attached to a particular IM username.

The IM server maintains the directory of user accounts and keeps track of who is online, and in most cases routes messages among users. The major difference between an IM server and a *Simple Mail Transfer Protocol* (SMTP) e-mail server is that the IM server operates in real time, sending messages back and forth between two users as they finish typing a line of text. The servers also pass information in real time as to the availability of various users in the directory, when they come online and change their “status” message.

Users can typically set their availability in one of many different modes:

- Online and ready to receive messages
- Away from the computer, in which case correspondents receive a message saying so (or whatever the user wishes to be displayed)
- Unavailable or offline
- Blocked from anyone's view for privacy reasons

This status message can be changed at the user's discretion and is one of the main attractions for teens and other hypercommunicators. You can actually track what people are doing (or at least, saying that they are doing), by monitoring their status messages. (I am at the beach, I am taking a nap, I am at lunch, I am having coffee, and so forth.) For my teenaged daughter, this is one way she documents her life and one way that her friends can keep track of her—having a cell phone is not enough! There are numerous third-party add-ins to enhance your away message with clever graphics, hyperlinks to various Websites, and other effluvia as well.

The combination of instant access and persistent status indicator is at the core of why IM is such a powerful application. In a single window on your computer, you have a list of all your correspondents and can quickly determine who is online and who is not.

The blocking ability for some systems works universally, meaning that your presence is cloaked for everyone, as well as for specific users that you do not wish to communicate with or know your particular status, such as ex-spouses or ex-colleagues.

In most IM networks, you can be signed on from only one computer at any given moment. If you attempt to sign on from a second machine, you get an error message or your first computer is automatically logged out of the system. This is one way for the network to keep track of where you are located, because you can be in only one place at any given time.

Each server uses the TCP/IP Internet infrastructure and communicates with its clients over an assigned port number across the Internet. These ports can be blocked or proxied to different numbers, depending on the network administrator's policies toward IM traffic. Typical port numbers follow:

IM System	Port Numbers
ICQ	4000
AIM	5190–3
XMPP	5222–3
MSNP (Microsoft)	1863
YMSG (Yahoo)	5050
Skype	80, 443, and others

Notice an interesting thing about Skype's protocol: there is no single assigned port number. Users can set one of the ports in its configuration settings, but Skype uses a series of ports to communicate.^[3] This setup suggests several concerns, which we address next.

The Dark Side

Although these are all compelling reasons for the rise of IM across the corporate network, all is not constructive with IM. This section discusses problems specifically germane to Skype and problems with all IM products in general.

When the Skype client is installed on a computer, it picks a random port to communicate with other Skype computers, using what is believed to be a form of *Request for Comments* (RFC) 3489^[4]. This process is similar to many network-based games and peer-to-peer file-sharing products—no surprise because the developers of Skype worked on the Kazaa music file-sharing software. Because of its programming model, Skype is adept at traversing *Network Address Translation* (NAT) routers and can usually find a communications path to the outside world. Skype also encrypts all its message traffic, and this fact coupled with random port usage and its peer-to-peer programming model makes it look very similar to some malicious code that is unleashed across your network.

This is part of its charm and its challenge: network administrators who want to block Skype usage usually have a very difficult time figuring out how to do so[5], and may have to resort to third-party blocking products or clever configurations. One of the papers listed in [3] shows a way to block Skype using the popular open-source Squid caching proxy: not only do you have to prevent outbound *User Datagram Protocol* (UDP) connections over port 443, but you also must prevent connections to numeric IP addresses.

Although Skype has its own problems because of the way it is designed, there are several significant drawbacks to widespread adoption and deployment of any IM application. IM is not immune to infections, and just as its popularity is on the increase, so are ways to send malicious payloads and attacks too. What makes matters worse with IM versus say, e-mail, is its very instant nature: an infection can easily spread across a network in a matter of seconds, given that users are logged in, have long lists of users, and tend to think that any message coming from their respondents is more trusted than the average e-mail. In addition, Internet chat has long been a mechanism for controlling large-scale bot-nets of zombie computers, whose owners are unaware of such usage. Numerous virus authors have used exploits in Internet Relay Chat, for example, to control their villains across the Internet.

To avoid these problems, many corporations have either designed their own or are using one of several commercial IM protection products to screen incoming messages for particular patterns and methods of attack. The IM protection products work just like antivirus products work with e-mail messages: they download pattern files on a regular basis from a central server, and perform deep packet inspection across a perimeter to determine what is malicious and what is not.

Interoperability

Each public IM system is an island unto itself: users on one cannot easily communicate with users of another, unless one of two things happens:

- A user runs one of the multisystem client programs that allows them to sign in to multiple systems concurrently. Still, using these types of products means that just the user can communicate with his or her “buddies” across systems. Many mostly free products that enable this are available^[6].
- A private IM operator can combine more than one protocol inside the IM server application. This approach means that clients need not know or care about other IM protocols, such as using Microsoft’s Live Communications Server 2005^[7].

But variables are changing on the interoperability scene to make life better for IM users. First, efforts are under way among the major operators to form better relationships with each other:

In October 2005, Yahoo and Microsoft announced plans to introduce interoperability between MSN and Yahoo Messenger by mid-2006, using *Session Initiation Protocols* (SIPs). In December 2005, AOL and Google announced a strategic partnership deal where Google Talk users can talk with AIM and ICQ users provided they have an identity at AOL.

Second, both Microsoft and Apple have made efforts to include multi-protocol IM clients as part of their desktop operating systems. Apple's iChat in its latest Mac OS 10.4 Tiger, as an example, now supports AIM, Google Talk, and Jabber. Microsoft has announced plans to support other networks in its next release of Windows Vista, expected later this year.

Finally, the private IM systems of Microsoft and Lotus both support multiple IM protocols, and are widening their support for others, making them more useful for corporations.

Still, with all this activity, the IM interoperability scene is pretty poor: think where e-mail was in the early 1990s with custom-crafted gateways and the like so that an MCIMail user could send messages to a CompuServe user.

Setting up two systems to talk to each other is neither simple nor obvious, and each pair of systems must be done separately. So to add Google Talk to Trillian, a user would need to provide the server host name (**talk.google.com**) and port number (5222). (By the way, GoogleTalk has the most helpful instructions on how to set up a variety of third-party applications to connect to its servers.)

But that is not all—even if a user follows these instructions to set up cross-system connections, most systems can exchange only plaintext messages. Video and voice chats between disparate systems are not generally supported, although Apple's iChat has done the best job so far in this arena. And even if users take the multiple-client approach, the structure of their buddy lists is not always maintained and sometimes is presented in a single group of buddies, rather than separated into the groups that were specified when initially setting up the IM account.

The other concern for cross-systems interoperability is a lack of support for privacy or online status. All of the IM systems have the ability to create blacklists, or lists of users that cannot view your online status. These blacklists are not necessarily preserved when running the multiple client systems.

The Rise of Open Source

There is hope on the interoperability scene, however, and that hope is spelled *open source*. The Jabber group of programmers is growing, and the community is aggressively establishing a more pluralistic IM society. These steps revolve around software using the protocol called the *Extensible Messaging and Presence Protocol* (XMPP), the IETF's formalization of the core protocols created by the Jabber open-source community in 1999, and contained in four RFCs^[8, 9, 10, and 11].

Jeremie Miller developed the original Jabber server in 1998. Now the project has reached critical mass. Notable is the wide number of different server and client formulations that support XMPP. Jabber.com sells a commercial license, along with a combination of *General Public License* (GPL)-based licensed servers and other commercial versions. The project has supported the efforts of dozens of client implementations^[12]. Last year, support reached a new milestone with Google Talk and more recently the Gizmo Project using these protocols.

Numerous efforts are under way with these clients to extend basic IM functions into new areas, including providing more sophisticated and secure communications, the ability to have multiple identities presented (`david@strom.com` for work colleagues, `dstrom@gmail.com` for personal communications) from the same IM client, and support for more interoperable communications between Jabber and private IM systems.

At the heart of XMPP is the *Extensible Markup Language* (XML) constructs and basic protocols. The core “transport” layer for XMPP is an XML streaming protocol that makes it possible to exchange fragments of XML between any two network endpoints. Authentication and channel encryption happen at the XML streaming layer using other IETF-standard protocols for *Simple Authentication and Security Layer*^[13] and *Transport Layer Security*^[14].

Servers can connect to each other for interdomain communications, using the form of address for each user as `<user@domain>`—similar to SMTP e-mail, and in many cases, the IM address is the same as one's Internet e-mail address to simplify things.

What is notable about using XMPP is that RFC 3921 also makes it possible to separate the messaging and presence functions if desired (although most deployments offer both). This feature is helpful when building applications-to-applications messaging that does not involve users typing text messages to each other, such as a server sending a network operator an alert when it detects a problem.

The Jabber Software Foundation develops extensions to XMPP through a standards process centered on *Jabber Enhancement Proposals* (JEPs), similar to the RFC process^[15]. Currently, more than 30 active proposals have been developed, extending IM into bookmarks, delayed messaging, and other areas.

What Microsoft Is Doing

Microsoft is heavily involved in the IM scene in three important areas. The company operates one of the larger public IM networks, it includes an IM client as part of its Windows operating system, and it sells a private IM server that has some powerful interoperability features called *Live Communications Server* (LCS). What does this mean for the IM community? All good things. Microsoft's MSN and Skype are the more popular IM services outside of North America, and having Skype now a part of eBay is making Microsoft add competitive features such as voice and video chats to its public IM service. Microsoft has actually led the way on IM interoperability with LCS, a fact that can only motivate its competitors to include more pluralist IM offerings of their own. Finally, building in more support for IM in future versions of Windows will help popularize these applications even further.

It was not always this way. Earlier versions of Windows included something called Windows Messenger that was woefully underfeatured and had many bugs. But like so many early Microsoft efforts, technology has improved over time, and now the built-in software that comes with Windows is actually quite competitive with the public IM clients from AOL, Yahoo, and Skype.

Certainly, having Microsoft on one side and open-source efforts on the other is a nice way to encourage development and innovation in the IM arena, and we should expect more here in the future.

Building IM Applications

For most of this article we have addressed the one-to-one aspect of IM. However, IM is evolving into a much more important role, and that is one-to-many communications, and communications between applications instead of actual people. Many vendors have begun selling products in this space, and it is more interesting for several reasons:

First, IM is replacing other means for applications communications. It used to be the case that many network management applications used the *Simple Network Management Protocol* (SNMP) or SMTP protocols to send out their alerts. Now, many applications are using IM messages and taking advantage of the real-time nature of the protocol.

Second, the origins of IM go back to group chat sessions, so group collaboration tools make sense for new IM applications.

Third, even the closed public IM vendors have begun to open their programming interfaces, making it is easier for corporations to build new and sophisticated applications that are presence-aware, in some cases between two computer programs to communicate their status. AOL this year is one such example of opening its IM *application programming interface* (API) kimono, and of course Jabber has always been an open system that has helped lead more of these innovations.

One illustration is with the automotive giant Reynolds and Reynolds, which is using Jabber servers to monitor its own software status at the numerous automotive dealerships around the world. The IT department at Reynolds can quickly see if the company's software is down and take steps to get it working again.

Accredited Home Lenders is using IM to provide its loan brokers a secure and reliable means of communicating in real time with loan specialists to resolve problems with loan applications. And Ecreation built a virtual disk jockey for a Dutch radio station that also broadcasts over the Internet, allowing the station to take requests from listeners around the world through Microsoft's IM network.

Even traders have embraced IM. NetEnergy has been using IM for the past three years, and now negotiates trades between buyers and sellers of oil contracts using IM, decreasing errors and enabling faster communications.

Finally, IM figures prominently helping deaf and hard-of-hearing users communicate. In the era before IM, deaf users required a telephone relay operator to type the message to them and speak to the hearing callers. Go America has built a gateway to IM for its **1711.com** Website, so that deaf users can send messages directly to the operator.

Summary

We have tried to paint a comprehensive picture of what IM is and where it is going. Certainly, the amount of messaging traffic using the various IM protocols is impressive, and will continue to grow as these new applications are created and as more people discover the advantages of using IM. In several instances IM has replaced voicemail for most internal communications, particularly at high-tech companies and places where real-time communications is important. Although IM is not without its problems, there are ways to protect networks from infection and abuse.

For Further Reading

- [1] Nielsen//NetRatings, August 2005 study.
- [2] Osterman Research survey:
http://www.ostermanresearch.com/results/surveyresults_0905.htm
- [3] More details about the underlying Skype protocols, mechanisms for blocking its use, and other helpful tips and tricks for network administrators can be found at this page maintained by Salman A. Baset:
<http://www1.cs.columbia.edu/~salman/skype/index.html>
- [4] J. Rosenberg, J. Weinberger, C. Huitema, and R. Mahy, "STUN—Simple Traversal of User Datagram Protocol (UDP) Through Network Address Translators (NATs)," RFC 3489, March 2003.

- [5] A dissection of the Skype protocol along with suggestions about how to block its use can be found in this paper by P. Biondi and F. Desclaux: “Silver Needle in the Skype.”
<http://www.blackhat.com/presentations/bh-europe-06/bh-eu-06-biondi/bh-eu-06-biondi-up.pdf>
- [6] Adium and iChat for the Mac, Gaim for Windows and Linux, Trillian Pro for Windows, WebMessenger for Windows Mobile/Palm, and others.
- [7] Microsoft’s Live Communications Server 2005 includes its Public IM connector for an additional charge. Lotus’ Sametime has had AIM connectivity for several years, and will support other IM networks later this year.
- [8] P. Saint-Andre, ed., “Extensible Messaging and Presence Protocol (XMPP): Core,” RFC 3920, October 2004.
- [9] P. Saint-Andre, ed., “Extensible Messaging and Presence Protocol (XMPP): Instant Messaging and Presence,” RFC 3921, October 2004.
- [10] P. Saint-Andre, “Mapping the Extensible Messaging and Presence Protocol (XMPP) to Common Presence and Instant Messaging (CPIM),” RFC 3922, October 2004.
- [11] P. Saint-Andre, “End-to-End Signing and Object Encryption for the Extensible Messaging and Presence Protocol (XMPP),” RFC 3923, October 2004.
- [12] A list of software clients that support Jabber protocols can be found at:
<http://www.jabber.org/software/clients.shtml>
- [13] J. Myers, “Simple Authentication and Security Layer (SASL),” RFC 2222, October 1997.
- [14] T. Dierks and C. Allen, “The TLS Protocol Version 1.0,” RFC 2246, January 1999.
- [15] Jabber Enhancement proposals are listed at:
<http://www.jabber.org/jeps/>

DAVID STROM has been writing about Internet protocols and applications for nearly 20 years. Founding editor-in-chief for *Network Computing* magazine, he was most recently the editor-in-chief for tomshardware.com and related Websites. Strom has written two books on Internet e-mail (with the doyenne of POP, Marshall T. Rose) and home networking and thousands of magazine articles for most of the leading trade magazines in the IT, computing, and networking fields. He can be reached by e-mail at david@strom.com, or by IM: **davidstrom** (AIM and Skype) or **dstrom** (Yahoo, Google Talk, and MSN).

Letters to the Editor

Dear Editor,

In Russ White's "Working with IP Addresses" article (IPJ Volume 9, Number 1), he presents an example subnetting problem ("The Hardest Subnetting Problem") together with a worked solution. While useful as a reinforcement exercise for the rest of the article, care should be exercised before using the steps in the solution "as-is" in a real-world network configuration.

The main problem is that by packing subnets tightly together as shown, growth is restricted in order to guarantee that no address space is wasted. Worse, growth of host numbers on all but the smallest subnet requires renumbering of the subnet or all the smaller subnets allocated after it.

For example, the /26 subnet with 58 hosts will not accommodate more than another four hosts, less than 10-percent growth, without being renumbered.

Since renumbering a network is a nontrivial task even with the tools at our disposal, it is desirable to make it as infrequent as possible.^[1]

Allowing for growth will likely but not necessarily waste some address space, but it is preferable to frequent renumbering. It turns out that this example has alternative arrangements of subnets that would permit growth of some subnets without the need to renumber and would lessen the amount of renumbering when it is required.

Using realistic estimates of future hosts rather than current numbers is a simple measure to decrease the frequency of renumbering required. This would also make it obvious that the entire allocation is close to exhaustion and can be exhausted by the need to accommodate as little as six hosts on two subnets that are near full capacity.

Constraints on the supply of IPv4 address space limits how much growth can be accommodated and requires taking a shorter-term rather than longer-term view of growth. For private RFC 1918^[2] IP allocations (such as the one used in the example), this applies in only very large organisations, allowing a long-term view to be accommodated.

Unfortunately, the future is hard to predict with any degree of accuracy. In most cases needs for subnet allocation become gradually known over time rather than all at once. The consequences of incorrect estimation can be minimised by using an allocation scheme that allows for as much growth as possible in existing subnets while leaving as much room as possible for future allocations.

This scenario can be achieved by distributing the subnets evenly, weighted by size, across the available address space. The larger the subnet, the more room that needs to be left between it and other large networks. This is particularly important for subnets that are near to capacity. At least the sum of the sizes of neighbouring networks should be allowed. Space close to a network should be reserved for it to grow into, and the remaining space between can be allocated to smaller networks in a recursive fashion. Any allocations in the areas of likely growth should be reclaimable, and preferably these networks should be sparsely populated in order to limit the impact of renumbering on these networks. Working with a diagram of the address space, for example, a linear graph or a binary tree of the address space is a helpful aid.

A more systematic way of distributing the subnets evenly is to use *mirror-image* (MI) counting for allocating subnet numbers. This process is described in RFC 1219^[3], but note that some aspects of subnet addressing have altered since this RFC was written (see RFC 1878^[4]), so the description of mirror-image counting there and procedure text exclude subnet numbers that are now valid.

Using mirror-image counting is like normal counting starting from zero, except that the binary digits of the number are reversed. These numbers can be allocated as subnet numbers, starting from the most significant bit. Contrary to the example in RFC 1219, leading zeros (including the solitary zero in zero itself) should always be removed before the number is reversed.

Simplifying greatly, new subnets are allocated by incrementing the subnet number until a number is reached where a subnet of the required size can be accommodated or the subnet prefix becomes so long no subnets of the required size remain. If the prefix matches a common but shorter prefix, the subnet may be able to be allocated if we can lengthen the mask of the matching subnet prefix, freeing space from a previous allocation by reducing its maximum possible size. If the longest mask is always used when allocating subnets it is sufficient to just skip matching prefixes. Note that the null prefix is common with all subsequent prefixes until its subnet mask is made smaller, extending the prefix.

The mask chosen is preferably the longest for the required subnet size—but can be as short as the length of the subnet prefix, because it can be adjusted later: made shorter if the subnetwork grows beyond its mask (if no later allocation has been made) or longer if a subnet sharing its prefix is allocated or increases size. The host number ignoring the subnet part must be allocated from 1.

As the number is incremented it grows from right to left, progressively enumerating subnets in smaller sizes. Since subnet numbers grow from right to left and host numbers from left to right, collision is delayed between the two. Allocating subnets in descending order of size is preferable in this procedure because it tends to reduce fragmentation of the address space.

The following table shows an example allocation using the sorted number of hosts in the example:

MI Number	Subnet Prefix	Network Size	Network Number	Prefix	Last Host Number	Max Host Number
(null)	00	64	0	/26	58	62
1	10	64	128	/26	177	190
01	010	32	64	/27	93	94
11	1100	16	192	/28	206	206
001 matches subnet prefix 00						
101 matches subnet prefix 10						
011	01100	8	96	/29	99	102

Note that the /28 and the /29 can grow simply by changing their netmask. A better allocation is possible if the third and fourth hosts in the sorted list are interchanged. In this case the three smallest networks would be able to grow without renumbering. Shortening a netmask is a much simpler operation than renumbering.

Of course in the real world, needs for subnet allocation do not conveniently arrive sorted in ascending order. If it happened that one of the two largest subnets was the fifth requiring allocation, fragmentation of the address space would require renumbering one of the three smallest networks to recover an address block of the necessary size.

Another point that may be worth mentioning is that most modern hosts and routers allow for multiple subnets to share the same physical subnet, allowing two smaller subnets to cover a range of addresses that would otherwise receive a single larger allocation. For example, a 40-host subnet can be allocated a /27 and a /28 rather than a /26.

—Andrew Friedman, Sydney, Australia
rbns-w-ipj@yahoo.com.au

Ed: Readers may wish to also peruse RFC 3531^[5].

- [1] P. Ferguson and H. Berkowitz, “Network Renumbering Overview: Why Would I Want It and What Is It Anyway?” RFC 2071, January 1997.
- [2] Y. Rekhter, B. Moskowitz, D. Karrenberg, G. J. de Groot, and E. Lear, “Address Allocation for Private Internets,” RFC 1918, February 1996.
- [3] P. F. Tsuchiya, “On the Assignment of Subnet Numbers,” RFC 1219, April 1991.
- [4] T. Pummill and B. Manning, “Variable Length Subnet Table for IPv4,” RFC 1878, December 1995.
- [5] M. Blanchet, “A Flexible Method for Managing the Assignment of Bits of an IPv6 Address Block,” RFC 3531, April 2003.

The author responds:

Andrew is correct in stating that it is often better to try to account for future growth when assigning address space. There are many viable ways to allow for growth when allocating address spaces; hopefully, this topic will be covered more fully in a future article. I used the method in the article to illustrate how to employ the technique for working with IP addresses, rather than as an absolute best practice for allocating addresses.

—Russ White, Cisco Systems
riw@cisco.com

Dear Editor,

Russ White's article titled "Working with IP Addresses" was a nice refresher on how complicated working with IPv4 addresses has become. It should remind us all how we have gotten used to dealing with the operational expense of IPv4 address scarcity. The story about putting a frog in a pot of cold water comes to mind.

In any case, at the end of the article in the section titled "Working with IPv6 Addresses," I think the author tries too hard to fit the IPv6 address structure into the model for IPv4. Actually, it is a lot simpler.

The IPv6 address structure and textual representation was designed to avoid most of the complexities encountered in IPv4. The big differences follow:

- Addresses are represented in groups of hexadecimal digits instead of decimal digits. Hexadecimal avoids the need to convert the decimal digits to octal to find subnet boundaries. In hexadecimal there are four bits per character. This makes it easy to find the subnet boundary in an address; in many cases it is at a character boundary.
- Subnet prefix lengths are listed directly in decimal. There are no decimal subnet masks. This eliminates the need to convert decimal addresses to octets, convert the subnet masks to octets, apply the mask, and convert the result back to decimal—or to use the table and division methods described in the article.

The combination of these changes makes it much easier to work with IPv6 addresses. They are, of course, longer. The length has a few advantages besides a much larger Internet.

A byproduct of the larger address space is that most of the common subnet boundaries fall on hexadecimal digit boundaries; for example, using the example address in the article:

2002:FF10:9876:DD0A:9090:4896:AC56:0E01

The most common subnet boundary is 64 bits. The address and prefix is represented as:

2002:FF10:9876:DD0A:9090:4896:AC56:0E01/64

The subnet itself then follows:

2002:FF10:9876:DD0A::/64

The current common prefix allocated to a site is a /48. The site prefix is then:

2002:FF10:9876::/48

The current default allocation to an ISP is a /32. The ISP prefix is then:

2002:FF10::/32

These common prefix lengths can be derived directly without any need for decimal-to-octal conversions, table lookups, divisions, etc.

One of the other benefits of the larger addresses and a byproduct of IPv6 autoconfiguration is that the low-order 64 bits of an IPv6 address are reserved for the host address (called Interface Identifier in IPv6 terminology). This means that “The Hardest Subnetting Problem” described in the article is avoided completely. You can have as many hosts on a specific segment as you want in IPv6. There is no need to do this kind of calculation. This makes an initial network design trivial and, more importantly, makes later changes very easy. There is no need to redesign a subnet architecture because a few hosts need to be added to a subnet.

—Bob Hinden, Nokia
bob.hinden@nokia.com

The author responds:

Bob brings up many interesting points about IPv6, and the use of the IPv6 address space. While most IPv6 address spaces have prefix lengths that break on even octet boundaries today, we can’t always count on this, for all time, so it is always good to have techniques to work with situations where the prefix length is not on an octet boundary when they do occur. As for the last problem, it is true that in all cases the subnet is the set of octets excluding the last 64 bits. But if we move the problem up one level, and ask: “What is the most efficient way to allocate out an existing /48 so customer A can get 10 subnets, customer B can get 20 subnets, etc. ?” we can see the same problem could occur at the next higher level.

—Russ White, Cisco Systems
riw@cisco.com

Corrections

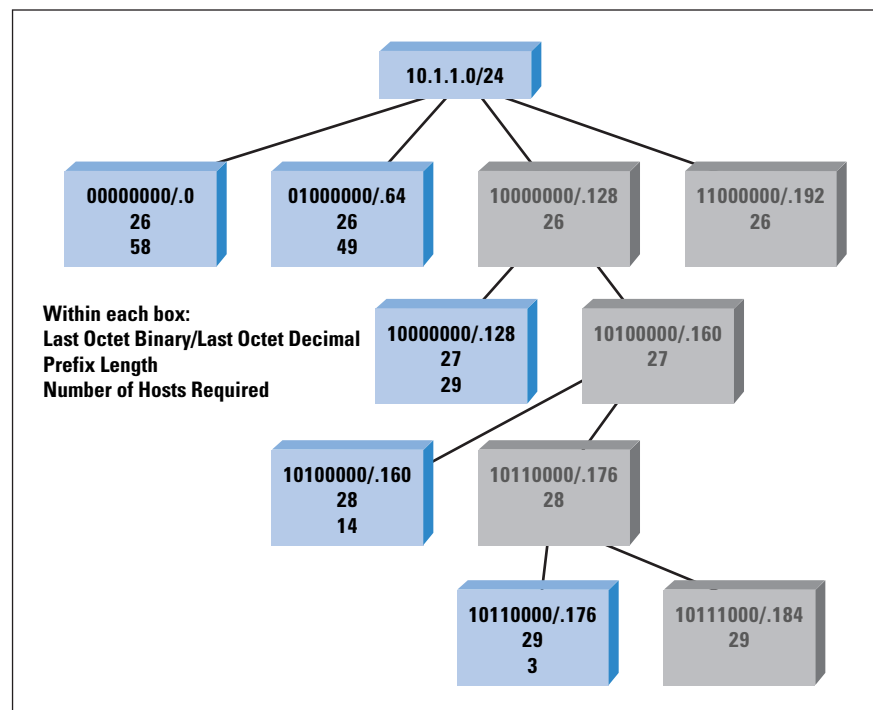
A few of our eagle-eyed readers have pointed us to some errors in IPJ, Volume 9, Number 1. The text below Figure 6 on page 29 and continuing at the top of page 30 should read as follows:

The figure shows four hosts with the addresses 10.1.0.1, 10.1.0.2, 10.1.0.3, and 10.1.0.4. Router A advertises 10.1.0.0/24, meaning: “Any host within the address range 10.1.0.0 through 10.1.0.255 is reachable through me.” Note that not all the hosts within this range exist, and that is okay—if a host within that range of addresses is reachable, it is reachable through Router A. In IP, the address that A is advertising is called a *network address*, and you can conveniently think of it as an address for the wire to which the hosts and router are attached, rather than a specific device.

For many people, the confusing part comes next. Router B is advertising 10.1.1.0/24, which is another network address. Router C can combine—or *aggregate*—these two advertisements into a single advertisement. Although we have just removed the correspondence between the wire and the network address, we have not changed the fundamental meaning of the advertisement itself. In other words, Router C is saying: “Any host within the range of addresses from 10.1.0.0 through 10.1.1.255 is reachable through me.” There is no wire with this address space, but devices beyond Router C do not know this, so it does not matter.

Also, Figure 8 on page 32 is reproduced here in its corrected form:

Figure 8: Subnet Chart



Book Review

Wireless Networking

Wireless Networking in the Developing World: A practical guide to planning and building low-cost telecommunications infrastructure, by Rob Flickenger et al., ISBN 1-4116-7837-0, 234 pages, Limehouse Book Sprint Team, January 2006. <http://wndw.net>

To quote from the book's Website:

"This book was created by a team of individuals who each, in their own field, are actively participating in the ever-expanding Internet by pushing its reach farther than ever before. Over a period of a few months, we have produced a complete book that documents our efforts to build wireless networks in the developing world."

Even though I don't live and work in what is commonly regarded as part of the developing world, I found this to be a unique and informative book, as its practical descriptions of wireless networking have application in many environments.

Given the widespread availability of the raw materials of computers, open-source software, Wi-Fi equipment, various pieces of recycled kitchenware, scrap metal, and plastic, and a wealth of online information resources, it is possible to construct inexpensive high-speed wireless network systems almost anywhere these days. However, perhaps the most visible missing component of the overall picture, but also the most valuable, is a practical path through this wealth of information on how to construct wireless networks, and a path that is based on the recent experiences of others who have constructed cost-effective and practical wireless networks in communities in the developing world. This book sets out to meet that goal.

Organization

The book starts with a description of radio physics covering the basics of the topic. It builds upon this a description of the typical radio design trade-offs between information capacity and radio penetration, and describes the commonly encountered factors of absorption, reflection, diffraction, and interference. I found the practical approach to Fresnel zone calculation and the description of the relationship between distance and antenna height so well done that I was tempted to embark on the design of a neighborhood Wi-Fi straightaway!

The chapter on network design is somewhat of a hybrid section, covering a mix of physical layout of a wireless network and TCP/IP considerations. There were the usual summaries of IP address structure and an introduction to routing.

Study of the deployment of the *Optimized Link State Routing* (OLSR) protocol is, however, more detailed. This is a link state routing protocol that is open-source, supportable by Linux-based access points, and accommodates link quality metrics into the routing protocol metric. I found the consideration of the link budget in this section a useful practical description of the considerations that are unique to the wireless world, and the worked examples are excellent, together with some useful references to online tools. This chapter is relatively dense, and many topics are covered in a relatively short space. I suspect that an interested reader would want to drill down further before feeling confident enough to manage a service network, but some carefully chosen references to further reading are there, so that the reader can follow up this introductory material with more specialized references.

The section on antennas and transmission lines was also well-structured. I had heard of using cylindrical cans as Wi-Fi antennas, but knew little of the detail of how to actually do it. This book not only explains their design, but provides a step-by-step illustrated guide to their construction. It also provides a good description of what is involved in outdoor installation of wireless equipment. The consideration of commercial solutions as compared to the do-it-yourself approach was carefully presented, as was the section devoted to security considerations.

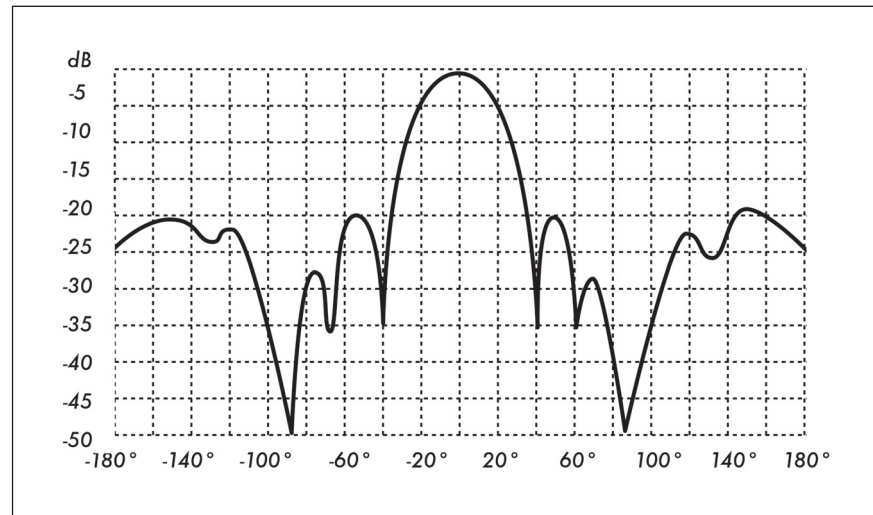
Aside from the technical considerations, the book also has some very interesting case studies of wireless networking projects, and was careful to include both success and failure stories. The issues in the developing world about combining technical capability with practical business solutions for communities that can be financially self-sustaining are indeed challenging, as the case studies show. They provide not only useful information about related experiences in setting up such network services, but also show how such projects can be assessed in a constructive manner.

Thoughtfully Written

Having spent some time working in this area myself as part of the ISOC *Developing Countries Workshop* training team, I have developed an appreciation of what constitutes truly useful and valuable training material, and this book is perhaps the best example I've seen yet. It is practical, helpful, technically accurate, and relatively complete in terms of coverage of material. Where the book does not dive into fine detail it provides useful references for further reading. The book is thoughtfully written in a simple non-nonsense style and does not hide behind technical jargon. Above all, it is material that can instill confidence that these networks can readily be built and operated by people like you and me.

I certainly would not call myself an expert after reading this book, but the next time a radio technician arrives in the office and starts talking about radiation patterns, front-to-back ratios, and the relative merits of omnis and yagis, at least I'll have an idea of what he is talking about. Even better, I might even be able to show him my own modest efforts in do-it-yourself Wi-Fi networking by then!

Rectangular plot of a Yagi Radiation Pattern from Chapter 4 of the book



Publishing Model

This is not a conventional technical book in the sense that it does not come with a conventional technical book price tag. The book is published in a manner as to be readily available in the developing world, so an online publication model has been used here. The PDF is freely available under a *Creative Commons* Attribution-Share-Alike 2.5 license at <http://wndw.net>, and they have managed to squeeze all 254 pages into an impressively small 1.92-MB file. You can find related resources and ways that you can assist in this project at <http://wndw.net>.

—Geoff Huston, APNIC
gih@apnic.net

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. Contact us at ipj@cisco.com for more information.

Fragments

Internet Governance

The *World Summit on the Information Society* (WSIS) was held in two phases. The first phase took place in Geneva in December 2003, and the second phase took place in Tunis in November 2005. The so-called “WSIS Outcome Documents” are now available at:

<http://www.itu.int/wsisis/promotional/outcome.pdf>

The follow-on to WSIS is called the *Internet Governance Forum* (IGF). The forum will hold its first meeting in Athens, Greece October 30th to November 2nd, 2006. For more information visit:

<http://www.intgovforum.org/>

The *Internet Society* (ISOC) played an active part in the WSIS process. You will find background information here:

<http://www.isoc.org/isoc/conferences/wsisis/index.shtml>

DNS Root Name Servers Explained

Daniel Karrenberg of RIPE NCC has written two “Member Briefings” on the subject of DNS root servers that can be found on the ISOC Website:

<http://www.isoc.org/briefings/019/>

<http://www.isoc.org/briefings/020/>

Internationalized Domain Names

Internationalized Domain Names (IDNs) are, according to the ICANN Website, “...domain names represented by local language characters. Such domain names could contain letters or characters from non-ASCII scripts (for example, Arabic or Chinese). Many efforts are ongoing in the Internet community to make domain names available in character sets other than ASCII.” ICANN has established an information area on its Website with links to more information about IDNs. See:

<http://icann.org/topics/idn/>

The ISP Column

Geoff Huston is well known to readers of this journal. He also hosts *The ISP Column* that can be found here:

<http://www.isoc.org/pubs/isp/index.shtml>

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 2006 Cisco Systems Inc.
All rights reserved.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA
--

The Internet Protocol Journal

September 2006

Volume 9, Number 3

A Quarterly Technical Publication for
Internet and Intranet Professionals

In This Issue

From the Editor	1
Wireless LAN Switches.....	2
IPv6 Internals.....	16
Book Reviews	30
Fragments	34
Call for Papers.....	35

FROM THE EDITOR

One of the most successful networking technologies of recent years has been IEEE 802.11 or, as it is commonly known, “Wi-Fi.” Wireless networks have seen widespread deployment within organizations as well as in public “hotspots” all over the world. As a frequent traveler, I am very pleased with this development. It has been a long time since I had to resort to a modem and phone line in order to access e-mail or use the Web. Wireless networks have truly changed the way we use the Internet. Our first article, by T. Sridhar, explores the emerging use of *Wireless LAN Switches* in wireless access networks.

IPv6 is a technology that perhaps should have been widely deployed by now, but wide deployment has not happened yet, for numerous reasons. This journal has covered many aspects of IPv6. This time, Iljitsch van Beijnum looks at some of the details you need to be aware of when considering a move to IPv6. The article is adapted from his book *Running IPv6*, which was reviewed in our December 2005 issue.

In previous editions of IPJ we have pointed you to other sources of information, such as *The IETF Journal*, Geoff Huston’s *ISP Column*, and other documents available from the Internet Society Website at <http://www.isoc.org>. This time I want to make you aware of an article that originally appeared in *Apster*, the newsletter of the *Asia Pacific Network Information Centre* (APNIC), one of the five *Regional Internet Registries* (RIRs). The article is entitled “IP Addressing in China and the Myth of Address Shortage,” and you will find the URL for it in our “Fragments” section. If you want to further explore the work of the RIRs, you can start by visiting the *Number Resource Organization* (NRO) at <http://nro.net>.

You may have read that both of our sister publications, *Packet* and *IQ Magazine*, are publishing their final issues this September. Naturally, this has led to some of our readers asking what is in store for IPJ. We want to reassure you that we intend to continue publishing IPJ in both its paper and online forms. Plans are also under way to enhance our Website to provide you with more tools and resources. If you have suggestions for the Website, please send us a note at ipj@cisco.com.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Wireless LAN Switches — Functions and Deployment

by T. Sridhar, Flextronics

Deployment of *Wireless LAN* (WLAN) switches is increasing in enterprise networks. These devices, which can be standalone switches or integrated into a blade on an enterprise class switch, are useful for the management and control of WLAN access points. Although their deployment is a relatively new phenomenon, such control and configuration functions have existed before in WLAN controller devices.

WLAN switches connect to the WLAN *access points* (APs) through wired connections (through a switch port). They also connect to the enterprise network through their other switch ports. The switches are the “gateway” to the wired enterprise—all frames from WLAN clients have to pass through the WLAN switches to the enterprise network.

To understand the motivation for WLAN switches and their operation in the network, it is useful to view the WLAN network architecture and the functions of the access points. We can view the WLAN switch as the control function and the APs as the wireless termination function.

This article presents the function of WLAN switches and controllers by detailing WLAN network architectures along with functions of the AP and controller. It also presents the various functions on the controller to AP interface. Subsequently, it outlines variables related to Layer 2/3 mobility in the centralized architecture and concludes by presenting some common myths and reality about these architectures.

This article uses the term *Wireless Termination Point* (WTP) to refer generically to APs and the term *Access Controller* (AC) to refer generically to the WLAN control function (whether implemented on a WLAN switch or standalone controller).

WLAN Network Architectures

Three types of WLAN network architectures are commonly deployed:

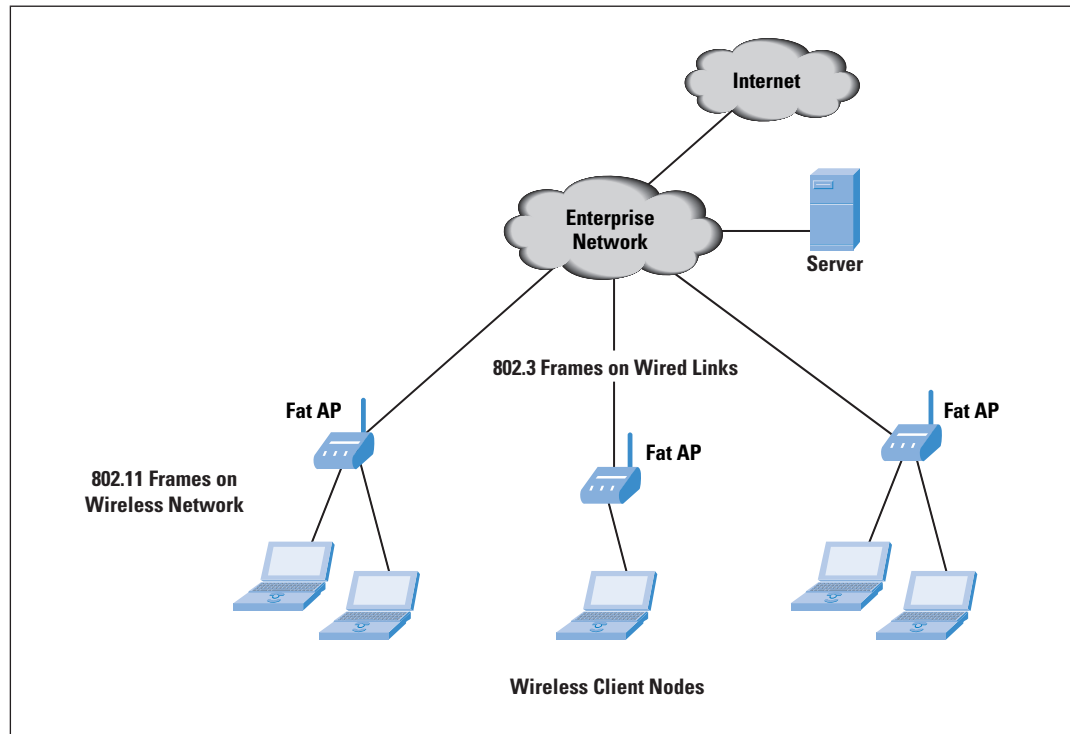
1. Autonomous Architecture
2. Centralized Architecture
3. Distributed Architecture

The following sections describe these architectures in greater detail.

Autonomous Architecture

In the autonomous architecture, the WTPs completely implement and terminate the 802.11 function so that frames on the wired LAN are 802.3 frames. Each WTP can be independently managed as a separate network entity on the network. The access point in such a network is often called a “Fat AP” (see Figure 1).

Figure 1: FAT APs in Autonomous WLAN Network Architecture

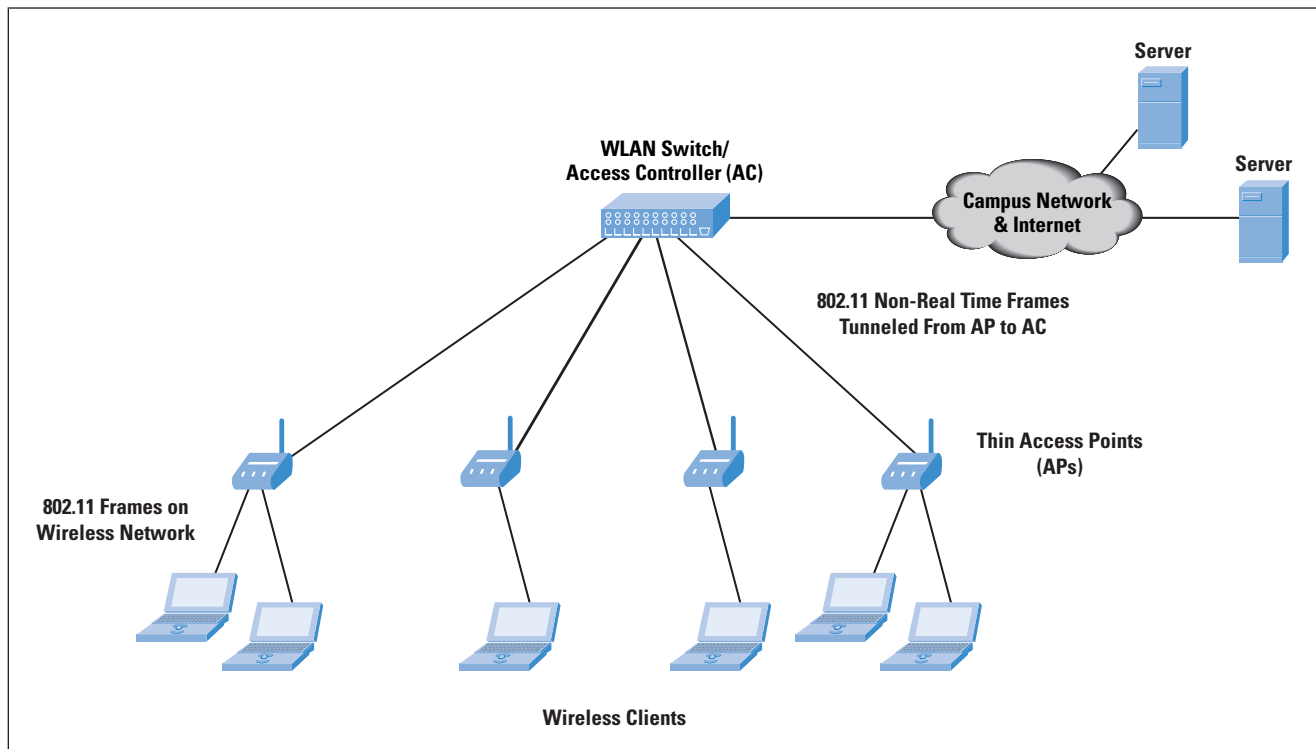


During the initial stages of WLAN deployment, most APs were autonomous APs, and manageable as independent entities in the network. During the past few years, centralized architectures (discussed next) with ACs and WTPs have gained popularity. The primary advantage of the centralized architecture is that it provides network administrators with a structured and hierarchical mode of control for multiple WTPs in the enterprise.

Centralized Architecture

The centralized architecture is a hierarchical architecture that involves a WLAN controller that is responsible for configuration, control, and management of several WTPs. The WLAN controller is also known as the *Access Controller* (AC). The 802.11 function is split between the WTP and the AC. Because the WTPs in this model have a reduced function as compared to the autonomous architecture, they are also known as “Thin APs.” Some of the functions on the APs are variable, as discussed in the following section (see Figure 2).

Figure 2: Thin APs in Centralized WLAN Network Architecture



Distributed Architecture

In the distributed architecture, the various WTPs can form distributed networks with other WTPs through wired or wireless connections. A mesh network of WTPs is one example of such an architecture. The WTPs in the mesh can be linked with 802.11 links or wired 802.3 links. This architecture is often used in municipal networks and other deployments where an “outdoor” component is involved. This article does not address the distributed architecture.

WTP Functions – Fat, Thin, and Fit APs

To understand the autonomous and centralized architecture, it is useful to look at the functions performed by the APs. We start with the Fat APs, which form the core of the autonomous architecture, followed by the Thin APs, which were specified as part of the WLAN switch- or controller-based centralized architecture. The article will then outline the functions of a new variant called the “Fit AP,” an optimized version of the AP for centralized architectures.

Fat Access Points

Figure 1 shows an example of an autonomous network with a fat access point. The AP is an addressable node in the network with its own IP address on its interfaces. It can forward traffic between the wired and wireless interfaces. It can also have more than one wired interface and can forward traffic between the wired interfaces—similar to a Layer 2 or Layer 3 switch. Connectivity to the wired enterprise can be through a Layer 2 or Layer 3 network.

It is important to understand that there is no “backhauling” of traffic from the Fat AP to another device through tunnels. This aspect is important and is addressed when discussing the other AP types. In addition, Fat APs can provide “router-like” functions such as the *Dynamic Host Configuration Protocol* (DHCP) server capabilities.

Management of the AP is done through a protocol such as the *Simple Network Management Protocol* (SNMP) or the *Hypertext Transfer Protocol* (HTTP) for Web-based management and a *Command-Line Interface* (CLI). To manage multiple APs, the network manager has to connect to each AP through one of these management schemes. Each AP shows up on the network map as a separate node. Any aggregation of the nodes for management and control has to be done at the *Network Management System* (NMS) level, which involves development of an NMS application.

Fat APs also have enhanced capabilities such as *Access Control Lists* (ACLs), which permit filtering of traffic for specific WLAN clients. Another significant capability of these devices is configuration and enforcement of *Quality of Service* (QoS)-related functions. For example, traffic from specific mobile stations might need to have a higher priority than others. Or, you might need to insert and enforce IEEE 802.1p priority or *Differentiated Services Code Point* (DSCP) for traffic from mobile stations. In summary, these APs act like a switch or router in that they provide many of the functions of such devices.

The downside of such APs is complexity. Fat APs tend to be built on powerful hardware and require complex software. These devices are expensive to install and maintain because of the complexity. Nevertheless, the devices have uses in smaller network installations.

Some Fat AP installations still use a controller at the back end for control and management functions. These controllers lead to a slightly scaled-down version of the Fat AP, called, not surprisingly, a Fit AP, discussed later.

Thin Access Points

As their name indicates, Thin APs are intended to reduce the complexity of APs. An important motivation for this reduction is the location of APs. In several enterprises, APs are plenum-mounted (and thus in hard-to-reach areas) so that they can provide optimum radio connectivity for end stations. In environments like warehouses, this is even more evident. For such reasons, network managers prefer to install APs just once and not have to perform complex maintenance on them.

Thin APs are often known as “intelligent antennas,” in that their primary function is to receive and transmit wireless traffic. They backhaul the wireless frames to a controller where the frames are processed before being switched to the wired LAN (see Figure 2).

The APs use a (typically secure) tunnel to backhaul the wireless traffic to the controller. In their most basic form, Thin APs do not even perform WLAN encryption such as *Wired Equivalence Privacy* (WEP) or *WiFi Protected Access* (WPA/WPA2). This encryption is done at the controller—the APs just transmit or receive the encrypted wireless frames, thereby keeping the APs simple and avoiding the necessity to upgrade their hardware or software.

The introduction of WPA2 necessitated encryption on the controller. Although WPA was hardware-compatible with WEP and required only a firmware upgrade, WPA2 was not backward-compatible. Instead of replacing APs across the enterprise, network managers could just backhaul the wireless traffic to the controller where the WPA2 decryption was done, and the frames were sent on the wired LAN.

The protocol between the AP and the controller for carrying the control and data traffic was proprietary. Also, there is no capability to manage the AP as a single entity on the Layer 2/3 network—it can be managed only through the controller, to which the NMS can communicate through HTTP, SNMP, or CLI/Telnet. A controller can manage and control multiple APs, implying that the controller should be based on powerful hardware and often be able to perform switching and routing functions. Another important requirement is that the connectivity and tunnel between the AP and the AC should ensure low delay for packets between those two entities.

With Thin APs, QoS enforcement and ACL-based filtering are handled at the controller—not a problem because all the frames from the AP have to pass through the controller anyway. Centralized control functions for ACLs and QoS are not new—they were implemented in networks with Fat APs too. Such installations have controllers that act as the gateway for managing traffic from APs to the wired network. However, the controller function takes on a new dimension with Thin APs, especially with respect to the data plane and forwarding functions. The controller function subsequently was integrated into Ethernet switches that connected the wireless and wired LANs—the motivation for the family of devices known as WLAN switches.

The Wireless MAC architecture in this scenario is known as the *Remote MAC* architecture. The entire set of 802.11 MAC functions is offloaded to the WLAN controller, including the delay-sensitive MAC functions.

Fit Access Points

Fit APs are gaining in popularity in that they try to take advantage of the best of both worlds—that is, the Fat APs and the Thin APs. A Fit AP provides the wireless encryption while using the AC for the actual key exchange. This approach is used for newer APs that use the latest wireless chipsets supporting WPA2. The management and policy functions reside on the controller that connects to multiple APs through tunnels.

Also, Fit APs provide additional functions such as DHCP relay for the station to obtain an IP address through DHCP. In addition, Fit APs can perform functions such as VLAN tagging based on the *Service Set Identifier* (SSID) that the client uses to associate with the AP (when the AP supports multiple SSIDs).

Two types of MAC implementations are possible with Fit APs, known as the *Local MAC* and the *Split MAC* architectures. Local MAC is where all the wireless MAC functions are performed at the AP. The complete 802.11 MAC functions, including management and control frame processing, are resident on the APs. These functions include time-sensitive functions (also known as *Real Time MAC* functions).

The Split MAC architecture divides the implementation of the MAC functions between the AP and the controller. The real-time MAC functions include functions such as beacon generation, probe transmission and response, control frame processing (for example *Request to Send* and *Clear to Send*—RTS and CTS), retransmission, and so on. The non-real time functions include authentication and deauthentication; association and reassociation; bridging between Ethernet and Wireless LAN; fragmentation; and so on.

Vendors differ in the type of functions that are split between the AP and the controller, and in some cases, even about what constitutes real time. One common implementation of a Fit AP involves local MAC at the AP and control and management functions at the AP.

Access Controller and Control Functions

The next critical component of the Centralized WLAN Architecture is the *Access Controller* (AC). For the following discussion, we consider the controller function to be implemented on a WLAN switch and call the function an AC. We also use the term “WTP” to refer to APs (fat, thin, or fit).

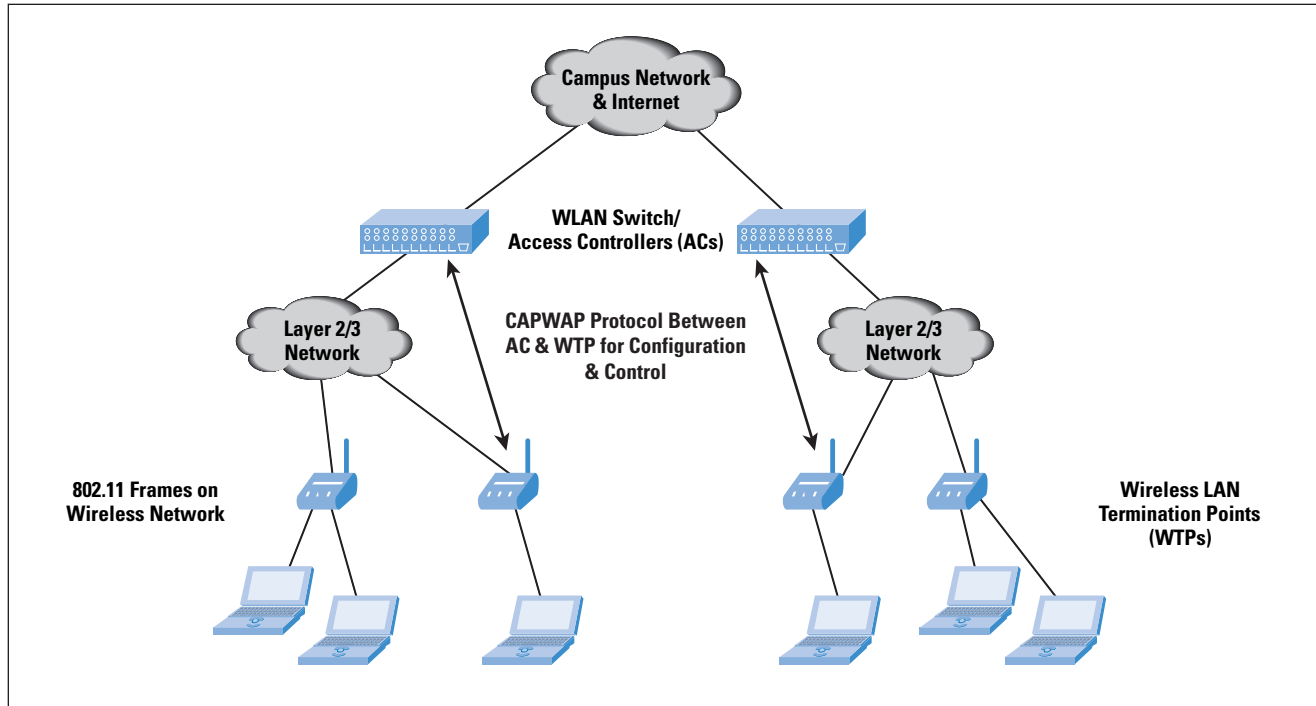
The *Control and Provisioning of Wireless Access Points* (CAPWAP) Working Group in the IETF is working on defining the interface and protocol between an AP and its controlled WTP. This section uses the CAPWAP framework to detail the interface between the AC and the WTP. [3,4,5]

Figure 3 shows an enterprise network with multiple ACs and WTPs. The WTPs can be connected to the ACs through a Layer 2 (switched) or Layer 3 (routed) network. The interface between the WTP and the AC is responsible for the following:

- Discovery and selection of an AC by WTP
- Firmware download to the WTP by the AC—upon startup and upon triggering by the WTP
- Capabilities negotiation between the WTP and the AC

- Mutual authentication between the WTP and the AC
- Configuration, status, and statistics exchange between the WTP and the AC
- QoS mapping across the wired and wireless segments

Figure 3: Centralized WLAN Architecture with Multiple ACs, WTPs and CAPWAP Protocol Context



In addition, although CAPWAP does not explicitly define all the details, the AC performs functions such as *Radio Resource Management* (RRM) and rogue AP detection based on configuration and monitoring of the various access points in its domain of control. The extent of these functions varies according to the vendor implementation. Another important function provided by ACs is mobility management. The following sections provide more detail about these functions, with specific reference to CAPWAP. Note that the CAPWAP protocol, which is based on the Cisco *Lightweight Access Point Protocol* (LWAPP), is still under development in the IETF, as of the writing this article (March 2006).

Discovery and AC Selection

A WTP discovers an AC to connect to through discovery request messages, to which one or more ACs can respond (depending on the network topology). Communication between the AC and the WTP is through the *User Datagram Protocol* (UDP). The WTP determines which AC to connect to and then tries to establish a secure session with the AC. Subsequent CAPWAP packets are sent over the secure session.

Subsequently a configuration exchange takes place between the AC and WTP. This exchange includes:

- IEEE SSID
- Security parameters (for WEP, WPA, and WPA2)
- Data rate that is to be advertised (11 or 54 Mbps)
- Radio channels to be used

CAPWAP Functions

CAPWAP control messages include the following message types:

- Discovery
- WTP configuration—used to push a specific configuration to the WTP and also to retrieve statistics from a WTP; statistics includes information such as:
 - Number of fragmented frames, multicast frames transmitted and received
 - Number of transmit retries, excessive retries (failed count)
 - Number of successfully transmitted and failed *Requests to Sends* (RTS)
 - Number of errored frames: duplicate frames, failed acks, decryption errors, *frame-check-sequence* (FCS) error count, etc.

Configuration includes information such as beacon period, maximum transmit power level, *Orthogonal Frequency Division Multiplexing* (OFDM) control, antenna control, supported rates, QoS, encryption, and so on.

- *Mobile session management*—to push specific mobile policies to the WTP

ACs can add policy information about specific mobile devices that can include security parameters that the WTP should apply for that mobile device. It can indicate whether the WTP should forward or discard traffic for that mobile device.

- *Firmware management*—used to push a specific firmware image to the WTP

AC and WTP Interaction

The WTP provides information such as hardware, software, or boot version; maximum number of radios; radios in use; encryption capabilities; type of radio (802.11b/g/a/n); type of MAC (local, split, or both); tunneling modes; and frame type between AC and WTP (for example, local bridging or native bridging—that is, encapsulating all user payloads as native wireless frames).

The AC information includes hardware or software version, number of mobile stations currently associated with the AC, number of WTPs currently attached to the AC, maximum numbers for each of these, security parameters (authentication credentials) between AC and WTP, control IPv4 or IPv6 address, and so on.

Because the WTPs fall under the category of “Fit APs,” they can also be configured with an IP address from the AC. Another parameter that can be configured is ACLs at the MAC address level.

Rebooting (reset) of the WTP can be done by the AC at any time. Independently, the WTP can request a new image through an *Image Data Request*, which is followed by an *Image Data Response* and the image data itself.

Events are sent by the WTP when it determines that it has important information to send to the AC. Such information can include data transfer messages that can be used to deliver debug information from the WTP to the AC.

Radio Resource Management

Radio resource management is a generic term used to describe the control and configuration of radios on the AP. The type of control includes reducing and increasing the strength automatically or on user input—for example, if two WTPs controlled by an AC are interfering with each other, the AC can send a signal to one of the APs to reduce its strength. It can also do this based on user configuration.

Several WTPs are designed to also be used as “Air Monitors;” that is, they can monitor channels when not transmitting. Opinion is still divided on whether this mode of using WTPs is efficient—some vendors use dedicated air monitors instead of having their WTPs do double duty. With dedicated air monitors, it is much easier to scan and monitor all channels without having to worry about degrading the service for client stations.

Air monitors can forward information about other access points to the AC. The AC can determine if the information is for a valid WTP (that is, one that is supposed to be on the network and has, in fact, registered with the AC) or for a “rogue” access point. If it is for a rogue access point, the AC can perform multiple steps to prevent clients from attaching to this AP—for example, it can instruct the air monitor to “jam” this rogue AP by increasing the transmit power on the same channel.

Mobility Management

Mobility management can take two forms—Layer 2 and Layer 3 mobility. Consider a client moving from one WTP to another, a scenario that can happen when a user with a laptop moves between two conference rooms within the same building. The client station reassociates with the new WTP, after which authentication is performed. Note that the association with the previous AP is “broken” before the association with the new AP is “made;” thus handoff in WLANs is known as “break before make.” Although this approach can lead to potential traffic disruption (and retransmissions), it is chosen over “make before break” (used in cellular telecommunications) to keep the client radio simple and less expensive.

One way to envision Layer 2 and Layer 3 mobility is to treat Layer 2 mobility as movement between APs under the control of the same AC (that is, Layer 3 network), whereas Layer 3 mobility is movement between APs under the control of different ACs.

Layer 2 Mobility

Layer 2 mobility means that when the station moves from one WTP to another, there is no impact on the IP addressability, effectively meaning that all the APs are on the same Layer 2 network and implying that they are connected to the same AC (see Figure 4). To prevent loss of data destined to the Layer 2 client, the WLAN switch must now forward client data to the new WTP. After the client association, the new WTP sends out an Ethernet frame to the AC with the client's MAC address as the source address. The switch now associates the client's MAC address to the port on which the new WTP is connected.

Although this process works well with Layer 2 (switched network) connectivity between the APs and the AC, it requires a slightly different approach when tunnels are used between them. The AC moves the mapping of the client to a different tunnel (that is, a virtual port) when it receives the MAC frame from the new WTP.

Another concern to be addressed with Layer 2 handoff is the buffering of data at the WTP. In normal circumstances, the switch or AC is not aware of the handoff until it hears from the new WTP. However, with enhanced statistics available at the WTP, it can determine that the specific client has moved away from the old WTP and stop forwarding data to the old WTP. These statistics can include maximum retry attempts on the *Carrier Sense Multiple Access/Collision Avoidance* (CSMA/CA) MAC layer protocol on the wireless link. The switch does not need to buffer the data because it is not clear when the handoff to the new WTP will occur. This approach helps avoid wasteful traffic on the link between the old WTP and the AC.

Some vendors have approached this problem differently with Fat APs. There, the APs might buffer the traffic until they see a frame from the switch indicating that the client is now on a different switch port. These APs then send the buffered traffic to the switch, which forwards that to the new WTP. Because our intent is to lower the complexity of the WTPs, this approach is not a preferred one in the Centralized AC + WTP architecture.

Another important feature of Layer 2 roaming is preauthentication that needs to be done on the new WTP. Through 802.11i, clients can preauthenticate with neighboring WTPs so that roaming to a different WTP does not involve the lengthy authentication process of *Pairwise Master Keys* (PMKs) being sent to the new WTP. (The *Pairwise Transient Keys* (PTKs) still need to be derived.)

When the AC maintains the PMK for a specific client (through interaction with a RADIUS server), this process is automatic—that is, the AC can send the client-specific PMK to the new WTP. The encryption of 802.11 frames is still done by the old and new WTPs with the new PTKs.

Layer 3 Mobility

Layer 3 mobility involves the client retaining the same IP address while moving across multiple APs. This often happens when the client has published its IP address to multiple nodes. Such a scenario is likely in peer-to-peer communications and when the mobile station needs to act as a server for some function. It is desirable that the correspondent nodes communicating with the mobile node not have to change their configuration whenever the mobile node moves to a new Layer 3 network.

This problem of Layer 3 mobility is solved by *Mobile IP*^[6]. We do not discuss the details of Mobile IP here except to indicate that it has three distinct components. The *Home Agent* (HA) on the client's home network is responsible for the address of the client. All packets destined to the client's (invariant) IP address are sent to the Home Agent. If the client is on the home network, the HA forwards the packets directly to the client. If it is on a foreign or visited network, the HA forwards the packets to a *Foreign Agent* (FA) that is on the visited network.

To do this, it has to set up a tunnel to the FA—which is usually a *Generic Routing Encapsulation* (GRE) or IP-in-IP tunnel.

After stripping out the original packet from the tunnel, the FA is responsible for forwarding the packet to the client. This description is a simplification—numerous other steps are involved here. The important factor in a wireless LAN scenario for Layer 3 client mobility is where the Mobile IP endpoint resides. Some client stations include a software stack for a MIP client.

This *Client MIP* (CMIP) software:

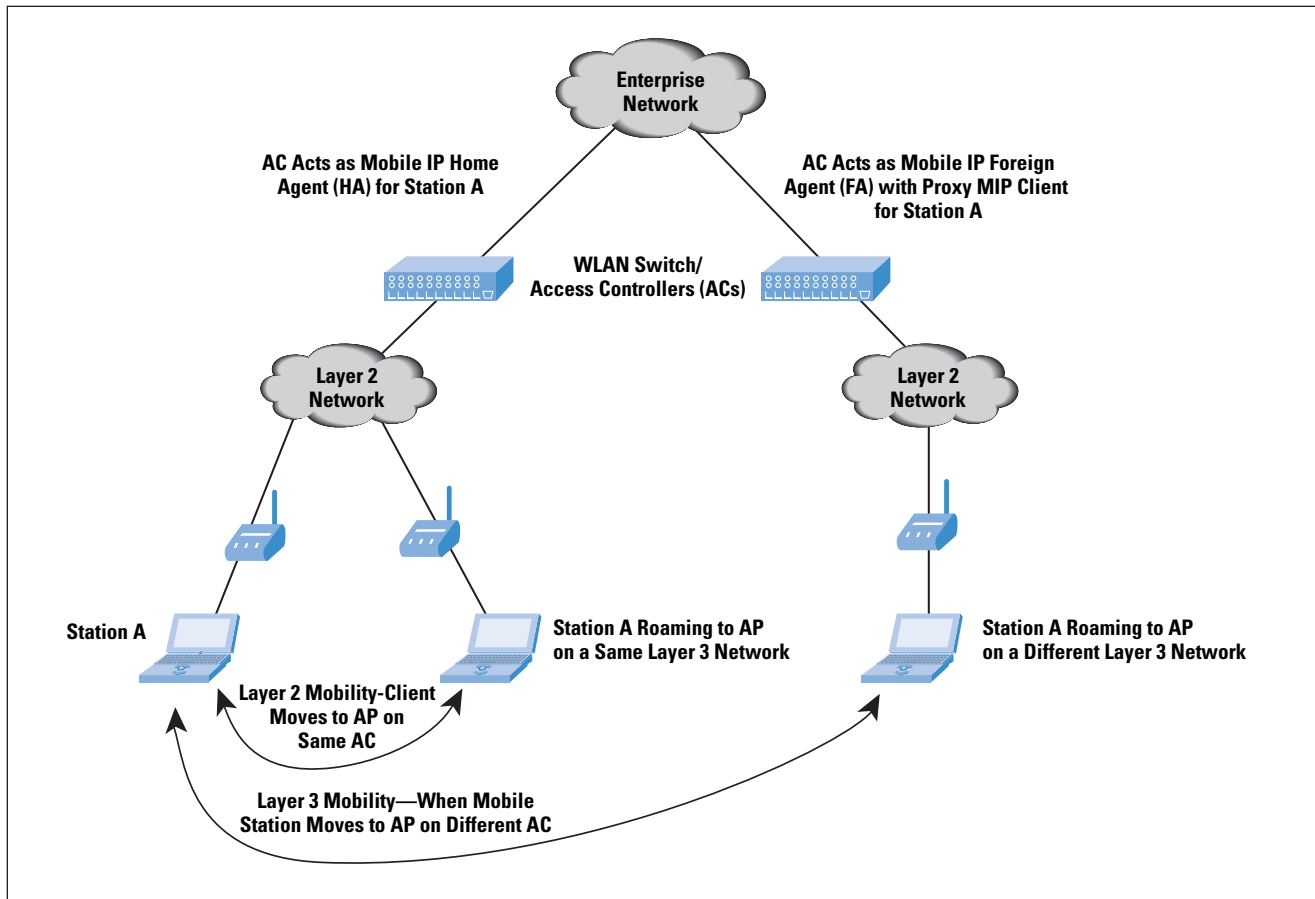
- Strips out the MIP header in the packet
- Inserts a new header to spoof the client's higher-layer applications into believing that the packets were destined for the client's IP address on the foreign network

The CMIP approach was the recommended approach for implementing MIP. However, it has the disadvantage of having to add a MIP client to every mobile station in the network—a setup that can become cumbersome when there are a large number of mobile stations.

The Centralized AC + WTP architecture offers a way of alleviating this problem. Some AC/WLAN switch vendors have implemented the MIP function on the AC so that the client never needs to be changed. Some implementations call this a *Proxy MIP* function.

The AC acts as an FA to terminate the tunnel from the HA and also performs the translation of the packets to the client's address on the visited network when forwarding packets to the client. When the client sends Layer 3 packets out, it sends them through the AC, which, in turn, modifies the headers for the source IP address and tunnels the packets to the HA. This process is called “reverse tunneling” (see Figure 4).

Figure 4: Layer 2 and Layer 3 Mobility in Centralized WLAN Network Architecture



When you consider a large enterprise network topology with multiple ACs and APs, you can envision the MIP tunnels to be established between the various ACs. (That is, they act as Foreign Agents for one set of users and as Home Agents for another set.) From a scalability perspective, it is important that the ACs have the necessary horsepower and switching capability (switching between tunnels from the APs to the ACs to the tunnels between the ACs).

WLAN Switches and Centralized Architectures – Common Myths

Previous sections considered various aspects of the Centralized AC + WTP architecture and some of the implementation factors. This section outlines some common myths about these architectures and implementations. The intent is to examine this still-evolving area to facilitate clarity.

1. *Myth 1: ACs need to perform switching functions—hence the name WLAN switches.*

There is no such requirement. In fact, the earliest ACs were appliances (and in some cases, PCs running Linux). The control function is the important part of the implementation—the switching is often included to accelerate the forwarding of traffic to and from the APs.

2. *Myth 2: Rogue WTP detection is a standard function of ACs.*

This is a desired function in several implementations but is not necessarily “standard.” One reason is that this is an area of differentiation among vendors (for example, the algorithms they use to classify a WTP as a rogue WTP). Another reason is that the ACs have to rely on APs or air monitors, and this reliance varies according to implementation.

3. *Myth 3: The delineation between Fat, Thin, and Fit APs is clearly defined.*

There are several types of implementations of AP (and AC) functions, so this myth is not necessarily true. For a sample of the taxonomy (snapshot) of WTP and AC implementations, see RFC 4118^[4].

4. *Myth 4: Layer 2 and Layer 3 mobility are standard in AC+ WTP architectures.*

This is not really true. The Proxy MIP implementation for Layer 3 mobility is a step in this direction, but most AC vendors rely on proprietary mechanisms for AC-AC communication and Layer 3 mobility.

5. *Myth 5: Security functions such as firewall, intrusion detection, and so on are not a function of ACs.*

Some vendors have debunked this argument and implemented such functions in their AC. This is an area for vendor differentiation.

Summary

This article has provided the functions and deployment of WLAN switches by detailing the architectures that rely on a centralized controller managing a set of wireless termination points. It outlined some major aspects of the CAPWAP control functions and the concerns related to Layer 2 and Layer 3 mobility while implementing an AC + WTP architecture. Although protocol standardization is being done in the IETF for this emerging area, there is still sufficient scope for vendor differentiation.

References

- [1] “IEEE 802.11i and Wireless Security,” David Halasz, www.embedded.com, August 25, 2004.
- [2] Rich Seifert, *The Switch Book: The Complete Guide to LAN Switching Technology*, ISBN 0471345865, Wiley, 2000.
- [3] B. O’Hara, et al., “Configuration and Provisioning for Wireless Access Points (CAPWAP): Problem Statement,” RFC 3990, February 2005.
- [4] L. Yang, et al., “Architecture Taxonomy for Control and Provisioning of Wireless Access Points (CAPWAP),” RFC 4118, June 2005.
- [5] P. Calhoun, Editor, “CAPWAP Protocol Specification,” (work in progress), Internet Draft, **draft-ietf-capwap-protocol-specification-00**, February 24, 2006.
- [6] C. Perkins, Editor, “IP Mobility Support for IPv4,” RFC 3344, August 2002.
- [7] Edgar Danielyan, “IEEE 802.11,” *The Internet Protocol Journal*, Volume 5, No. 1, March 2005.
- [8] Gregory R. Scholz, “Securing Wireless Networks,” *The Internet Protocol Journal*, Volume 5, No. 3, September 2002.

T. SRIDHAR is Vice President of Technology at Flextronics in San Jose, California. He received his BE in Electronics and Communications Engineering from the College of Engineering, Guindy, Anna University, Madras, India, and his Master of Science in Electrical and Computer Engineering from the University of Texas at Austin. He can be reached at T.Sridhar@flextronics.com

IPv6 Internals

by Iljitsch van Beijnum

This article discusses some of the protocol details you should be aware of when planning a transition from IPv4 to IPv6. Although it is not intended as a complete step-by-step guide, this article explains the differences between IPv4 and IPv6 as they relate to actually operating a network. Vendor- and operating system-specific details can be found in the book from which this text was adapted, and further information is available in the references.

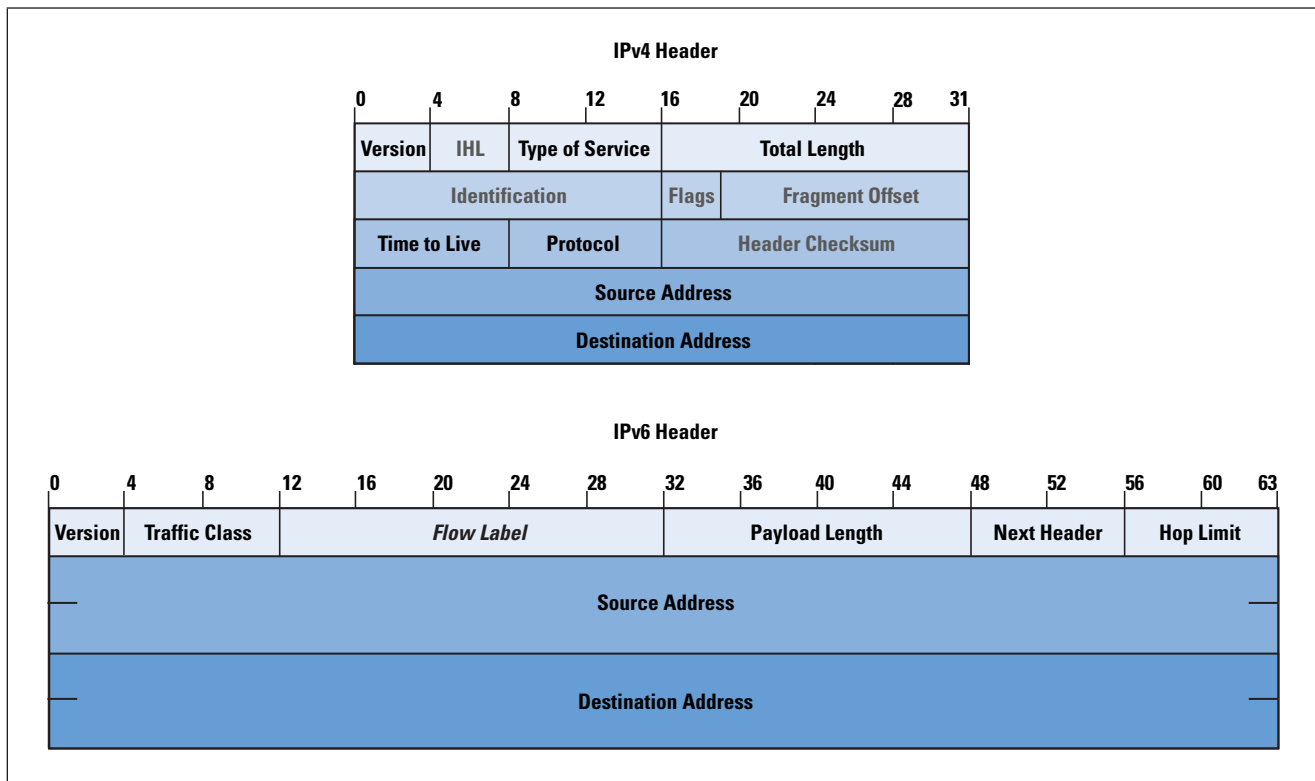
The easiest way to observe—in action—the mechanisms discussed in this article is to set up an IPv6 router on the local subnet and enable IPv6 on the operating system of your choice, if it is not enabled by default. If “native” IPv6 connectivity is not possible, you can set up automatic IPv6 tunneling or use a manually configured IPv6-in-IPv4 tunnel. Getting portable IPv6 address space from a *Regional Internet Registry* (RIR)^[1] is a topic worthy of its own article, but *6to4*^[2] creates 65,536 IPv6 subnets from a single IPv4 address, and service providers that provide IPv6 connectivity—either natively or over manually configured tunnels—are usually quite generous with IPv6 address space. However, you need to renumber when changing *Internet Service Providers* (ISPs), or when changing IPv4 addresses with 6to4. Most router vendors currently support IPv6 routing, but all widely used general-purpose operating systems can also route IPv6.

When you have IPv6 connectivity, the browser that comes with your system should be able to work over IPv6 (visit <http://www.kame.net/>), and there are v6 versions of *ping* and *traceroute* (called *ping6* and *traceroute6*) to determine IPv6 connectivity. More and more applications work over IPv6, but many still do not.

Differences Between IPv4 and IPv6

All knowledge about IPv6 begins with studying the IPv6 header format and the ways in which it is different from the IPv4 header format. Even though at the time the IPv6 specifications were written 64-bit CPUs were rare, the IPv6 designers elected to optimize the IPv6 header for 64-bit processing. For this reason, I have drawn the IPv6 header 64 bits wide in Figure 1, a little different from the way it is usually depicted. Because 64-bit CPUs can read one 64-bit-wide memory word at a time, it is helpful that fields that are 64 bits (or a multiple of 64 bits) wide start at an even 64-bit boundary. Because every 64-bit boundary is also a 32-bit boundary, 32-bit CPUs aren't affected negatively by 64-bit optimization. The IPv4 header is presented in the usual form that highlights its 32-bit background.

Figure 1: The IPv4 and IPv6 headers



The fields in the IPv4 header that are not present in the IPv6 header have gray text; the field that is present in IPv6 but not in IPv4 is shown in *italic*. The changes from IPv4 to IPv6 follow:

- *Version* now always contains 6 rather than 4.
- The *Internet Header Length* (IHL) field that indicates the length of the IPv4 header is no longer needed because the IPv6 header is always 40 bytes long.
- *Type of Service* is now *Traffic Class*. The original semantics of the IPv4 Type of Service field have been superseded by the *diffserv* semantics per RFC 2474^[3]. However, in IPv4, both interpretations of the field are in use (although most routers either cannot or are not configured to look at the field anyway). The IPv6 RFCs do not mandate a specific way to use the Traffic Class field, but generally the RFC 2474 *diffserv* interpretation is assumed.
- The *Flow Label* is new in IPv6. The idea is that packets belonging to the same stream, session, or flow share a common flow label value, making the session easily recognizable without having to look “deep” into the packet. Recognizing a stream or session is often useful in *Quality of Service* mechanisms. Although few implementations actually look at the flow label, most systems do set different flow labels for packets belonging to different TCP sessions. A zero value in this field means that setting a flow label per session is either not supported or not desired.

- The *Total Length* is the length of the IPv4 packet including the header, but in IPv6, the *Payload Length* does not include the 40-byte IPv6 header, thereby saving the host or router receiving a packet from having to check whether the packet is large enough to hold the IP header in the first place—making for a small efficiency gain. Despite the name, the Payload Length field includes the length of any additional headers, not just the length of the user data.
- The *Identification*, *Flags*, and *Fragment Offset* fields are used when IPv4 packets must be fragmented. Fragmentation in IPv6 works very differently (explained later), so these fields are relegated to a header of their own.
- *Time to Live* (TTL) is now called *Hop Limit*. This field is initialized with a suitable value at the origin of a packet and decremented by each router along the way. When the field reaches zero, the packet is destroyed. This way, packets cannot circle the network forever when there are loops. Per RFC 791^[4], the IPv4 TTL field should be decremented by the number of seconds that a packet is buffered in a router, but keeping track of how long packets are buffered is too difficult to implement, regardless of buffering time. The new name is a better description of what actually happens.
- The *Protocol* field in IPv4 is replaced by *Next Header* in IPv6. In both cases, the field indicates the type of header that follows the IPv4 or IPv6 header. In most cases, the value of this field would be 6 for TCP or 17 for the *User Datagram Protocol* (UDP). Because the IPv6 header has a fixed length, any options such as source routing or fragmentation must be implemented as additional headers that sit between the IPv6 header and the higher-layer protocol such as TCP, forming a “protocol chain.”
- The IPv4 *Header Checksum* was removed in IPv6.
- The *Source Address* and *Destination Address* serve the same function in IPv6 as in IPv4, except that they are now four times as long at 128 bits.

All IPv6 hosts and routers are required to support a maximum packet size of at least 1280 bytes. For lower-layer protocols that cannot support a *Maximum Transmission Unit* (MTU) of 1280 bytes, the relevant “IPv6 over ...” standard must have a mechanism to break up and reassemble IPv6 packets so that the minimum of 1280 bytes can be accommodated. In IPv4, the official minimum size is 68 bytes—too low to be workable.

Checksums

In IPv4, the IP header is protected by a header checksum, and higher-layer protocols generally also have a checksum. The checksum algorithm for the IPv4 header, *Internet Control Message Protocol* (ICMP), ICMPv6, TCP, and UDP is the same one’s complement addition, except that in IPv4, UDP packets may forego checksumming and simply set the checksum field to zero. In IPv6, this practice is no longer allowed: UDP packets must have a valid checksum.

The TCP, UDP, and ICMPv6 checksums are computed over a “pseudoheader” and the TCP, UDP, or ICMPv6 header, and user data, respectively. The pseudoheader consists of the source and destination addresses, the upper-layer packet length, and the protocol number. Including this information in the checksum calculation ensures that TCP, UDP, or ICMPv6 do not process packets that were delivered incorrectly, for instance, because of a bit error in the IP header.

IPv6 no longer has a header checksum to protect the IP header, meaning that when a packet header is corrupted by transmission errors, the packet is very likely to be delivered incorrectly. However, higher-layer protocols should be able to detect these problems, so they are not fatal. Also, lower layers almost always employ a *Cyclic Redundancy Check* (CRC) to detect errors.

Extension Headers

To allow special processing along the way, IPv4 allows extension of the IP header with one or more options. These options are rarely used today, both because they do not really solve common problems and because packets with options cannot be processed in the “fast path,” and many routers and firewalls block some or all options. Not unlike the checkout counters at a grocery store, many routers have several “paths” that packets may follow: a fast one, implemented in hardware or highly optimized software, that supports only the most common operations (no checks), and one or more slower paths that use more advanced but slower software code that supports less common operations such as looking at IP options. However, many modern routers have only a fast path, so using additional features does not lead to a performance penalty.

Because the header is of fixed length in IPv6, options cannot be tagged onto the IP header as in IPv4. Instead, they are put in a header of their own that sits between the IPv6 header and the TCP or UDP (or other higher-level protocol) header. The most common extension headers follow:

- *Hop-by-Hop Options*: See the section that follows.
- *Routing*: This header is similar to the *Source Route* option in IPv4.
- *Fragment*: This header is used for fragmentation; see later in this article.
- *Authentication*: This header authenticates the user data and most header fields.
- *Encapsulating Security Payload* (ESP): This header encrypts or authenticates user data.
- *Destination Options*: See the section that follows.

The Hop-by-Hop Options and Destination Options headers are container headers: they have room for multiple suboptions. The Hop-by-Hop Options are processed by all routers along the way. All other options are normally ignored by routers and processed only by the destination. Obviously firewalls, or routers configured to perform filtering, may also look at these options. The Hop-by-Hop Options, Routing, Fragment, and Destination Options extension headers are defined in RFC 2460^[5]. The Authentication and ESP extension headers are part of *IP Security* (IPsec).

Note that there is no standard extension header format, meaning that when a host encounters a header that it does not recognize in the protocol chain, the only thing it can do is discard the packet. Worse, firewalls and routers configured to filter IPv6 have the same problem: as soon as they encounter an unknown extension header, they must decide to allow or disallow the packet, even though another header deeper inside the packet would possibly trigger the opposite behavior. In other words, an IPv6 packet with a TCP payload that would normally be allowed through could be blocked if there is an unknown extension header between the IPv6 and TCP headers.

ICMPv6

The IPv6 version of the ICMP generally serves the same purposes as its IPv4 counterpart, but there are some changes. In IPv4, when a router or the destination host cannot process the packet properly, it sends back an ICMP error message along with the original IP header and the first 8 bytes of the higher-layer header. For UDP and TCP, this is enough for the source of the original host to see which TCP session or UDP association generated the offending packet. Because IPv6 supports an arbitrary number of extension headers between the IPv6 header and the higher-layer header, ICMPv6 returns as much of the original packet as will fit in the minimum MTU size of 1280 bytes. In addition to error messages, which are recognizable by an ICMP type of 127 or lower, there are also informational messages, with a type of 128 or higher. Because informational messages are not the result of an error, they do not include an original packet or part thereof. The most common ICMPv6 message types follow:

- 1:** Destination unreachable
- 2:** Packet too big
- 3:** Time exceeded
- 4:** Parameter problem
- 128:** Echo request
- 129:** Echo reply
- 130:** Multicast listener query
- 131:** Multicast listener report
- 132:** Multicast listener done
- 133:** Router solicitation
- 134:** Router advertisement
- 135:** Neighbor solicitation
- 136:** Neighbor advertisement
- 137:** Redirect message

ICMP and ICMPv6 messages also include a “code” that indicates the exact nature of the ICMP message within a certain type. As with ICMP, ICMPv6 calculates a checksum over the control message, but unlike ICMP, the ICMPv6 checksum calculation also includes a pseudoheader. Another departure from IPv4 is the fact that hosts and routers are required to limit the number of ICMPv6 messages they send. So if a router receives 100 packets per second toward an unreachable destination, it is not supposed to send back ICMPv6 packets at the same rate of 100 per second. The ICMPv6 redirect message works slightly different from the ICMP redirect message in IPv4. Like its IPv4 counterpart, the ICMPv6 redirect can be used by a router to inform a host that it should use a different router to reach the destination in question. But routers can also use the IPv6 Redirect to tell a host that the destination is reachable on the local subnet. Thus two hosts that have addresses in different prefixes can communicate directly after receiving redirects from a router.

Neighbor Discovery

When a system wants to send an IPv6 packet to another system connected to the same subnet or link, it needs to know what MAC address (or “link address” in the new IPv6 terminology) it should address the packet to, unless the interface in question is a point-to-point interface. Neighbor discovery allows systems to discover each other’s MAC addresses, similar to *Address Resolution Protocol* (ARP) on Ethernet with IPv4.

Each IPv6 system joins the “solicited node” multicast group that corresponds to each of its addresses. Because the solicited node group address consists of the prefix **ff02:0:0:0:0:1:ff00::/104** followed by the bottom 24 bits of the address in question, addresses in different prefixes based on the same interface identifier (including the link-local address) all map to the same solicited node address.

Whenever a system needs to find out the link address for another system residing on the same link, it sends a neighbor solicitation to the solicited node address to which the IPv6 address of the remote system maps. The source host includes its own MAC address in the neighbor solicitation, so the neighbor knows where to send the reply.

Neighbor Unreachability Detection

RFC 2461^[6] specifies a procedure for neighbor unreachability detection. IPv6 hosts and routers actively track whether their neighbors are reachable by periodically sending neighbor discovery messages directly to the neighbor. If the neighbor answers, it is reachable; if it does not, there must be some kind of problem, and the system discards the neighbor’s MAC address and tries a regular multicast neighbor discovery procedure, allowing IPv6 systems to detect dead neighbors and neighbors that change their MAC address. But it is most useful to detect dead routers. On a subnet with more than one router, a host can simply install a default route toward another router when the router that it has been using becomes unreachable.

If a router loses its IPv6 address and no longer runs IPv6, Windows XP, Linux, MacOS, and FreeBSD all switch over to another router without incident. However, turning off the active router has much more severe effects: at the very least, ongoing downloads stall for a while, and in some cases, the session breaks. I have no explanation for this difference in behavior.

Stateless Address Autoconfiguration

Hosts and routers always configure link-local addresses on every interface on which IPv6 is enabled. The link-local address is nearly always derived from the interface MAC address, but to guarantee uniqueness, it is necessary to perform *Duplicate Address Detection*, which is discussed later.

When a host has a link-local address, it can obtain one or more global IPv6 addresses by using RFC 2462^[7, 12], *Stateless Address Autoconfiguration*. IPv6 routers send out *Router Advertisement* (RA) packets (ICMPv6 type 134) periodically and in response to router solicitations. The information in RAs includes:

- An 8-bit *cur hop limit* field that tells hosts what value to use in the Hop Limit field of outgoing IPv6 packets
- The *managed address configuration* (M) flag—This flag is not well-defined, but the basic idea is that when it is set, hosts use a stateful mechanism (presumably *Dynamic Host Configuration Protocol Version 6* [DHCPv6]) to configure their addresses, and when the flag is not set, they use stateless address autoconfiguration.
- The *other stateful configuration* (O) flag—This flag is similar to the M flag, but indicates that the host should use a stateful mechanism to discover nonaddress configuration information.
- A 16-bit *router lifetime* value in seconds—This value tells hosts how long the default route that was created as the result of this RA should remain valid.
- The 32-bit *reachable time* value in milliseconds—This value indicates how long a neighbor should be considered reachable after receiving a “reachability confirmation,” which is generally a neighbor advertisement but could be any packet.
- The 32-bit *retrans timer* value in milliseconds—The retrans timer tells hosts how long they should wait before retransmitting neighbor solicitation messages when there is no answer.

When fields that determine a value are set to zero, this means the value is not specified in the RA, so hosts must discover that value through other means. In addition to the preceding, router advertisements may also contain one or more options, such as:

- *Source link-layer address*, the router MAC address
- *MTU*, the maximum packet size that should be used on this subnet
- *Prefix information*, which specifies the prefixes used on the subnet and their properties

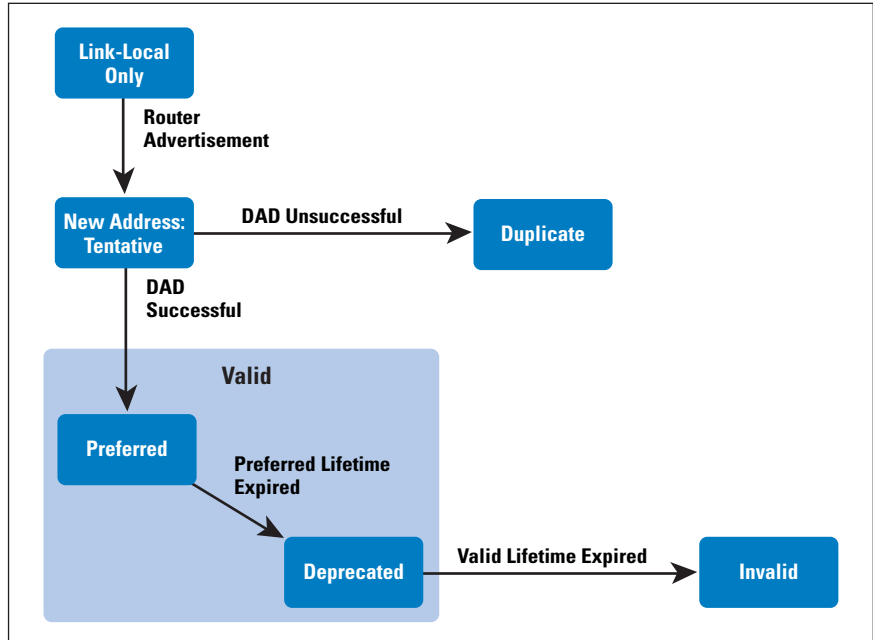
The prefix information option, in turn, has its own list of attributes:

- The *address prefix* itself and its length—For stateless address auto-configuration to work, the prefix must be 64 bits long.
- The *on-link* flag—This flag tells hosts that the prefix is “on-link,” so systems with addresses within this prefix are reachable on the subnet in question without help from a router.
- The *autonomous address configuration* flag—This flag tells hosts that they can create an address for themselves by combining this prefix with an interface identifier.
- A 32-bit *valid lifetime* in seconds—This value indicates how long the prefix should be considered on-link and how long autoconfigured addresses using the prefix can be used.
- A 32-bit *preferred lifetime* in seconds—This flag tells hosts how long autoconfigured addresses using this prefix are preferred.

Duplicate Address Detection

To avoid the situation where two IPv6 systems use the same address, systems perform *Duplicate Address Detection* for (nearly) all new IPv6 addresses before they are used. Duplicate address detection is done for global unicast addresses—and not just for those created using stateless address autoconfiguration, but also for link-local addresses. For obvious reasons, there is no duplicate detection for anycast addresses, because the whole point of anycast is that multiple systems have the same address.

Figure 2: The Lifecycle of an IPv6 Address



As depicted in Figure 2, a host starts with only a link-local address. Duplicate address detection is also done for the link-local address, but this is not shown in the figure.

When a host receives a router advertisement that contains one or more prefixes with the autonomous address configuration flag set, the host creates addresses with interface identifiers derived from the IEEE 64-bit *Extended Unique Identifier* (EUI-64) and possibly also a randomly generated one, if the host uses RFC 3041^[8] address privacy. The host marks the resulting addresses as “tentative” and proceeds to execute the duplicate address detection procedure by joining the solicited node multicast group for the address in question and sending out one or more neighbor solicitation messages for the address. (If the number of duplicate address detection retries is configured to be zero, no duplicate detection is performed.) Only when there is no answer is the address used. If there is a conflict, the system is supposed to log the error and wait for manual intervention.

Address Lifetime

After successfully maneuvering past the duplicate address detection hurdle, addresses configured through stateless address autoconfiguration can be used until the “preferred lifetime” from the router advertisement message expires. In most cases, the lifetime does not expire because new RAs refresh the timers. But if there are no more RAs, eventually the preferred lifetime elapses and the address becomes “deprecated.” New sessions should not use deprecated addresses but should choose “preferred” (nondeprecated) addresses, if available. However, existing sessions will continue to use the deprecated address. Eventually, the “valid lifetime” also runs out, and the deprecated address is removed from the interface, breaking any sessions that are still using the address.

Renumbering

Having different preferred and valid timers for the router advertisement itself and also for any prefixes contained in it makes it possible to do two things: renumber easily and cause more problems. It is even possible to do both at the same time. With stateless autoconfiguration, renumbering is easy: you simply give the router an address in the new prefix and set the preferred lifetime for the old prefix to zero, making hosts create one or more new addresses and deprecate any existing ones in the old prefix as soon as they receive the resulting router advertisement. After that, all new communication should start using the new address immediately. Existing TCP sessions and UDP associations continue to use the same address as before. After some time, all communication that started before the change should have stopped so that the old addresses can be removed safely.

This process is slightly more complex than it seems at first glance: as a precaution against attackers, hosts are not supposed to trust a valid lifetime of less than 7200. So make sure that the hosts have received at least one RA after setting the valid lifetime to 7200, and then set both the lifetimes to zero and remove the autonomous address configuration flag for the prefix. Two hours later, all hosts should have removed the addresses in this prefix, so you can remove the prefix from the router.

Beware that when you renumber because you are switching from one ISP to another, it is unavoidable that at some point, packets with source addresses in address space from ISP A end up at ISP B, or the other way around. If ISP B employs antispoofing or ingress filtering, it will not allow these packets through, so reduced connectivity will result. You can ask one ISP to remove the filters temporarily and then send out all your outgoing traffic over that ISP (or one that did not filter in the first place). However, do not expect too much cooperation from your ISP unless you are a valued customer.

Address Prefix and Router Lifetime Mismatch

Earlier, I mentioned the potential for causing more problems because router advertisements and the prefixes they contain have independent lifetimes. This scenario allows for four permutations:

- The RA lifetime is valid, and the prefix lifetime is valid: IPv6 works.
- The RA lifetime is invalid, and the prefix lifetime is invalid: IPv6 is disabled.
- The RA lifetime is valid, but the prefix lifetime is invalid: The system has an IPv6 default route but no global IPv6 address.
- The RA lifetime is invalid, but the prefix lifetime is valid: The system has a global IPv6 address but no IPv6 default route.

When a host has no global addresses but does have an IPv6 default route (case 3), it cannot reach the rest of the IPv6 Internet. Unfortunately, FreeBSD and MacOS hosts do not know that: they try anyway, with long delays as a result. Only after trying all the remote destination IPv6 addresses and timing out, the system falls back on IPv4 (for applications that try more than one address). Linux, on the other hand, does not install the IPv6 default route or ignores it when no global IPv6 addresses are present, so the timeout is immediate.

Windows XP does install the default route but magically manages to avoid lengthy timeouts anyway. On the other hand, Windows XP suffers timeouts when it has an IPv6 address but no default route (case 4) because Windows implements the on-link assumption: it first performs neighbor discovery on the local subnet for any IPv6 addresses. Only after neighbor discovery times out does Windows revert to IPv4. FreeBSD and MacOS, however, do not implement the on-link assumption, so they immediately notice that the IPv6 destination address is unreachable and revert to IPv4—if an IPv4 address is available and the application cycles through all addresses. With Linux, the default route does not seem to expire even though the timers eventually reach zero and lower. But addresses do expire and are removed when the lifetime for the associated prefix times out.

Address Selection

Choice is good, but it comes with problems of its own. The explicit support for multiple addresses in IPv6 requires the system or applications to choose which address to use for a given communication session. The coexistence of IPv4 and IPv6 in the same host makes this situation even more pressing. RFC 3484^[9] provides guidelines in this area—it lists no fewer than 10 rules for choosing a destination address and 8 rules for selecting a source address. Most of these rules are fairly obvious, such as preferring a nondeprecated address over a deprecated one and not using a link-local source address to communicate with a destination that has a global address. It gets more interesting with the “policy table.” On systems that support this mechanism, such as Windows XP and FreeBSD 5.4, the administrator can instruct the system to prefer certain address ranges over others.

Path MTU Discovery and Fragmentation

Because routers cannot fragment IPv6 packets, *Path MTU Discovery* (PMTUD) is mandatory in all cases where links with MTUs larger than 1280 bytes are used for IPv6, so it is imperative that routers generate ICMPv6 packet-too-big messages and that these messages make it back to the source of the offending packet. Filtering out these ICMPv6 messages makes it impossible to communicate reliably.

If you decide that you must filter ICMPv6 packet-too-big messages, you *must* use an MTU equal to the IPv6 mandatory minimum of 1280 bytes across your network so there is no need for PMTUD.

Upon reception of a packet-too-big message, TCP reduces its packet size to accommodate the smaller MTU on the path in question. However, protocols that run over UDP often cannot arbitrarily reduce their packet size. In IPv4, UDP packets are generally sent without the “don’t fragment” bit set, so routers fragment them if necessary. In IPv6, this setup is not possible; if the packet is too large, the source host has to fragment it. The source host does this by first splitting the packet into unfragmentable and fragmentable parts. The IPv6 header and any headers that must be processed by routers along the way make up the unfragmentable part; the payload data and any headers that have to be processed only on the destination host are the fragmentable part. The fragmentable part is then split into as many parts as required to fit in the path MTU, and each part is transmitted as a packet containing the unfragmentable part, a fragment header, and one of the fragments of the fragmentable part.

The fragment header is 8 bytes, and except for a “next header” field and two reserved fields, it contains the same fragment offset, more fragments, and identification fields as the IPv4 header. The identification field is now 32 bits long and is used to indicate which fragments belong to the same original packet. All fragments except the last one have the “more-fragments” bit set and are multiples of 8 bytes.

After receiving the first fragment (which is not necessarily the first fragment of the original packet), a host waits up to 60 seconds for all other fragments to come in and, if they do, reassembles the original packet by combining all the fragments with the same source and destination addresses and identification field into a single packet. If one or more fragments is lost, the packet cannot be reassembled, so the entire packet is lost.

Note that IPv6 fragmentation has the same problem as IPv4 fragmentation: the TCP or UDP port numbers are available only in the first fragment, making it hard for firewalls and the like to filter fragmented packets. Common solutions are to reassemble the packet prior to filtering or to discard all fragments.

DHCPv6

DHCPv6 (RFC 3315^[10]) is the IPv6 version of the DHCP. Because IPv6 has stateless address autoconfiguration, DHCP occupies a very different part of the landscape in IPv6 compared to IPv4. Although the details are different in the by-now-expected places (address length, use of multicasts, some streamlining), the DHCPv6 protocol itself is quite similar to the IPv4 version of DHCP. The more important differences are the way in which the protocol is used. DHCPv6 has three purposes:

- *Address configuration*: Giving out addresses to individual hosts
- *Nonaddress configuration*: Giving out other configuration information, such as DNS resolver addresses and domain search lists
- *Prefix delegation*: Giving out entire prefixes to routers (RFC 3633^[11])

A DHCPv6 client interested in an address or other configuration information sends out a *solicit* message indicating its needs to the link-local scope multicast address **ff02::1:2**, port 547. (Server-to-client messages are addressed to port 546.) DHCPv6 servers that receive the *solicit* message either directly or forwarded by a relay and can accommodate the request respond with an *advertise* message. The client considers the offers in the various *advertise* messages and directs a *request* message to the server of its choice. The server then replies with a *reply* message, confirming the address or configuration information. Alternatively, if the client wants to receive only configuration information and no addresses or prefixes, it can send a *request-information* message, and the server immediately sends back a reply message, so only half the messages are exchanged and the whole process completes much faster. The client can also use the “rapid commit” option to indicate that it wants to use the expedited procedure for address or prefix assignment if it is fairly certain that it will take up the offer from the first DHCPv6 server that responds.

As expected, IPv6 addresses assigned with DHCPv6 come with a preferred and a valid lifetime. Sometime before this timer expires, the client sends a *renew* message, asking the server if it can continue to use the address. When it has no more use for the address, the client sends a *release* message. Less common situations have other messages.

To allow servers to recognize clients, each device that implements DHCPv6 has *DHCP Unique Identifier* (DUID). In IPv4, DHCP clients use a MAC address or user-supplied string as a Client Identifier. In DHCPv6 this client identifier is always the DUID. Devices can create their DUID in various ways, as long as the DUID is unique and not subject to change, if at all possible.

DHCPv6 supports an authentication mechanism that allows clients and servers to interact in a secure way, so third parties cannot inject false DHCP messages or modify legitimate ones. However, this mechanism must be preconfigured manually on all servers and clients, partially negating the advantages of DHCP over manual configuration.

An interesting use of DHCPv6 is *Prefix Delegation* (PD). With DHCPv6 PD, routers request a prefix that they then use to number one or more of their interfaces, supporting stateless address autoconfiguration for hosts connected to that interface. By creatively borrowing the DHCP timers and reusing them in router advertisements, a whole site can be renumbered by changing a single setting in a DHCPv6 configuration on a DHCPv6 server or a router functioning as a DHCPv6 PD server.

Ed.: This article is adapted from chapter 8 of *Running IPv6* by Iljitsch van Beijnum, published by Apress in 2005, ISBN 1590595270. The article differs from the chapter in that it has been edited for size and the vendor-specific examples have been removed. Used with permission. For information about the book, see:

<http://www.apress.com/book/bookDisplay.html?bID=10026>

References

- [1] Karrenberg D., Ross G., Wilson P., and Nobile L., “Development of the Regional Internet Registry System,” *The Internet Protocol Journal*, Volume 4, No. 4, December 2001.
- [2] Carpenter, B., Fink, B., and Moore, K., “Connecting IPv6 Routing Domains Over the IPv4 Internet,” *The Internet Protocol Journal*, Volume 3, No. 1, March 2000.
- [3] Nichols, K., Blake, S., Baker, F., and Black, D., “Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers,” RFC 2474, December 1998.

- [4] Postel, J., “Internet Protocol,” RFC 791, September 1981.
- [5] Deering, S. and Hinden, R., “Internet Protocol, Version 6 (IPv6) Specification,” RFC 2460, December 1998.
- [6] Narten, T., Nordmark, E., and Simpson, W., “Neighbor Discovery for IP Version 6 (IPv6),” RFC 2461, December 1998.
- [7] Narten, T. and Thomson, S., “IPv6 Stateless Address Autoconfiguration,” RFC 2462, December 1998.
- [8] Narten, T. and Draves, R., “Privacy Extensions for Stateless Address Autoconfiguration in IPv6,” RFC 3041, January 2001.
- [9] Draves, R., “Default Address Selection for Internet Protocol Version 6 (IPv6),” RFC 3484, February 2003.
- [10] Droms, R., Ed., Bound, J., Volz, B., Lemon, T., Perkins, C., and Carney, M., “Dynamic Host Configuration Protocol for IPv6 (DHCPv6),” RFC 3315, July 2003.
- [11] Troan, O. and Droms, R., “IPv6 Prefix Options for Dynamic Host Configuration Protocol (DHCP) Version 6,” RFC 3633, December 2003.
- [12] François Donzé, “IPv6 Address Autoconfiguration,” *The Internet Protocol Journal*, Volume 7, No. 2, June 2004.

ILJITSCH VAN BEIJNUM holds a Bachelor of Information and Communication Technology degree from the Haagse Hogeschool in The Hague, Netherlands. In 1995, he found himself in the emerging Internet Service Provider business. There he learned about system administration, IP networking, and especially routing. After first starting a small ISP with four others and working as a senior network engineer for UUNET Netherlands, he became a freelance consultant in 2000. Not long after that, he started contributing to the IETF Multihoming in IPv6 working group. He wrote the book *BGP: Building Reliable Networks with the Border Gateway Protocol*, ISBN 0-596-00254-8, published by O'Reilly in 2002, and *Running IPv6*, ISBN 1590595270, published by Apress in 2005. E-mail: iljitsch@muada.com

Book Reviews

Electronic Brains

Electronic Brains, Stories from the Dawn of the Computer Age, by Mike Hally, ISBN 0-309-09630-8, Joseph Henry Press, 2005.

Electronic Brains is a personal account from the early days of computing that describes the childhood of a technology that is little more than 50 years old. The book originated as a BBC radio programme, still accessible at <http://www.bbc.co.uk/radio4/science/electronicbrains.shtml>. Mike Hally traveled over the globe looking for the first “computers” and the stories from the dawn of a new age. This book contains the results of the investigation, giving a first-hand testimony of hard work, passion, and amazing developments that shaped the second half of the last century.

Organization

Chapter 1, “From ABC to ENIAC,” presents the development of what is commonly accepted as the first computer, the ENIAC, a computer that replaced calculating machines and people making the operations in ballistic trajectories analysis by hand. Credit is given to John Atanassof and Clifford Berry, the developers of ABC, possibly the first operational computer in the world.

Development of the UNIVAC, the computer famed by predicting the result of the 1952 U.S. presidential election, is presented in Chapter 2. Designed by Eckert and Mauchly, the developers of ENIAC, UNIVACs were commercial computers used for processing census data and so well marketed that the term “UNIVAC” was used as a synonym for “computer.”

Chapter 3 looks at the development of the Rand 409, maybe the first mass-produced computer. The 409 was a medium-sized computer, with a price tag of US\$100,000 that compared favorably against UNIVAC’s \$1 million, achieving a sell rate of one per week.

“Computing in Great Britain” is the focus of Chapter 4, where credit is given to Maurice Wilkes and Alan Turing. A worthy detail that gives a glimpse of the technical difficulties overcome is the description of memory based on mercury delay-lines, where binary data was stored using sound pulses on tubes filled with mercury engineered in such way that the delay from transmitter to receiver allowed the electronics to do the calculation before the data in memory was needed at the receiver side.

Perhaps the strangest computer development is set forth in Chapter 5. The *Lyons Electronics Office* (LEO) was a computer developed by a large catering company to expedite its clerical operations. LEO was possibly the first commercial computer in the world, so successful that the catering company began to produce and sell it to other corporations.

Chapter 6 describes the efforts by USSR scientists to develop computing technology. More than one development was made; it is not clear which was the first soviet computer, and the developments were secret—in some cases very specialized, such as a computer with ternary logic instead of the currently used binary logic (ENIAC used decimal logic).

Chapter 7 focuses on computing developments in Australia, work that did not last because the funds were scarce and sometimes the budget was assigned to other sciences, such as radiophysics. Here we can see that computers were used for purposes totally different than their uses in cold-war countries; for example, they were used to answer crossword puzzles—strange if we consider that the disk had a capacity of 3 KB.

A strange computer, formally known as *Hydraulic Economics Computer*, is described in Chapter 8. It was not a typical computer—it was a system developed to show the interrelation between macroeconomic variables using colored water, pumps, and valves. Universities, central banks, and Ford bought the computer, and four of them survive in different parts of the world. The emergence of IBM is the subject of Chapter 9, which presents IBM as a late adopter of computing technology that eventually became the leader of the computer age. We learn that the first computer produced by IBM was the IBM 701; after that came the IBM 1401 and then the IBM 360—the system that consolidated IBM as the ruler in the computing world.

Summary

From the ABC to the well-known ENIAC and UNIVAC, *Electronic Brains* is a testimony to the people who worked day and night to accomplish something that few others understood. Motivated mainly by passion and with little to no economic support, team spirit is a common factor in all the computer developments: “...it was like a brotherhood! We would help each other in case someone got stuck on a particular activity. I would have gone anywhere with those guys. I’ve never had such unified job environment. We knew we were pushing back the frontiers.”

Electronic Brains is an enjoyable book that I recommend to any person with interest in computers and technology. Computer historians could scoff at the rather simple analysis of technical details, but this is not a technical book. The value of *Electronic Brains* is the first-hand account of early undertakings and the multiple-country investigation that is presented. With many anecdotes, this book will serve as a witness to the pioneers of a new era, the computing era.

—Claudio Gutiérrez

claudio.gutierrez.m@gmail.com

Business 2010

Business 2010—Mapping the Commercial Landscape, by Ian Pearson and Michael Lyons, ISBN 1-84439-105-1, Published by Spiro Press, <http://www.spiropress.com/>

This interesting book explores how trends in technology, economic factors, social changes, and evolving attitudes to technology will reshape the business landscape by the year 2010. The book describes its subject matter in terms that are understandable and interesting to both technical and nontechnical audiences. It is valuable to technologists because it expands their perception of the future beyond that which is available through traditional sources such as vendor roadmap sessions by linking closely commercial, technical, and social trends.

Organisation

The book is divided into three main sections. The first looks at the major influences on future business: technological progress, changing attitudes, social forces, and economics. The implications of these factors are then examined, and finally the application of the analysis to business strategy is examined. These ideas are then pulled together in a succinct and easily understood conclusion.

Pearson and Lyons focus on the effect of particular techniques. Some of these, such as self-organising systems and the mimicking of natural phenomena (“biomimetrics”), are fairly unconventional, but others, such as increased miniaturisation, wireless devices, low-cost computing and networking, the semantic Web, and artificial intelligence will be more familiar. The Internet and its potential effect on financial transactions and taxation features heavily. The authors note that attitudes to technology are changing and adoption cycles are reducing, describing the impact that technology has had on the physical labour market and the likely future impact on knowledge workers. The authors consider the economic implications of the exploitation of information, looking at the relative cost of creation and reproduction when compared with more traditional goods and services.

The next three chapters look at the implications of this analysis, starting by looking at numerous trade-offs and counter-balancing forces, such as the effect of the “browser wars” and the relationships between customers and producers. The importance of customer and worker information to a commercial organisation and the problems arising from its exploitation are described. The discussion then considers how the knowledge economy changes the importance of physical assets and commercial relationships, followed by an examination of the political and organisational implications of technology.

Finally the authors look at the business effects, starting with the ease of transferring information between systems. They note that corporate intranets make both the devolving of authority through outsourcing and the imposition of increased command and control through micro-management easier.

Pearson and Lyons suggest that new technology alters the value chains that influence businesses, leading to more temporary business relationships, their replacement by “value-nets,” and the rise of the virtual company. This section concludes by looking at globalisation—how goods and services are paid for and some of the implications for taxation.

The authors ask the question—how can business adapt? They start their analysis by examining the interactions between the physical and mental worlds and cyberspace, noting that a strategic analysis works only if the forces acting on a business do not change too rapidly. As change becomes more rapid, there will be no time to develop business cases, because first-mover advantage will be the only advantage a business can have. Pearson and Lyons conclude that the critical factors in allowing cyber-economy to grow are ease of navigation and the effective use of branding. They conclude by examining who will be the winners and losers in business in the year 2010—and why.

Synopsis

This book is succinct and well-written, covering a complex but interesting field in just under 200 pages. The authors paint a convincing description of future business trends, exploring the technical, commercial, economic, and political pressures that will influence them. Their cause, effect, and potential response treatment leads the reader through the subject in a way that is both interesting and instructive. The authors are not afraid to be controversial and at times they take the reader into some very unfamiliar territory, adding extra spice to the book.

While other books are available that look at the future from a more technologically orientated perspective, this book is one of the few that manages to couple the developments in the commercial and technical worlds, thereby giving a more comprehensive viewpoint. In an age when technologists are increasingly being asked to take more of a commercial view, this can only be a good thing. The approach taken has much in common with that taken by Alvin Toffler in his books *Future Shock* and *The Third Wave*. An updated treatment like this is to be welcomed.

The Authors

Ian Pearson works for British Telecom (BT) as its chief futurologist; he is a well-known speaker on future technology trends and has published extensively in this field. Michael Lyons also works for BT and has more than 30 years of research experience in the telecoms industry. He has recently been working in the fields of decision support systems and long-term research issues, leading a research team in BT's Research and Venturing department. Pearson is described as an “unfettered thinker” and Lyons as a “pragmatic modeller,” characteristics which give the book its balanced view.

—Edward Smith, BT, UK

edward.a.smith@btinternet.com

ICANN Ratifies Global Policy for Allocation of IPv6 Address Space

On September 7, 2006, the ICANN Board ratified the *Global Policy for Allocation of IPv6 Address Space*. This policy provides for the allocation of IPv6 address space from ICANN to the *Regional Internet Registries* (RIRs).

On July 13, 2006, the Secretary of the *Address Supporting Organization* (ASO) *Address Council* (AC) forwarded to ICANN the proposed global policy for allocation of IPv6 address space. This proposed global policy had been submitted to the ASO AC by the Executive Council of the *Number Resource Organization* (NRO) on June 6, 2006, and adopted by the ASO AC on July 12, 2006. Each RIR community individually discussed the policy and approved its adoption via their own policy development processes. The IPv6 Allocation Policy document is available from the ASO Website:

<http://aso.icann.org/docs/aso-global-ipv6.pdf>

See also:

<http://www.icann.org/announcements/announcement-11sep06.htm>

<http://www.nro.net>

IP addressing in China and the Myth of Address Shortage

In recent years, various sources have repeated a myth that the IPv4 address pool is close to exhaustion. Many of these stories also falsely claim that there are fewer IPv4 addresses allocated to China than to some individual US universities. The *Asia Pacific Network Information Centre* (APNIC) is committed to countering this myth and has published an article in its newsletter *Apster* on this topic. The article is available here:

<http://www.apnic.net/news/hot-topics/internet-gov/ip-china.html>

Calendar of Internet-related Events

The *Internet Society* (ISOC) maintains an online list of meetings and conferences, see:

<http://geneva.isoc.org/events/>

Don't forget to tell us if you move!

We receive quite a lot of IPJ return mail marked as “undeliverable.” If you change your address please let us know by either using the IPJ subscription tool or sending an e-mail with the new information to ipj@cisco.com. Your cooperation is much appreciated.

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Cisco, Cisco Systems, and the Cisco
Systems logo are registered
trademarks of Cisco Systems, Inc. in
the USA and certain other countries.
All other trademarks mentioned in this
document are the property of their
respective owners.*

*Copyright © 2006 Cisco Systems Inc.
All rights reserved.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA
--

The Internet Protocol Journal

December 2006

Volume 9, Number 4

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
SYN Flooding Attacks.....	2
XML Networking	17
Letters to the Editor.....	33
Book Review	37
Fragments	40
Call for Papers.....	43

Internet security and stability are topics we keep returning to in this journal. So far we have mainly focused on technologies that protect systems from unauthorized access and ensure that data in transit over wired or wireless networks cannot be intercepted. We have discussed security-enhanced versions of many of the Internet core protocols, including the *Border Gateway Protocol* (BGP), *Simple Network Management Protocol* (SNMP), and the *Domain Name System* (DNS). You can find all these articles by visiting our Website and referring to our index files. All back issues continue to be available in both HTML and PDF formats. In this issue, Wesley Eddy explains a vulnerability in the *Transmission Control Protocol* (TCP) in which a sender can overwhelm a receiver by sending a large number of SYN protocol exchanges. This form of *Denial of Service* attack, known as *SYN Flooding*, was first reported in 1996, and researchers have developed several solutions to combat the problem.

Speaking of Internet stability, at 12:26 GMT on December 26, 2006, an earthquake of magnitude 6.7 struck off Taiwan's southern coast. Six submarine cables were damaged, resulting in widespread disruption of Internet service in parts of Asia. We hope to bring you more details and analysis of this event in a future issue of IPJ. The topic will also be discussed at the next *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT), which will take place in Bali, Indonesia, February 21 through March 2, 2007. For details see: <http://www.apricot2007.net>

The design and operation of systems that use Internet protocols for communication in conjunction with advanced applications—such as an e-commerce system—require the use of a certain amount of “middleware.” This software, largely hidden from the end user, has been the subject of a great deal of development and standardization work for several decades. An important component of today's Web systems is the *Extensible Markup Language* (XML). Silvano Da Ros explains how XML networking can be used as a critical building block for network application interoperability.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Defenses Against TCP SYN Flooding Attacks

by Wesley M. Eddy, Verizon Federal Network Systems

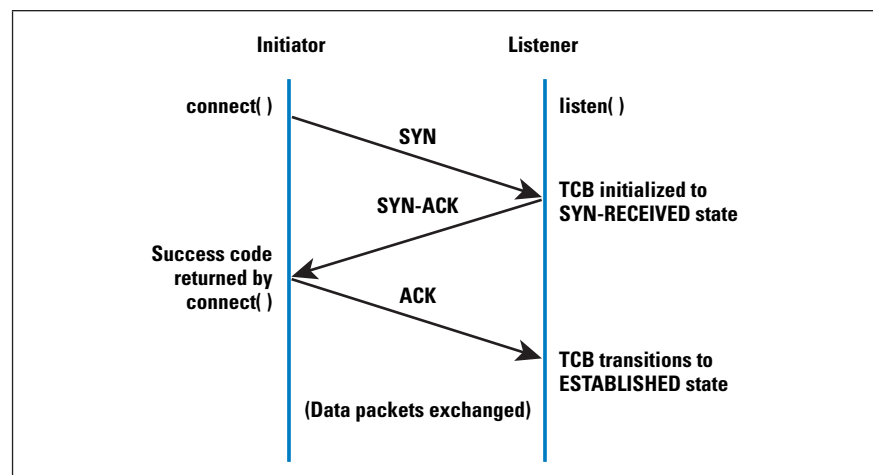
This article discusses a specific *Denial of Service* (DoS) attack known as *TCP SYN Flooding*. The attack exploits an implementation characteristic of the *Transmission Control Protocol* (TCP), and can be used to make server processes incapable of answering a legitimate client application's requests for new TCP connections. Any service that binds to and listens on a TCP socket is potentially vulnerable to TCP SYN flooding attacks. Because this includes popular server applications for e-mail, Web, and file storage services, understanding and knowing how to protect against these attacks is a critical part of practical network engineering.

The attack has been well-known for a decade, and variations of it are still seen. Although effective techniques exist to combat SYN flooding, no single standard remedy for TCP implementations has emerged. Varied solutions can be found among current operating systems and equipment, with differing implications for both the applications and networks under defense. This article describes the attack and why it works, and follows with an overview and assessment of the current tactics that are used in both end hosts and network devices to combat SYN flooding attacks.

Basic Vulnerability

The SYN flooding attack became well-known in 1996, when the magazines *2600* and *Phrack* published descriptions of the attack along with source code to perform it^[1]. This information was quickly used in attacks on an Internet service provider's (ISP's) mail and Telnet servers, causing outages that were widely publicized in *The Washington Post* and *The Wall Street Journal* (among other venues). CERT quickly released an advisory on the attack technique^[2].

Figure 1: Normal TCP 3-Way Handshake



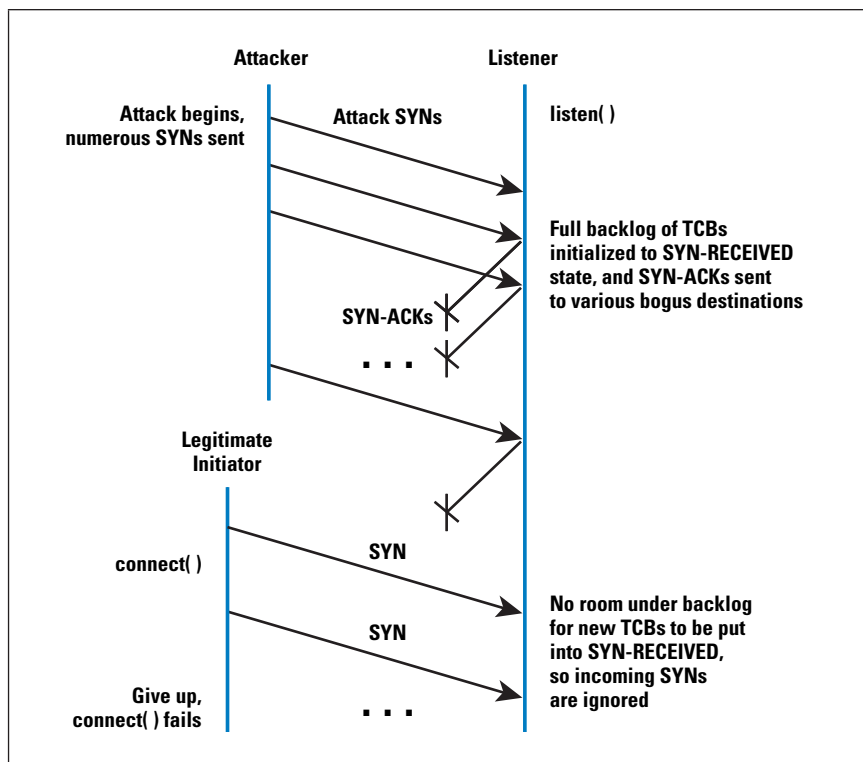
The basis of the SYN flooding attack lies in the design of the 3-way handshake that begins a TCP connection. In this handshake, the third packet verifies the initiator's ability to receive packets at the IP address it used as the source in its initial request, or its return reachability. Figure 1 shows the sequence of packets exchanged at the beginning of a normal TCP connection (refer to RFC 793 for a detailed description of this process).

The *Transmission Control Block* (TCB) is a transport protocol data structure (actually a set of structures in many operations systems) that holds all the information about a connection. The memory footprint of a single TCB depends on what TCP options and other features an implementation provides and has enabled for a connection. Usually, each TCB exceeds at least 280 bytes, and in some operating systems currently takes more than 1300 bytes. The TCP SYN-RECEIVED state is used to indicate that the connection is only half open, and that the legitimacy of the request is still in question. The important aspect to note is that the TCB is allocated based on reception of the SYN packet—before the connection is fully established or the initiator's return reachability has been verified.

This situation leads to a clear potential DoS attack where incoming SYNs cause the allocation of so many TCBs that a host's kernel memory is exhausted. In order to avoid this memory exhaustion, operating systems generally associate a "backlog" parameter with a listening socket that sets a cap on the number of TCBs simultaneously in the SYN-RECEIVED state. Although this action protects a host's available memory resource from attack, the backlog *itself* represents another (smaller) resource vulnerable to attack. With no room left in the backlog, it is impossible to service new connection requests until some TCBs can be reaped or otherwise removed from the SYN-RECEIVED state.

Depleting the backlog is the goal of the TCP SYN flooding attack, which attempts to send enough SYN segments to fill the entire backlog. The attacker uses source IP addresses in the SYNs that are not likely to trigger any response that would free the TCBs from the SYN-RECEIVED state. Because TCP attempts to be reliable, the target host keeps its TCBs stuck in SYN-RECEIVED for a relatively long time before giving up on the half connection and reaping them. In the meantime, service is denied to the application process on the listener for legitimate new TCP connection initiation requests. Figure 2 presents a simplification of the sequence of events involved in a TCP SYN flooding attack.

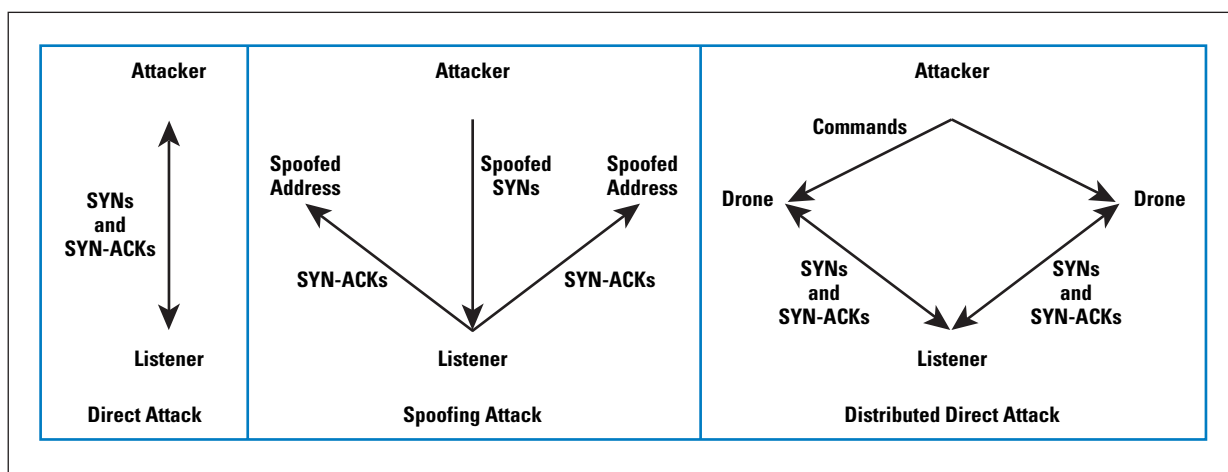
Figure 2: Attack Demonstration: Enough illegitimate TCBs are in SYN-RECEIVED that a legitimate connection cannot be initiated.



Attack Methods

The scenario pictured in Figure 2 is a simplification of how SYN flooding attacks are carried out in the real world, and is intended only to give an understanding of the basic idea behind these types of attacks. Figure 3 presents some variations that have been observed on the Internet.

Figure 3: Some Variants of the Basic Attack



Direct Attack

If attackers rapidly send SYN segments without spoofing their IP source address, we call this a *direct attack*. This method of attack is very easy to perform because it does not involve directly injecting or spoofing packets below the user level of the attacker's operating system. It can be performed by simply using many TCP *connect()* calls, for instance. To be effective, however, attackers must prevent their operating system from responding to the SYN-ACKS in any way, because any ACKs, RSTs, or *Internet Control Message Protocol* (ICMP) messages will allow the listener to move the TCB out of SYN-RECEIVED. This scenario can be accomplished through firewall rules that either filter outgoing packets to the listener (allowing only SYNs out), or filter incoming packets so that any SYN-ACKS are discarded before reaching the local TCP processing code.

When detected, this type of attack is very easy to defend against, because a simple firewall rule to block packets with the attacker's source IP address is all that is needed. This defense behavior can be automated, and such functions are available in off-the-shelf reactive firewalls.

Spoofing-Based Attacks

Another form of SYN flooding attacks uses IP address spoofing, which might be considered more complex than the method used in a direct attack, in that instead of merely manipulating local firewall rules, the attacker also needs to be able to form and inject raw IP packets with valid IP and TCP headers. Today, popular libraries exist to aid with raw packet formation and injection, so attacks based on spoofing are actually fairly easy.

For spoofing attacks, a primary consideration is address selection. If the attack is to succeed, the machines at the spoofed source addresses must not respond to the SYN-ACKS that are sent to them in any way. A very simple attacker might spoof only a single source address that it knows will not respond to the SYN-ACKS, either because no machine physically exists at the address presently, or because of some other property of the address or network configuration. Another option is to spoof many different source addresses, under the assumption that some percentage of the spoofed addresses will be unresponsive to the SYN-ACKS. This option is accomplished either by cycling through a list of source addresses that are known to be desirable for the purpose, or by generating addresses inside a subnet with similar properties.

If only a single source address is repetitively spoofed, this address is easy for the listener to detect and filter. In most cases a larger list of source addresses is used to make defense more difficult. In this case, the best defense is to block the spoofed packets as close to their source as possible.

Assuming the attacker is based in a “stub” location in the network (rather than within a transit *Autonomous System* (AS), for instance), restrictive network ingress filtering^[7] by stub ISPs and egress filtering within the attacker’s network will shut down spoofing attacks—if these mechanisms can be deployed in the right places. Because these ingress/egress filtering defenses may interfere with some legitimate traffic, such as the Mobile IP triangle routing mode of operation, they might be seen as undesirable, and are not universally deployed. *IP Security* (IPsec) also provides an excellent defense against spoofed packets, but this protocol generally cannot be required because its deployment is currently limited. Because it is usually impossible for the listener to ask the initiator’s ISPs to perform address filtering or to ask the initiator to use IPsec, defending against spoofing attacks that use multiple addresses requires more complex solutions that are discussed later in this article.

Distributed Attacks

The real limitation of single-attacker spoofing-based attacks is that if the packets can somehow be traced back to their true source, the attacker can be easily shut down. Although the tracing process typically involves some amount of time and coordination between ISPs, it is not impossible. A distributed version of the SYN flooding attack, in which the attacker takes advantage of numerous drone machines throughout the Internet, is much more difficult to stop. In the case shown in Figure 3, the drones use direct attacks, but to increase the effectiveness even further, each drone could use a spoofing attack and multiple spoofed addresses.

Currently, distributed attacks are feasible because there are several “botnets” or “drone armies” of thousands of compromised machines that are used by criminals for DoS attacks. Because drone machines are constantly added or removed from the armies and can change their IP addresses or connectivity, it is quite challenging to block these attacks.

Attack Parameters

Regardless of the method of attack, SYN flooding can be tuned to use fewer packets than a brute-force DoS attack that simply clogs the target network by sending a high volume of packets. This tuning is accomplished with some knowledge of the listener’s operating system, such as the size of the backlog that is used, and how long it keeps TCBs in SYN-RECEIVED before timing out and reaping them. For instance, the attacker can minimally send a quick flight of some number of SYNs exactly equal to the backlog, and repeat this process periodically as TCBs are reclaimed in order to keep a listener unavailable perpetually.

Default backlogs of 1024 are configured on some recent operating systems, but many machines on the Internet are configured with backlogs of 128 or fewer. A common threshold for retransmission of the SYN-ACK is 5, with the timeout between successive attempts doubled, and an initial timeout of 3 seconds, yielding 189 seconds between the time when the first SYN-ACK is sent and the time when the TCB can be reclaimed.

Assuming a backlog of 128 and that an attacker generates 40-byte SYN segments (with a 20-byte TCP header plus a 20-byte IP header), the attacker has to send only 5.12 kilobytes (at the IP layer) in order to fill the backlog. Repeated every 189 seconds, this process gives an average data rate of only 27 bytes per second (easily achievable even over dialup links). This data rate is in stark contrast to DoS attacks that rely on sending many megabits per second of attack traffic. Even if a backlog of 2048 is used, the required data rate is only 433 bytes per second, so it is clear that the ease of attack scales along with increases to the backlog—and more sophisticated defenses are needed.

Lessons Learned

The protocol flaw in TCP that makes SYN flooding effective is that for the small cost of sending a packet, an initiator causes a relatively greater expense to the listener by forcing the listener to reserve state in a TCB. An excellent technique for designing protocols that are robust to this type of attack is to make the listener side operate statelessly^[3] until the initiator can demonstrate its legitimacy. This principle has been used in more recent transport protocols, such as the *Stream Control Transmission Protocol* (SCTP)^[4], which has a 4-way handshake, with listener TCB state being created only after the initiator echoes back some “cookie” bytes sent to it by the listener. This echo proves to some extent that the initiator side is at the address it appears to be (that is, it has return reachability) and is not attempting a SYN flooding style of attack.

Outside of transport protocols and TCBs, security protocols also commonly use this defense technique. For instance, the *Internet Key Exchange Version 2* (IKEv2)^[5] component of IPsec does not create state for a new *Security Association* until it can verify that initiators are capable of responding to packets sent to the address they claims to be using. There are other security protocols in which the listener sends out “puzzles” in response to initiation attempts and grants services or state only when puzzle solutions are returned^[6]. This tactic not only verifies the addresses of initiators but also implies a computational burden that causes them to further demonstrate their genuine willingness to communicate productively.

Countermeasures

During the initial Panix attack, random spoofed source addresses were being used, but it was noted that the attack TCP SYNs all used the same source port number. A filter that denied incoming packets from this port was temporarily effective, but easy for the attacker to adapt to, and the attack segments began using random ports. Panix was able to isolate which of its ingress routers the attack was coming from and null-route packets destined for its servers coming through that router, but this solution was obviously a heavy-handed one, and seems to have also been overcome when the attacker started sending packets that were routed through a different upstream provider. Panix had mixed success in getting its providers to assist in tracing and blocking the attack, and the networking community was spurred into devising other solutions.

Two broad classes of solutions to SYN flooding attacks have evolved, corresponding to where the defenses are implemented. The first class of solutions involves hardening the end-host TCP implementation itself, including altering the algorithms and data structures used for connection lookup and establishment, as well as some solutions that diverge from the TCP state machine behavior during connection establishment, as described in RFC 793.

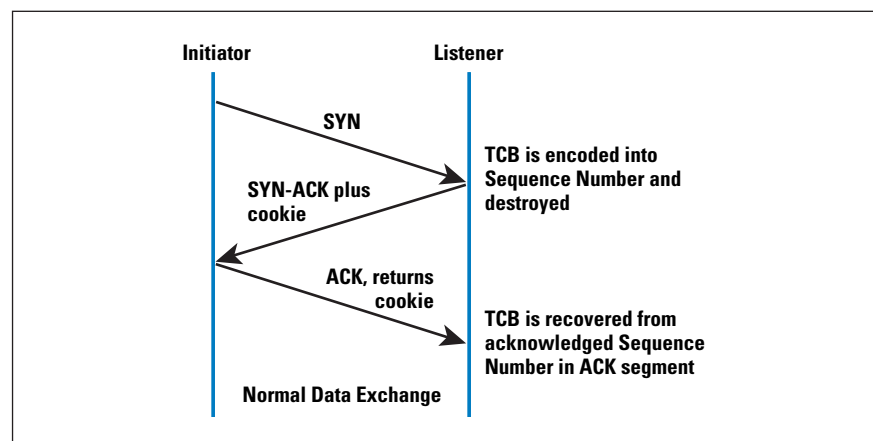
The second class involves hardening the network, either to lessen the likelihood of the attack preconditions (an army of controlled hosts or the propagation of IP packets with spoofed source addresses), or to insert middleboxes that can isolate servers on the networks behind them from illegitimate SYNs.

End-Host Countermeasures

Increasing TCP Backlog: Because the basic attack mechanism relies on overflowing a host's backlog of connecting sockets, an obvious end host-based solution is to simply increase the backlog, as is already done for very popular server applications. In at least some popular TCP implementations, this solution is known to be a poor one because of the use of linear list traversal in the functions that attempt to free state associated with stale connection attempts. Increasing the backlog is typically possible through altering the *listen()* call of an application and setting an operating system kernel parameter named `SOMAXCONN`, which sets an upper bound on the size of the backlog that an application can request. This step by itself should not be seriously considered as a means to defend against SYN flooding attacks—even in operating systems that can efficiently support large backlogs—because an attacker who can generate attack segments will most likely be able to scale to larger orders than the backlog supportable by a host.

Reducing the SYN-RECEIVED Timer: Another simple end host-based mechanism is to put a tighter limit on the amount of time between when a TCB enters the SYN-RECEIVED state and when it may be reaped for not advancing. The obvious disadvantage to this mechanism is that in cases of aggressive attacks that impose some amount of congestion loss in either the SYN-ACK or handshake-completing ACK packets, legitimate connection TCBs may be reaped as hosts are in the process of retransmitting these segments. Furthermore, there is only a linear relationship between the reduction that an administrator makes in the SYN-RECEIVED timer and the corresponding increase in packet rate that the adversary must make in order to continue attacking the server. Other alternative end-host solutions make it much more difficult for an attack to remain viable. For these reasons, a reduction in the SYN-RECEIVED timer is not an advisable defense against SYN flooding attacks.

Figure 4: Connection Establishment with SYN Cookies



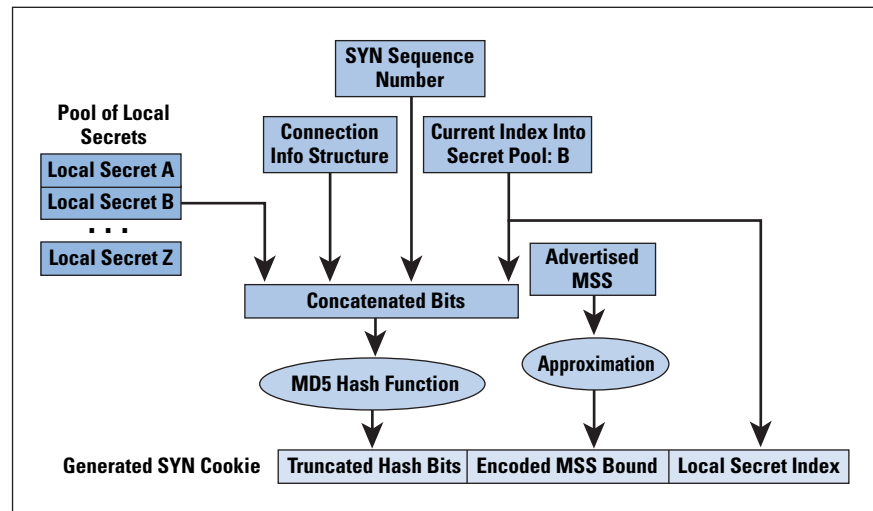
SYN Caches: Two end-host defenses, called SYN caches and SYN cookies (described later), operate by reducing the amount of state allocated initially for a TCB generated by a received SYN, and putting off instantiating the full state^[8]. In a host that uses a SYN cache, a hash table with a limited amount of space in each hash bucket is used to store a subset of the data that would normally go into an allocated TCB. If and when a handshake completing ACK is received, this data can be moved into a full TCB; otherwise the oldest bucket at a particular hash value can be reaped when needed. In Lemon's FreeBSD example^[8], the SYN cache entry for a half connection is 160 bytes, versus 736 bytes for a full TCB, and 15359 entries in the SYN cache are supported.

The SYN cache data structure is robust to attackers attempting to overflow its buckets because it uses the initiator's local port number and some secret bits in the hash value. Because stacks are a more effective data structure to search than a simple linked list, stacks that use a SYN cache can have improved speed, even when not under attack. Under Lemon's tests, during an active attack a host using a SYN cache was able to establish legitimate connections with only about a 15-percent increase in latency.

SYN Cookies: In contrast to the SYN cache approach, the SYN cookies technique causes absolutely zero state to be generated by a received SYN. Instead, the most basic data comprising the connection state is compressed into the bits of the sequence number used in the SYN-ACK. Since for a legitimate connection, an ACK segment will be received that echoes this sequence number (actually the sequence number plus one), the basic TCB data can be regenerated and a full TCB can safely be instantiated by decompressing the Acknowledgement field. This decompression can be effective even under heavy attack because there is no storage load whatsoever on the listener, only a computational load to encode data into the SYN-ACK sequence numbers. The downside is that not all TCB data can fit into the 32-bit Sequence Number field, so some TCP options required for high performance might be disabled. Another problem is that SYN-ACKs are not retransmitted (because retransmission would require state), altering the TCP synchronization procedures from RFC 793.

Recent work by Andre Oppermann uses the TCP Timestamp option in conjunction with the Sequence Number field to encode more state information and preserve the use of high-performance options such as TCP Window Scaling, and TCP *Selective Acknowledgment Options* (SACK), and can also be used to preserve *TCP-Message Digest 5* (MD5) support with SYN cookies. This option is a step forward, in that it removes the major negative effect of previous SYN cookie implementations that disabled these features.

Figure 5: Process for Generation and Validation of TCP SYN Cookies.



The exact format of TCP SYN cookies is not an interoperability issue, because they are only locally interpreted, and the format and procedures for generation and validation can vary slightly among implementations. Figure 5 depicts the general process of SYN cookie generation and validation used by multiple implementations.

To compute the SYN-ACK sequence number (that is, the TCP cookie) when using TCP cookies, a host first concatenates some local secret bits, a data structure that contains the IP addresses and TCP ports, the initial SYN sequence number, and some index data identifying the secret bits. An MD5 digest is computed over all these bytes, and some bits are truncated from the hash value to be placed in the SYN-ACK sequence number. Because the sequence number is about a fourth the size of the full hash value, this truncation is necessary, but generally at least 3 bytes worth of the hash bits are used, meaning that there should still be close to a 2^{24} effort required to guess a valid cookie without knowing the local secret bits. In addition to the hash output, some of the cookie bits indicate a lower bound on the *Maximum Segment Size* (MSS) that the SYN contained, and the index bits identifying the local secret used within the hash.

To validate a SYN cookie, first the acknowledgement number in an incoming ACK segment is decremented by 1 to retrieve the generated SYN cookie. The valid value for the set of truncated hash bits is computed based on the IP address pair, TCP port numbers, segment sequence number minus one, and the value from the secret pool corresponding to the index bits inside the cookie. If these computed hash bits match those within the ACK segment, then a TCB is initialized and the connection proceeds. The encoded MSS bound is used to set a reasonable-sized MSS that is no larger than what was originally advertised. This MSS is usually implemented as three bits whose code points correspond to eight “commonly advertised” MSS values based on typical link *Maximum Transmission Units* (MTUs) and header overheads.

Hybrid Approaches: A hybrid approach combines two or more of the single defense techniques described previously. For instance, some end-host operating systems implement both a large backlog and SYN cookies, but enable SYN cookies only when the amount of the backlog that is occupied exceeds some threshold, allowing them to normally operate without the disadvantages of SYN cookies, but also allowing them to fail over to the SYN-cookie behavior and be strongly protected when an attack occurs.

Network-Based Countermeasures

Filtering: The most basic network-level defense is application of the filtering techniques described in RFC 2827^[7]. Using ingress filtering, an ISP refuses to further route packets coming from an end site with IP source addresses that do not belong to that end site. Ingress filtering would be highly effective at preventing SYN flooding attacks that rely on spoofed IP packets. However, it is not currently reliable because ingress filtering policies are not universally deployed. Ingress filtering is also wholly ineffective against SYN flooding attacks that use a distributed army of controlled hosts that each directly attack. Ingress filtering is also a mechanism that an end site wishing to defend itself most often has no control over, because it has no influence upon the policies employed by ISPs around the world.

Firewalls and Proxies: A firewall or proxy machine inside the network can buffer end hosts from SYN flooding attacks through two methods, by either spoofing SYN-ACKs to the initiators or spoofing ACKs to the listener^[9].

Figure 6 shows the basic operation of a firewall/proxy that spoofs SYN-ACKs to the initiator. If the initiator is legitimate, the firewall/proxy sees an ACK and then sets up a connection between itself and the listener, spoofing the initiator's address. The firewall/proxy splits the end-to-end connection into two connections to and from itself. This splitting works as a defense against SYN flooding attacks, because the listener never sees SYNs from an attacker. As long as the firewall/proxy implements some TCP-based defense mechanism such as SYN cookies or a SYN cache, it can protect all the servers on the network behind it from SYN flooding attacks.

Figure 6: Packet Exchanges through a SYN-ACK spoofing Firewall/Proxy.

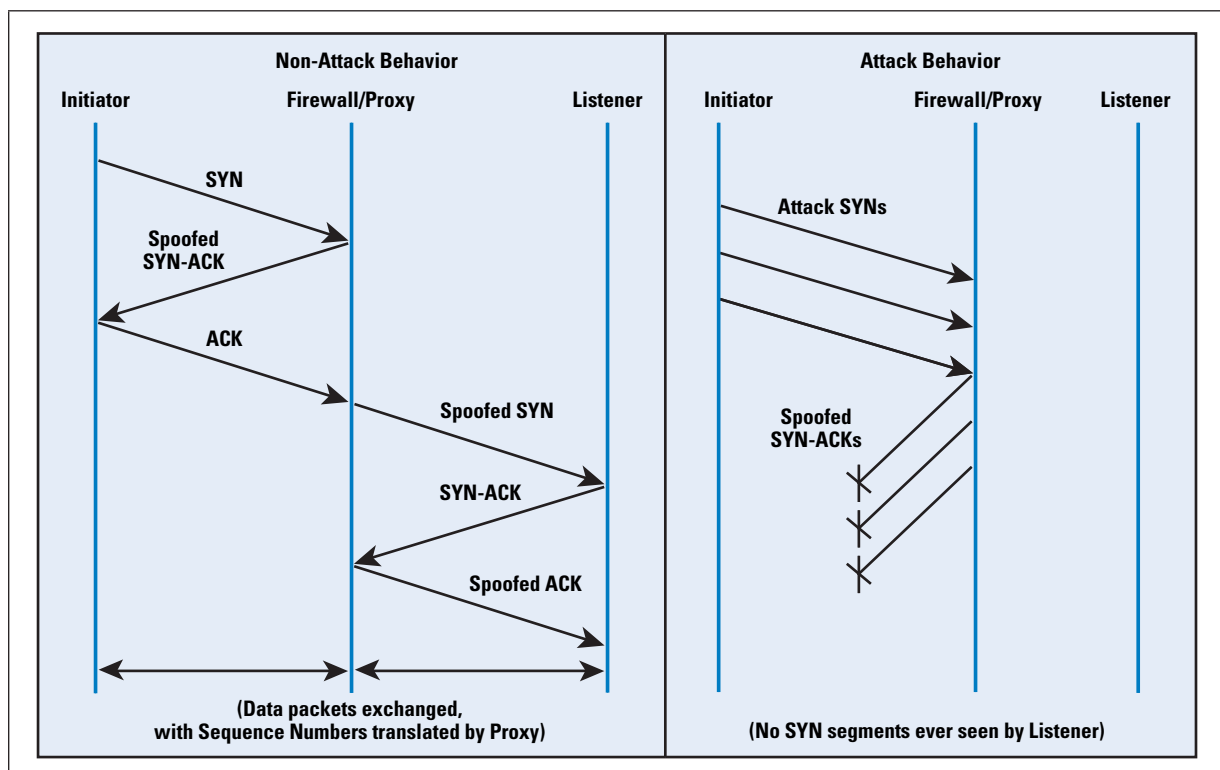
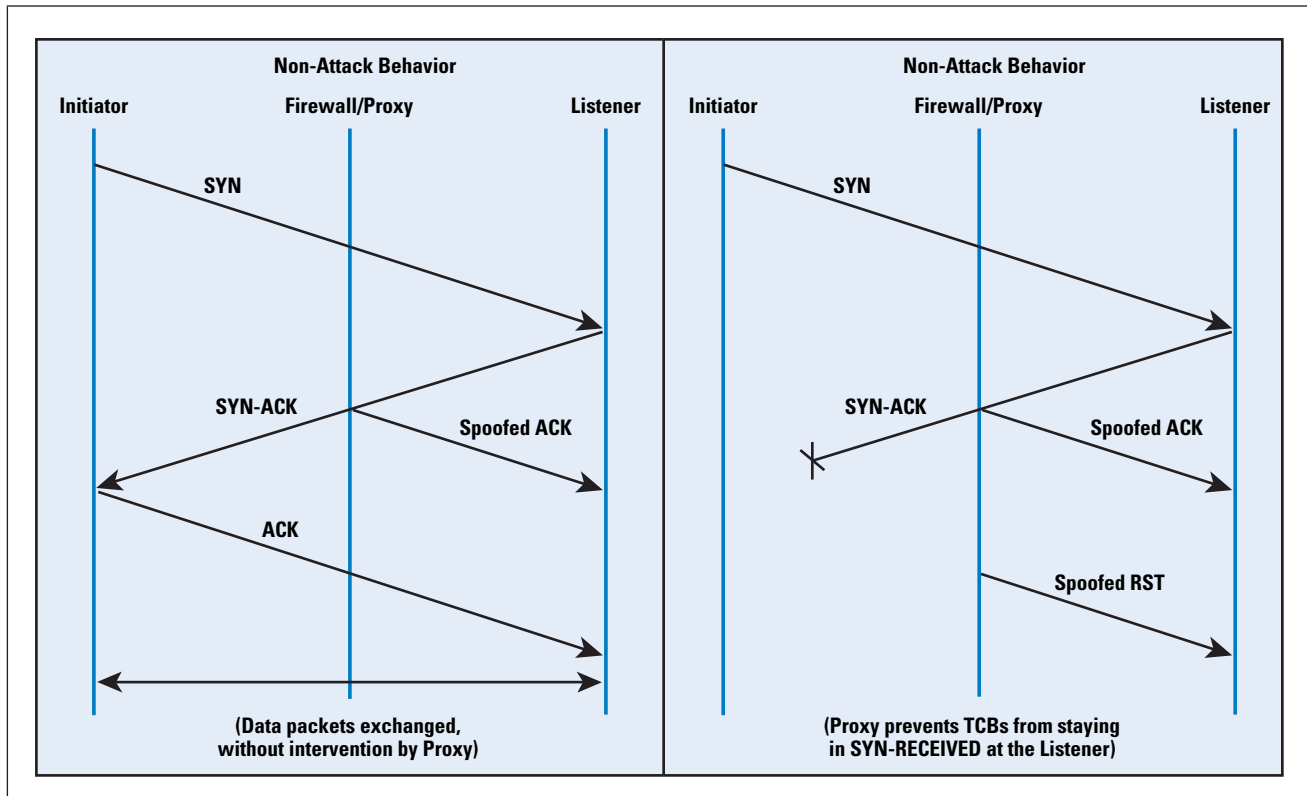


Figure 7 illustrates the packet exchanges through a firewall/proxy that spoofs ACKs to the listener in response to observed SYN-ACKs. This spoofing prevents the listeners TCBs from staying in the SYN-RECEIVED state, and thus maintains free space in the backlog. The firewall/proxy then waits for some time, and if a legitimate ACK from the initiator is not observed, then it can signal the listener to free the TCB using a spoofed TCP RST segment. For legitimate connections, packet flow can continue, with no interference from the firewall/proxy. This solution is more desirable than the mode of operation in Figure 5, where the firewall/proxy spoofs SYN-ACKs, because it does not require the firewall/proxy to actively participate in legitimate connections after they are established.

Figure 7: Packet Exchanges through an ACK-spoofing Firewall/Proxy.



Active Monitor: An active monitor is a device that can observe and inject traffic to the listener, but is not necessarily within the routing path itself, like a firewall is. One type of active monitor acts like the ACK-spoofing firewall/proxy of Figure 6, with the added capability of spoofing RSTs immediately if it sees SYNs from source addresses that it knows to be used by attackers^[9]. Active monitors are useful because they may be cheaper or easier to deploy than firewall-based or filtering solutions, and can still protect entire networks of listeners without requiring every listener's operating system to implement an end-host solution.

Defenses in Practice

Both end-host and network-based solutions to the SYN flooding attack have merits. Both types of defense are frequently employed, and they generally do not interfere when used in combination. Because SYN flooding targets end hosts rather than attempting to exhaust the network capacity, it seems logical that all end hosts should implement defenses, and that network-based techniques are an optional second line of defense that a site can employ.

End-host mechanisms are present in current versions of most common operating systems. Some implement SYN caches, others use SYN cookies after a threshold of backlog usage is crossed, and still others adapt the SYN-RECEIVED timer and number of retransmission attempts for SYN-ACKs.

Because some techniques are known to be ineffective (increasing backlogs and reducing the SYN-RECEIVED timer), these techniques should definitely not be relied upon. Based on experimentation and analysis (and the author's opinion), SYN caches seem like the best end-host mechanism available.

This choice is motivated by the facts that they are capable of withstanding heavy attacks, they are free from the negative effects of SYN cookies, and they do not need any heuristics for threshold setting as in many hybrid approaches.

Among network-based solutions, there does not seem to be any strong argument for SYN-ACK spoofing firewall/proxies. Because these spoofing proxies split the TCP connection, they may disable some high-performance or other TCP options, and there seems to be little advantage to this approach over ACK-spoofing firewall/proxies. Active monitors should be used when a firewall/proxy solution is administratively impossible or too expensive to deploy. Ingress and egress filtering is frequently done today (but not ubiquitous), and is a commonly accepted practice as part of being a good neighbor on the Internet. Because filtering does not cope with distributed networks of drones that use direct attacks, it needs to be supplemented with other mechanisms, and must not be relied upon by an end host.

Related Attacks

In addition to SYN flooding, several other attacks on TCP connections are possible by spoofing the IP source address and connection parameters for in-progress TCP connections^[10]. If an attacker can guess the two IP addresses, TCP port numbers, and a valid sequence number within the window, then a connection can be disrupted either through resetting it or injecting corrupt data. In addition to spoofed TCP segments, spoofed ICMP datagrams have the capability to terminate victim TCP connections.

Both these other attacks and SYN floods target a victim's TCP application and can potentially deny service to the victim using an attack rate less than that of brute-force packet flooding. However, SYN flooding and other TCP spoofing attacks have significant differences. SYN flooding denies service to new connections, without affecting in-progress connections, whereas other spoofing attacks disrupt in-progress connections, but do not prevent new connections from starting. SYN flooding attacks can be defended against by altering only the initial handshaking procedure, whereas other spoofing attacks require additional per-segment checks throughout the lifetime of a connection. The commonality between SYN flooding and other TCP spoofing attacks is that they are predicated on an attacker's ability to send IP packets with spoofed source addresses, and a similar defense against these attacks would be to remove this capability through more universal deployment of address filtering or IPsec.

Conclusion

At the time of this writing, the TCP SYN flooding vulnerability has been well-known for a decade. This article discussed several solutions aimed at making these attacks ineffective, some of which are readily available in commercial off-the-shelf products or free software, but no solution has been standardized as a part of TCP or middlebox function at the IETF level. The IETF's *TCP Maintenance and Minor Extensions* (TCPM) working group is in the process of producing an informational document that explains the positive and negative aspects of each of the common mitigation techniques^[10], and readers are encouraged to consult this document for further information.

In this author's opinion, some variant of the SYN cache technique should be a mandatory feature to look for in a server operating system, and the variant can be deployed in combination with other network-based methods (address-based filtering, ACK-spoofing firewalls, IPsec, etc.) in appropriate situations. It is encouraging to see that protocol designers have learned a lesson from the SYN flooding vulnerability in TCP and have made more recent protocols inherently robust to such attacks.

Acknowledgements

Several individual participants in the IETF's TCPM working group have contributed bits of data found in the group's informational document on SYN flooding^[11], some of which is replicated in spirit here.

References

- [1] daemon9, route, and infinity, "Project Neptune," *Phrack Magazine*, Volume 7, Issue 48, File 13 of 18, July 1996.
- [2] CERT, "CERT Advisory CA-1996-21 TCP SYN Flooding and IP Spoofing Attacks," September 1996.
- [3] Aura, T. and P. Nikander, "Stateless Connections," Proceedings of the First International Conference on Information and Communication Security, 1997.
- [4] Stewart, R., Xie, Q., Morneault, K., Sharp, C., Schwarzbauer, H., Taylor, T., Rytina, I., Kalla, M., Zhang, L., and V. Paxson, "Stream Control Transmission Protocol," RFC 2960, October 2000.
- [5] Kaufman, C., "Internet Key Exchange (IKEv2) Protocol," RFC 4306, December 2005.
- [6] Aura, T., Nikander, P., and J. Leiwo, "DOS-resistant Authentication with Client Puzzles," *Lecture Notes in Computer Science*, Volume 2133, revised from the 8th International Workshop on Security Protocols, 2000.

- [7] Ferguson, P. and D. Senie, “Network Ingress Filtering: Defeating Denial of Service Attacks which employ IP Source Address Spoofing,” BCP 38, RFC 2827, May 2000.
- [8] Lemon, J., “Resisting SYN Flood DoS Attacks with a SYN Cache,” BSDCON 2002, February 2002.
- [9] Schuba, C., Krsul, I., Kuhn, M., Spafford, E., Sundaram, A., and D. Zamboni, “Analysis of a Denial of Service Attack on TCP,” Proceedings of the 1997 IEEE Symposium on Security and Privacy, 1997.
- [10] Touch, J., “Defending TCP Against Spoofing Attacks,” Internet-Draft (work in progress), **draft-ietf-tcpm-tcp-antispoof-05**, October 2006.
- [11] Eddy, W., “TCP SYN Flooding Attacks and Common Mitigations,” Internet-Draft (work in progress), **draft-ietf-tcpm-syn-flood-00**, July 2006.

WESLEY M. EDDY works for Verizon Federal Network Systems as an onsite contractor at NASA’s Glenn Research Center, where he performs research, analysis, and development of network protocols and architectures for use in space exploration and aeronautical communications. E-mail: **weddy@grc.nasa.gov**

Boosting the SOA with XML Networking

by Silvano Da Ros

In the 1990s, the widespread adoption of *object-oriented programming* (OOP) and advancing network technologies fostered the development of distributed object technologies, including *Object Management Group's* (OMG's) *Common Object Request Broker Architecture* (CORBA) and Microsoft's *Distributed Common Object Model* (DCOM). Both CORBA and DCOM follow the OOP consumer-producer service model, where applications locally instantiate any number of objects and execute methods for the objects to obtain a service. However, with *distributed* object technologies, a local application can request a service from a remote application by instantiating a remote object and executing the methods of the object using *Remote Procedure Call* (RPC) over the network. The local application executes the methods of the remote object as if the object were an inherent part of the local application.

To push toward a simpler consumer-producer service model than distributed objects, the *Service-Oriented Architecture* (SOA) was created as a worldwide standards-based application interoperability initiative^[1]. SOA differs from distributed object technologies, because you no longer deal with object instantiation and method invocation to provide services between your applications^[2]. Instead, you can create *Extensible Markup Language* (XML)-based standard Web services to exchange XML documents between your applications using Internet-based application layer protocols, such as *Hyper Text Transfer Protocol* (HTTP) and the *Simple Mail Transfer Protocol* (SMTP).

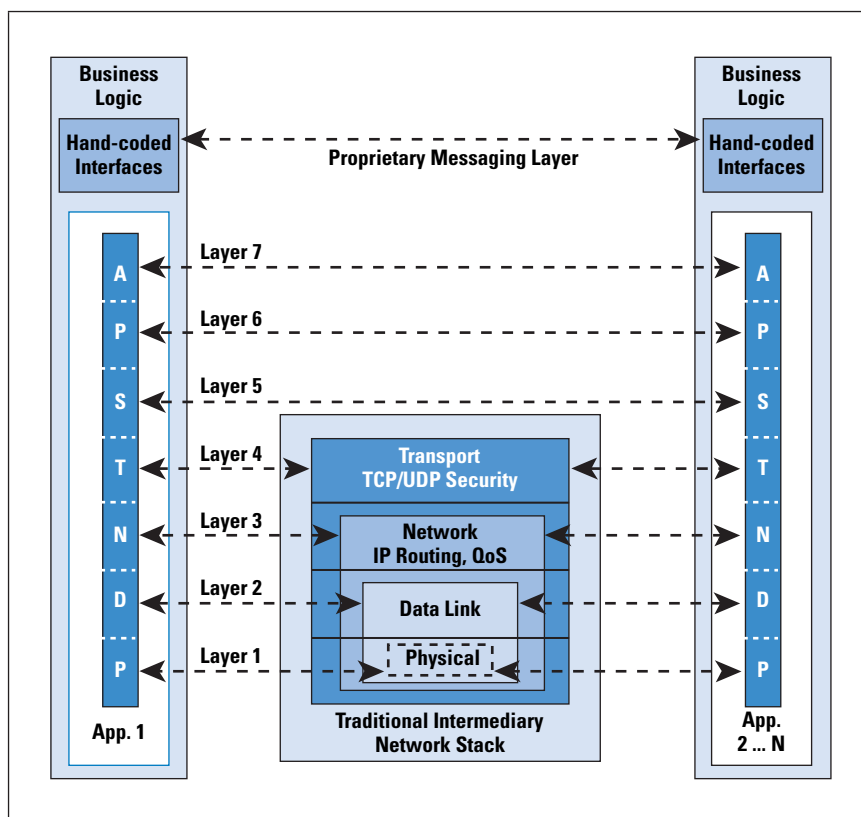
This is where XML networking comes into the picture. This article shows how to use SOA at the network edge, in conjunction with XML within the network, to help with the work required for enabling interoperability between your applications. The problem with SOA on its own is that to scale applications, hardware and software upgrades are required on the servers where your business logic resides. Because application integration using XML is CPU-intensive, it benefits from XML hardware built specifically for XML computations. However, the applications servers that run your business logic are effectively independent of the underlying XML processing. Therefore, to accelerate the SOA at the network level transparently to the application, XML networking technologies can be used. XML networking can provide SOA acceleration using a special middleware-enabled network layer, which this article explains. This special network layer also provides additional benefits to your applications that SOA alone cannot provide at the edge, such as dynamic message routing and transformation.

To help in the understanding of SOA acceleration with XML networking, the following section discusses SOA and its constituent technologies. Further sections explore the specifics of XML and XML-based network processing.

A Brief History of SOA

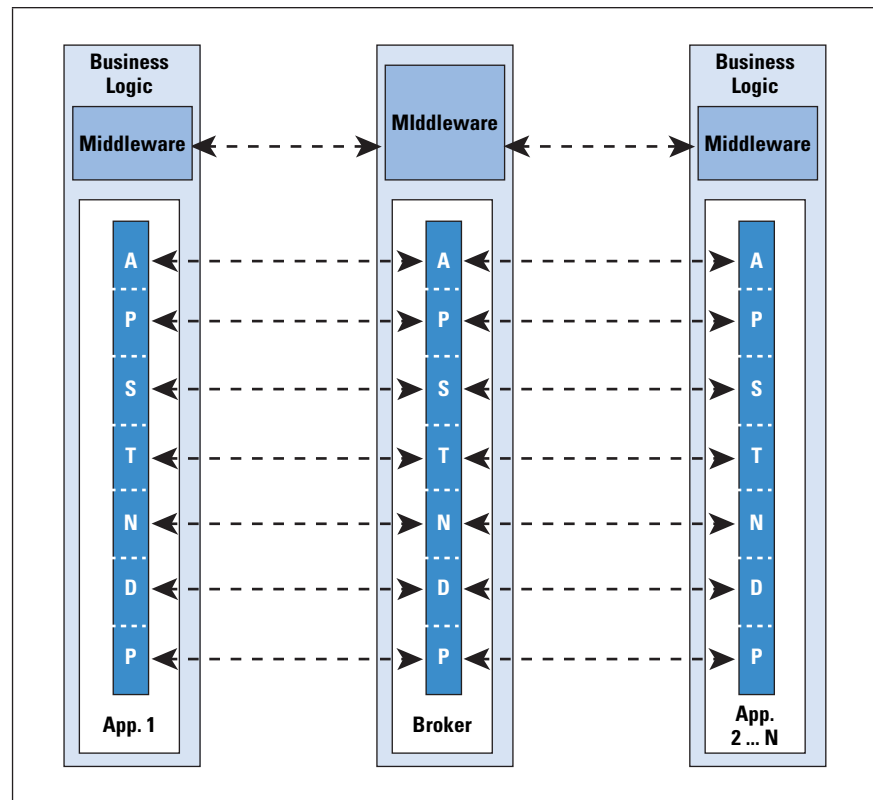
Traditionally, hand-coding proprietary interfaces were required to interoperate between your applications, as Figure 1 illustrates. This task is a trivial one if you have only a few applications, but if you have numerous disparate applications, all requiring interfaces into one another, the result is a complex, many-to-many web of connections between them. In the 1980s, *Electronic Data Interchange* (EDI) was developed to standardize the message formats that trading partners use to exchange text-based transaction data residing on mainframes, making it an early predecessor to SOA.

Figure 1: The Proprietary Messaging Layer



In the mid-1990s, standard middleware (or integration brokers) became available, such as CORBA and DCOM mentioned previously, to integrate advanced client-server applications. Figure 2 shows how integration brokers allow you to perform the translations between end systems over a standard messaging layer without creating application-specific interfaces between each system. During the same time, numerous software vendors, such as IBM WebSphere and TIBCO, also developed standard messaging layer protocols, which required adding vendor-specific adapters within the common integration brokers. Additionally, with newer application development environments being adopted, such as *Java 2 Sun Enterprise Edition* (J2EE) and Microsoft .NET, even more programming complexity is required when considering application interoperability without using the SOA. Fortunately, these new platforms currently support the SOA, allowing an application developed in one platform to tap into the data supplied by an application developed in the other.

Figure 2: The Standard Messaging Layer



Figures 1 and 2 illustrate how the proprietary and standard messaging layers sit above the network stack—at best, a traditional network device can operate only up to and including the transport layer. For example, by tracking TCP connection state information, a firewall device allows you to configure security services for your applications. Some firewalls can inspect the context of the application, but only to ensure the application behavior is RFC-compliant and not performing some sort of malicious activity. Additionally, at the next layer down the stack, you can configure Layer 3 *Quality of Service* (QoS) functions, such as *IP Precedence*, *Differentiated Services* (DiffServ), traffic shaping, and resource reservation, to ensure delivery of traffic to your critical applications. Although the network layers can provide these intelligent network services to your applications, they do not add any value toward accelerating your SOA.

Notice how the middleware portion in the proprietary messaging layer in Figure 1 takes up a larger portion of the application stack than the standard messaging layer from Figure 2. This situation occurs because the list of available messages that your standard messaging layer applications support is now much smaller—the broker takes care of the interfacing complexity on behalf of your applications. A reduced number of messages requires that you maintain much less middleware programming code on your applications than if every application in your network had to account for the messages of every other application.

Optimizing the SOA

Now that you understand SOA, you can better understand where XML networking fits into the scheme of things. Figure 3 illustrates how network equipment vendors can add specialized “application-aware” intelligence into Layers 5 through 7 of the OSI model.

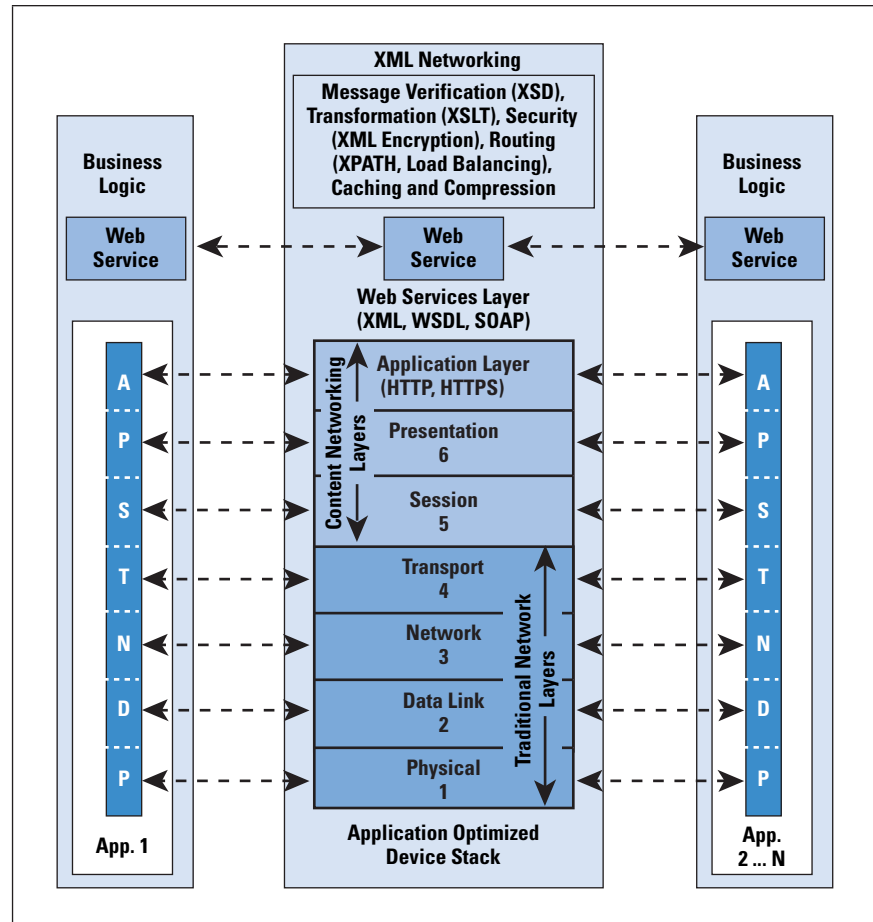
You can start by contrasting XML networking with traditional content networking technologies^[3]. As you can see in Figure 3, by incorporating content networking services into the network, such as *Server Load Balancing* (SLB), caching, and *Secure Sockets Layer* (SSL) acceleration, network vendors give you the ability to transparently accelerate your applications without the need of application hardware upgrades. However, by residing only within the OSI model, content networking services and protocols provide a “network-oriented” way to accelerate your applications. In order to achieve full application awareness, you must look not only into the application headers, but also into the application payload. Although the content networking protocols can inspect into the packet payload, they are meant for providing network layer services but *not* application integration services. For example, *Network-Based Application Recognition* (NBAR) allows you to mark the IP DiffServ field in packets containing high-priority application traffic by first detecting the behavior of the application. However, like the network layers, the content networking layers cannot fulfill SOA acceleration requirements either.

In contrast, XML networking provides integration services by inspecting the full context of the application transaction and adding XML standards-based intelligence on top of the TCP/IP stack. An XML-enabled network provides you greater control, flexibility, and efficiency for integrating your applications than integration brokers. Figure 3 shows how you can inspect the XML-based “Web services” layer to accelerate your applications developed within an SOA model without the need of an integration broker.

The most popular Web services protocol is *Simple Object Access Protocol* (SOAP)^[4]. With SOAP, your applications can request services from one another with XML-based requests and receive responses as data formatted with XML. Because SOAP uses XML, its Web services are self-descriptive and very simple to use.

You define your SOAP Web services with the XML-based *Web Services Description Language* (WSDL)^[5]. The WSDL binds the SOAP messages to the Web services layer, as discussed later in this article. You can then transport your SOAP messages over standard application layer protocols, such as HTTP and HTTPS, between your client-server applications.

Figure 3: The Web Services Layer

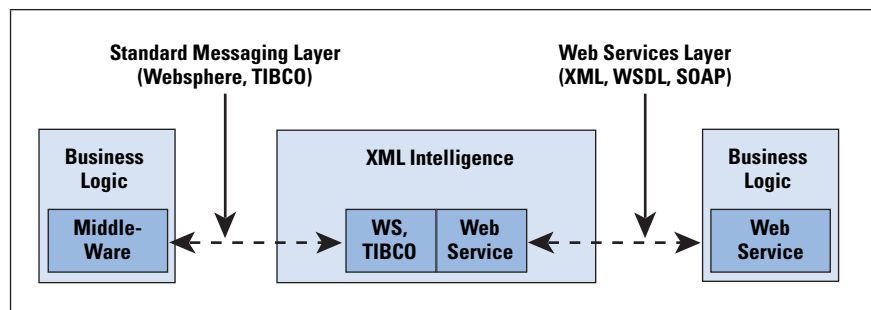


Similar to content networking technologies, XML networking can transparently add value to your applications within the network. When the XML network device receives the standard XML document from the Web services layer, you can configure the device to perform application-oriented services on the document. But because XML networking operates at the middleware layer and uses standard documents to integrate your applications, it provides you with fully standard functions using languages for:

- **Message Verification:** You can develop *World Wide Web Consortium (W3C) XML Schema Definitions (XSDs)* that your XML network can use to verify the syntax of your XML documents^[6].
- **XML Translation:** Using *XML Stylesheet Language Transformations (XSLT)*, you can translate XML documents to other non-XML formats, and conversely, directly within your network^[7].
- **Context-Based Routing:** Use *XML Path (XPath)* to route messages based on data stored within XML documents^[8]. A popular example is to route stock exchange quotes to a desired location when the value of the stock drops below a certain threshold.
- **High Availability:** Messages containing specific content can be load balanced across numerous identical origin servers.

- *Data Security*: You can accelerate XML encryption computations using either hardware or software XML-accelerated devices.
- *Compression and Caching*: You can cache frequently requested XML documents and compress XML documents to reduce network bandwidth. Like XML encryption, XML caching and compression can be performed using either hardware or software XML-accelerated devices.
- *Application-Layer Request Translation*: You can use XML networking to convert non-Web service requests into standard Web service requests. As with integration brokers, vendor-specific adapters are required to translate between WebSphere MQ, TIBCO, and SOA Web services. For example, Figure 4 shows how you can use network-level XML intelligence to translate between Websphere or TIBCO messages and Web services layer XML-based messages.

Figure 4: Intelligent Protocol Switching



Introducing XML Service Languages

Now that you have a general understanding of both SOA and XML networking, you can examine the specific XML technologies used for application interoperability, including:

- XML is used to format application data for storage and transmission.
- XSLT is used to translate between one XML format to another.
- XSD is used to describe, control, and verify an XML document format.
- XPath is a way to address items in an XML document hierarchy.
- SOAP is a messaging protocol used to encode information in Web service request and response messages before sending them over a network.

XML has its roots in the late 1960s from the *Generalized Markup Language* (GML), which was used to organize IBM's mainframe-based legal documents into a searchable form. The *Standard Generalized Markup Language* (SGML) was officially standardized in 1986 as an ISO international norm (ISO 8879). Since then, XML has become the predominant markup language for describing content. XML differs from HTML because it is not concerned with presenting or formatting content; instead, XML is used for describing the data using tags and attributes that you define yourself. Figure 5 is a sample XML document that organizes a police department's traffic ticket information.

Figure 5: A Traffic Ticket XML
Example

```
<?xml version="1.0"?>
<dept-tickets>
  <dept-chief>Greg Sanguinetti"/>
  <dept-id>12389289/>
  <ticket id="034567910" code="301">
    <offender>
      <name>John Smith</name>
      <license-number>10003887</license-number>
      <plate-number>9AER9876</plate-number>
    </offender>
    <offence-date>09/30/2005</offence-date>
    <location>
      <state>CA</state>
      <city>SJ</city>
      <intersection>West Tasman Dr.-Great America Pkwy.</intersection>
    </location>
    <officer>
      <officer-name>Paul Greene</officer-name>
      <officer-badge>7652323</officer-badge>
      <cruiser-plate-number>6TYX0923</cruiser-plate-number>
    </officer>
    <description>Failure to stop at red light</description>
    <fine>100</fine>
  </ticket>
  <ticket id="..." code="...">
    ...
  </ticket>
  <ticket id="..." code="...">
    ...
  </ticket>
</dept-tickets>
```

The XML in Figure 5 identifies the group of tickets for a police department by the department ID and the department chief's name. This example gives the data for a single traffic ticket as defined by the "ticket" element (or tag); however, you could include as many tickets as you want within the element "dept-tickets." The "ticket" element has two self-explanatory attributes (in dark blue), called "id" and "code," referring to the identification number for the individual ticket and the offense code, respectively. The sub-elements of the "ticket" element are also self-explanatory: "offender," "offence-date," "location," "officer," "description" and "fine."

In order to build a well-formed XML document, you must embrace the data for each element within its respective open and close tags (for example, <ticket>...</ticket>, or <ticket>.../>), properly nest all elements, and make sure all element names are in the proper case. You must also specify the XML version with the "<?xml version="1.0"?>" tag at the beginning of the XML document. HTML is less rigid than XML because it is case-insensitive and most Web browsers will allow you to leave out the close tags of an element. Because XML is very strict about following the correct syntax (that is, by making sure the XML is well-formed), XML processors are much easier to develop than traditional HTML Web browsers.

To verify that your documents are valid XML, you can check them against XSD files, which define the XML elements and their sequence, as discussed later in this article.

Transforming XML Using XSLT

You can use XSLT to translate one XML-based language into another. For example, you can translate standard XML into HTML. To translate the XML from Figure 5 into HTML for online viewing, you can use the XSLT file in Figure 6.

Figure 6: XSLT Translation – From XML to HTML

```
<?xml version="1.0"?>
<xsl:stylesheet version="1.0" xmlns:xsl="http://www.w3.org/1999/XSL/Transform">
  <xsl:template match="/">
    <html>
      <body>
        <br><b>Chief: </b><xsl:value-of select="dept-tickets/dept-chief"/></b>
        <br><b>Department No: </b><xsl:value-of select="dept-tickets/dept-id"/>
        </b>
        <table border="5">
          <!-- Output the HTML table headings -->
          <th>Ticket Number</th>
          <th>Offender's Name</th>
          <th>License Number</th>
          <th>State of Offense</th>
          <th>Officer's Name</th>
          <!-- Output the HTML table data -->
          <xsl:for-each select="dept-tickets/ticket">
            <tr>
              <td align="center"><xsl:value-of select="@id"/></td>
              <td align="left"><xsl:value-of select="offender/name"/></td>
              <td align="center"><xsl:value-of select="offender/license-number"/></td>
              <td align="center"><xsl:value-of select="location/state"/> </td>
              <td align="left"><xsl:value-of select="officer/officer-name"/></td>
              <td align="right">${<xsl:value-of select="fine"/></td>
            </tr>
          </xsl:for-each>
        </table>
      </body>
    </html>
  </xsl:template>
</xsl:stylesheet>
```

You must use a namespace to differentiate elements among the XML-based languages that you use in your XML document. As Figure 6 illustrates, the namespace is the string “xsl:”, which prefixes all of the XSLT elements. The particular application that parses the document (whether it is your XML device or a standalone XSLT parser^[9]) will know what to do with the specific elements based on the prefix. For example, an XSLT parser will look for the specific *Universal Resource Indicator* (URI) string constant that the W3C assigned to XSLT (that is, <http://www.w3.org/1999/XSL/Transform>) and perform the intended actions based on the elements in the document.

XML parsers do not use the URI of the namespace to retrieve a schema for the namespace—it is simply a unique identifier within the document. According to W3C, the definition of a namespace simply defines a two-part naming system (for example, “xslt:for-each”) and nothing else. After you define the namespace, the XML parser will understand the elements used within the document, such as “for-each” and “value-of” specified in Figure 6. For XSD documents, you must use a different namespace URI (that is, <http://www.w3.org/2001/XMLSchema>), as the next section discusses.

When you configure an XSLT parser or XML networking device to apply XSLT to an XML document, the parser starts at the top of the XSLT document by matching the root XML element within the source XML file. For example, the `<xslt:template match="/">` element in Figure 6 matches the “dept-tickets” root element from the XML file in Figure 5. The XSLT parser then creates the destination XML document (that is, a well-formed HTML file, in this example) and outputs the `<html>` and `<body>` tags to the new document. The XSLT parser then outputs the HTML table headers and loops through the XML document “ticket” elements, outputting selected items within the columns of the HTML table. The resulting HTML is given in Figure 7 for three sample tickets.

Figure 7: Resulting HTML Table – Source View

```
<html>
<body>
  <br><b>Chief: </b>Greg Sanguinetti<br><b>Department No: </b>12389289
  <table border="5">
    <th>Ticket Number</th><th>Offender's Name</th>
    <th>License Plate</th><th>State of Offense</th>
    <th>Officer's Name</th><th>Fine Amount</th>
    <tr>
      <td>034567910</td><td>John Smith</td><td>10003887</td>
      <td>CA</td><td>Paul Greene</td><td>100</td>
    </tr>
    <tr>
      <td>042562930</td><td>Gerald Rehnquist</td><td>11023342</td>
      <td>CA</td><td>Joel Patterson</td><td>200</td>
    </tr>
    <tr>
      <td>182736493</td><td>Jenny Barker</td><td>47281938</td>
      <td>CA</td><td>Emily Jones</td><td>120</td>
    </tr>
  </table>
</body>
</html>
```

Figure 8 illustrates the resultant HTML table that clients would see within a Web browser after the XSLT translation takes place.

Figure 8: Resulting HTML Table – Browser View

Ticket Number	Offender's Name	License Plate	State of Offense	Officer's Name	Fine Amount
034567910	John Smith	10003887	CA	Paul Greene	100
042562930	Gerald Rehnquist	11023342	CA	Joel Patterson	200
182736493	Jenny Barker	47281938	CA	Emily Jones	120

Verifying XML Using XSD

Because you can customize the structure and tags within an XML document, you should verify its syntax using XSDs. The XSD file in Figure 9 verifies the XML document given previously in Figure 5.

Figure 9: XSD File for Validating Traffic Ticket XML

```
<?xml version="1.0"?>
<xsd:schema xmlns:xsd="http://www.w3.org/2001/XMLSchema">
  <xsd:element name="dept-tickets">
    <xsd:complexType>
      <xsd:sequence>
        <xsd:element name="dept-chief"/>
        <xsd:element name="dept-id"/>
        <xsd:element name="ticket" maxOccurs="unbounded">
          <xsd:complexType>
            <xsd:sequence>
              <xsd:element name="offender">
                <xsd:complexType>
                  <xsd:sequence>
                    <xsd:element name="name"/>
                    <xsd:element name="license-number"/>
                    <xsd:element name="plate-number"/>
                  </xsd:sequence>
                </xsd:complexType>
              </xsd:element>
              <xsd:element name="offence-date"/>
              <xsd:element name="location">
                ...
              </xsd:element>
              <xsd:element name="officer">
                ...
              </xsd:element>
              <xsd:element name="description"/>
              <xsd:element name="fine"/>
            </xsd:sequence>
            <xsd:attribute name="id"/>
            <xsd:attribute name="code"/>
          </xsd:complexType>
        </xsd:element>
      </xsd:sequence>
    </xsd:complexType>
  </xsd:element>
</xsd:schema>
```

You must define an XSD namespace with the URI “<http://www.w3.org/2001/XMLSchema>” and prefix all the XSD elements that you use in the XSD file, such as “element,” “complex-type,” and “attribute,” with this namespace. At the top of your XSD file, you must specify the root XML element; the remaining elements within your XML document can be defined within the root element. Using the “complex-type” XSD element, you can specify elements that contain child elements (in contrast, “simple-type” indicates that the element does not contain any child elements). In this example, the “dept-tickets” element may contain a sequence of one or more child elements (as represented by the <xsd:sequence> element), including “dept-chief,” “dept-id,” and any number of element “ticket.”

Routing Messages Using XPATH

XPATH was developed primarily to be used with XSLT to transform the XML tags within an XML document based on the path of the data. Previously, in Figure 6, you saw how to select the entire list of tickets using the XSLT “select” attribute:

xsl:value-of select="dept-tickets/ticket"

However, within an XML network, you can also use XPATH to search within an XML document to route XML messages based on the values of the document data. For example, a state government may need the headquarters police department to route unpaid tickets that are within a tolerable threshold amount to the motor vehicle department for processing—there, the driver’s license can be suspended until the ticket is paid. However, those unpaid tickets that exceed a maximum threshold amount must be routed to the court service government department for processing. The court may decide to press further charges, depending on the driver’s previous driving record. Additionally, severe infractions, such as drunken or reckless driving, must be routed automatically to the court, regardless of whether the ticket is paid or not. The XPATH expression “dept-tickets/ticket” given previously returns the entire list of traffic tickets. Alternatively, if you want only the unpaid tickets with a fine value of greater than \$100, you could use the XPATH expression:

dept-tickets/ticket[@paid='no' and fine>100]

The XPATH symbol “@” here indicates that an attribute is being selected, and not an element. To select tickets with codes 309 and 310 (that is, fictitious codes for severe infractions), you can use the following XPATH expression:

dept-tickets/ticket[@code=309 or @code=310]

Using SOAP Web Services

SOAP provides a standard way to send transaction information over TCP/IP application protocols, such as HTTP. For example, you could create a SOAP request-response operation over HTTP for exchanging traffic ticket information between two applications. As Figure 10 illustrates, the requesting client application sends a “getFineRequest” message to the server, which in turn responds with the appropriate fine amount within a “getFineResponse” message.

Figure 10: A Sample SOAP Request-Response Operation

```

Client Request :
POST /getticketfine HTTP/1.1
Host: www.example.com
Content-Type: application/soap+xml;

<?xml version="1.0"?>
<soap:envelope
  xmlns:soap="http://www.w3.org/2001/12/soap-envelope"
  soap:encodingStyle="http://www.w3.org/2001/12/soap-encoding">
  <soap:body>
    <tn:getFineRequest xmlns:tn="http://example.com/getticketfine">
      <tn:ticket-id>034567910</tn:ticket-id>
    </tn:getFineRequest>
  </soap:body>
</soap:envelope>

Server Response:
HTTP/1.1 200 OK
Content-Type: application/soap+xml;

<?xml version="1.0"?>
<soap:envelope
  xmlns:soap="http://www.w3.org/2001/12/soap-envelope"
  soap:encodingStyle="http://www.w3.org/2001/12/soap-encoding">
  <soap:body>
    <tf:getFineResponse xmlns:tf="http://example.com/getticketfine">
      <tf:fine>100</tf:fine>
    </tf:getFineResponse>
  </soap:body>
</soap:envelope>

```

You encapsulate each SOAP message within the “Envelope” SOAP element. Within Envelope, you need to prefix the SOAP elements with the SOAP namespace, called “soap:” in this example, which you define as an attribute within Envelope. The “encodingStyle” attribute of the Envelope element defines the data types in the SOAP document. You must also define a custom namespace (that is, “tf,” which stands for “ticket-fine”), with which you prefix all the application-specific elements.

To define the structure of the SOAP Web service running within your applications, you can use WSDL, which you develop so that your clients know the exact specification of the services that they can request, the types of responses they should expect to receive, and the protocols (for example, SOAP or HTTP) with which they should send messages.

For example, you can publish the WSDL to your clients, who may not be aware of the messages available within your Web services layer. The clients can retrieve the WSDL file and send the appropriate SOAP messages to the SOAP Web service running on your application. To publish the WSDL file to your clients, you can use a publicly available *Universal Description, Discovery and Integration* (UDDI) registry, such as XMethods^[10], or you could create your own UDDI registry^[11].

WSDL uses XSD to define your SOAP application data types. For example, for one application to request a fine amount (of XSD type `xs:integer`) for a given ticket ID (of XSD type `xs:string`) from your SOAP Web service called “ticketFineService,” you could use the WSDL in Figure 11.

Figure 11: WSDL for SOAP Request-Response Operation

```
<?xml version="1.0"?>
<definitions name="TicketInfo"
  targetNamespace="http://example.com/ticketinfo.wsdl"
  xmlns:tns="http://example.com/ticketinfo.wsdl"
  xmlns:soap="http://schemas.xmlsoap.org/wsdl/soap/"
  xmlns="http://schemas.xmlsoap.org/wsdl/">

  <message name="getFineRequest">
    <part name="ticket-id" type="xs:string"/>
  </message>

  <message name="getFineResponse">
    <part name="value" type="xs:integer"/>
  </message>

  <porttype name="ticketFine">
    <operation name="getTicketFine">
      <input message="tns:getFineRequest"/>
      <output message="tns:getFineResponse"/>
    </operation>
  </porttype>

  <binding name="ticketBinding" type="ticketFine">
    <soap:binding transport="http://schemas.xmlsoap.org/soap/http"/>
    <operation name="getTicketFine">
      <soap:operation soapAction="getTicketFine"/>
      <input>
        <soap:body use="encoded"/>
      </input>
      <output>
        <soap:body use="encoded"/>
      </output>
    </operation>
  </binding>

  <service name="ticketFineService">
    <documentation>WSDL File for ticketFineService</documentation>
    <port name="ticketFine" binding="ticketBinding">
      <soap:address location="http://example.com/getticketfine"/>
    </port>
  </service>
</definitions>
```

You start your WSDL file by declaring all the required namespaces. In order for the WSDL file to refer to element names that are defined within the same file (for example, “tns:getFineRequest” within the “porttype” element), you must use the “targetNamespace” element to define a custom URI that your custom namespace uses (that is, “tns,” meaning “this name space”).

You define the WSDL namespace for SOAP-specific elements with *xmlns:soap=http://schemas.xmlsoap.org/wsdl/soap*. For WSDL-only elements, you can use the default namespace *xmlns= http://schemas.xmlsoap.org/wsdl/*. Note that elements within the file that do not have a prefix use the default namespace.

After you create the namespaces for the WSDL file, you can then create the two messages for the transaction, “getFineRequest” and “getFineResponse,” using WSDL “message” elements. WSDL ports create the request-response transaction flow using the “operation” element, by specifying which message is the request (input) and which is the response (output). After you define the transaction, you must bind it to SOAP with WSDL using the WSDL “binding” element. Additionally, to set the transport to HTTP, you must use the “binding” SOAP-specific element. You then link the operation you created previously within the WSDL “port-type” element to SOAP using the “operation” subelement within the parent “binding” element.

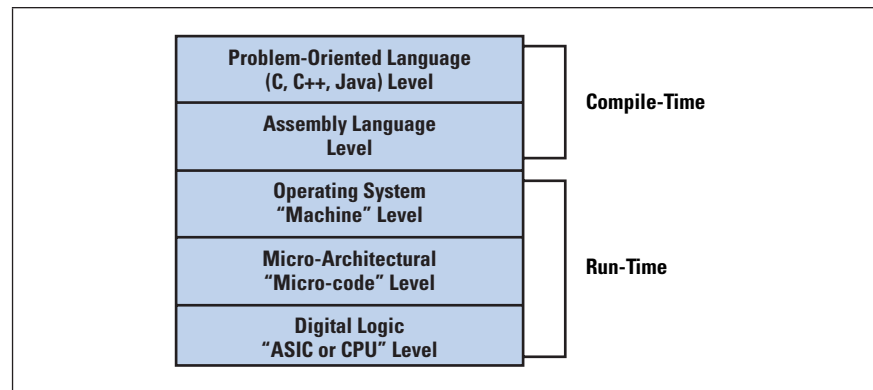
If you set the “use” element to “encoding,” you do not need to use an XSD “type” attribute for defining SOAP data types in your SOAP messages. However, you must specify the “encodingStyle” URI to **http://www.w3.org/2001/12/soap-encoding**, as you learned previously in Figure 10. Otherwise, if you set “use” to “literal,” then you would need to use the *type="xsd:string"* attribute in Figure 10 within the “tn:ticket-id” element when sending a request.

To define the SOAP Web service, you must use the WSDL “service” element. The SOAP element “address” within this element is the location where SOAP clients can send the “getTicketFine” requests, as Figure 10 illustrates.

Hardware vs. Software XML Acceleration

To help you understand the difference between hardware and software XML acceleration, Figure 12 illustrates a typical multilevel computer architecture^[12]. The highest level is where you would typically program your applications. When you compile your application, the compiler would typically “assemble” your various objects into assembly language prior to generating the machine-level code. This machine code is what is normally stored in an “.exe” file, which only your operating system can understand. When you execute the “.exe” at run time, the operating system converts the machine-level code into microcode, which the digital logic level within the CPU hardware can execute directly.

Figure 12: A Multilevel Computer Architecture



Examples of software-based network applications in the past that have transitioned to hardware acceleration include IP routing, encryption, firewalling, caching, and load balancing. XML networking is also a recent candidate for hardware acceleration; it is available by XML vendors that use XML *Application-Specific Integrated Circuits* (ASICs) or *Field Programmable Gate Arrays* (FPGAs) in their products^[13]. By programming the digital logic layer with the necessary circuits to perform intensive XML computations such as XSLT transformation, XML encryption, and XML schema validation, you can drastically increase the performance of the hardware platform.

However, some vendors have also found clever ways of accelerating XML computations on general-purpose hardware. Accelerating XML in software requires bypassing the additional machine-level step at run time. By “compiling” XML-based language instructions directly into microcode at the micro-architectural level, you can introduce XML computations to the underlying hardware directly at run time. That is, executing XML microcode at the digital logic level bypasses additional processing at the operating system “machine” level.

Summary

When a technology matures as a software agent running within an application, the need often arises to move the agent’s functions to the network. Indeed, this was the case with numerous software-based technologies of the past, such as IP routing, encryption, stateful firewall filtration, and server load balancing.

To facilitate the interoperability of diverse applications, SOA was developed as a prescription to complexity problems faced by commonly used distributed-object technologies. As SOA matures, the need to introduce XML-based functions to the network will grow. In order to streamline the responsibilities of an SOA-based application, you can transition your XML technologies, such as XML translation, validation, and security, from within the application to an XML-enabled network.

For Further Reading

- [1] Hao He, “What Is Service-Oriented Architecture?” O’Reilly, September 30, 2003,
<http://webservices.xml.com/pub/a/ws/2003/09/30/soa.html>
- [2] Werner Vogels, “Web Services Are Not Distributed Objects,” *Computing in Science and Engineering*, November-December 2003, <http://computer.org/internet/>
- [3] Christophe Deleuze, “Content Networks,” *The Internet Protocol Journal*, Volume 7, Number 2, June 2004.
- [4] “SOAP Version 1.2 Part 0: Primer,” W3C Recommendation, 24 June 2003,
<http://www.w3.org/TR/2003/REC-soap12-part0-20030624/>
- [5] “Web Services Description Language (WSDL) Version 2.0 Part 0: Primer,” W3C Working Draft, 3 August 2005,
<http://www.w3.org/TR/2005/WD-wsdl20-primer-20050803/>
- [6] “XML Schema Part 0: Primer Second Edition,” W3C Recommendation, 28 October 2004,
<http://www.w3.org/TR/xmlschema-0/>
- [7] “XSL Transformations (XSLT), Version 2.0,” W3C Recommendation, 16 November 1999,
<http://www.w3.org/TR/xslt>
- [8] “XML Path Language (XPath) Version 1,” W3C Recommendation, 16 November 1999, <http://www.w3.org/TR/xpath>
- [9] www.xmlspy.com
- [10] www.xmethods.net
- [11] www.uddi.org
- [12] Andrew S. Tanenbaum, *Structured Computer Organization*, (5th edition), ISBN 978-0131485211, Prentice Hall, 2005.
- [13] Michael John Sebastian Smith, *Application-Specific Integrated Circuits*, ISBN 978-0201500226, Addison-Wesley, June 1997.

SILVANO DA ROS currently works as a networking consultant in Toronto and has worked previously as a Systems Engineer for Cisco Systems. He is the author of *Content Networking Fundamentals*, published by Cisco Press. He holds a Bachelor of Computer Science and a Masters of Engineering (in Internetworking) from Dalhousie University. E-mail: sdaros@sympatico.ca

Time to Live

As I read the very fine article entitled “IPv6 Internals” (IPJ Volume 9, No. 3, September 2006), I was prompted to review the history of the *Time to Live* (TTL) as discussed in section 5.3.1 of RFC 1812. Being gray of head, little facts from other eras come quickly to mind. The *Xerox Network Systems* (XNS) Internet Transport on which Novell Netware was based required that no router ever store a packet in queue longer than 6 seconds. Requirements of RFC 791 were also softened in RFC 1812; rather than *requiring* the TTL to be decremented at least once and additionally once per second in queue, that document requires that the TTL be treated as a hop count and—reluctantly—reduces the treatment of TTL as a measure of time to a suggestion.

The reason for the change is the increasing implementation of higher-speed lines. A 1,500-byte datagram occupies 12,000 bits (and an asynchronous line sends those as 15,000 bits), which at any line speed below 19.2 kbps approximates or exceeds 1 second per datagram. Any time there are several datagrams in queue, the last message in the queue is likely to sit for many seconds, a situation that in turn can affect the behavior of TCP and other transports. However, 56-kbps lines became common in the 1980s, and T1 and T3 lines became common in the 1990s. Today, hotels generally offer Ethernet to the room; we have reports of edge networks connected to the Internet at 2.5 Gbps, and residential broadband in Japan and Europe at 26 Mbps per household. At 56 kbps, a standing queue of five messages is required to insert a 1-second delay, and at T1 it requires a queue depth of more than 100 messages. At higher speeds, the issue becomes less important.

That is not to say that multisecond queues are now irrelevant. Although few networks are being built today by concatenating asynchronous links, in developing countries—and on occasion even in hotels here in Santa Barbara, California—people still use dialup lines. In Uganda, some networks that run over the instant messaging capacity of GSM [*Global System for Mobile Communications*], which is to say using 9,600-bps datagrams, have been installed under the supervision of Daniel Stern and UConnect.org (www.uconnect.org). Much of the world still measures *round-trip times* (RTTs) in seconds, and bit rates in tens of kbps.

The TCP research community, one member of which recently asked me whether it was necessary to test TCP capabilities below 2 Mbps, and the IETF community in general would do well to remember that the ubiquity of high bandwidth in Europe, North America, Australia, and Eastern Asia in no sense implies that it is available throughout the world, or that satellite communications and other long-delay pipelines can now be ignored.

—Fred Baker, Cisco Systems
fred@cisco.com

The author responds:

Although to the casual observer the evolution of the Internet seems one of continuously increasing speed and capacity, reality is slightly different. The original ARPANET used 50-kbps modems in the late 1960s. In the next three decennia or so, the maximum bandwidth of a single link increased by a factor 200,000 to 10 Gbps. Interestingly enough, the minimum speed used for Internet connections went down to a little under 10 kbps, so where once the ARPANET had a uniform link speed throughout the network, the difference between the slowest and the fastest links is now six orders of magnitude. The speed difference between a snail and a supersonic fighter jet is only five orders of magnitude. Amazingly, the core protocols of the Internet—IP and TCP—can work across this full speed or bandwidth gamut, although changes were made to TCP to handle both extremes better, most notably in RFCs 1144 and 1323.

Even though I don't think keeping track of the time that packets are stored in buffers, as suggested in the original IPv4 specification, makes much sense even in slow parts of the network, Fred makes a good point: many Internet users still have to deal with speeds at the low end of the range; some of us only occasionally when connecting through a cellular network, others on a more regular basis. Even in Europe and the United States many millions of Internet users connect through dialup. For someone who is used to having always-on multimegabit connectivity, going back to 56 kbps or worse, 9,600 bps can be a bizarre experience. Many of today's Websites are so large that they take minutes to load at this speed. Connecting to my mail server using the *Internet Mail Access Protocol* (IMAP) takes 15 minutes. And one of my favorite relatively new applications, podcasting, becomes completely unusable: downloading a 50-minute audio program takes hours at modem speeds.

And that's all IPv4. It is possible to transport IPv6 packets over the *Point-to-Point Protocol* (PPP) that is used for almost all low-speed connections, but in practice this isn't workable because there are no provisions for receiving a dynamic address from an ISP [*Internet Service Provider*]. With IPv4, Van Jacobson did important work to optimize TCP/IP for low-speed links (RFC 1144). By reducing the *Maximum Transmission Unit* (MTU) of the slow link and compressing the IP and TCP headers, it was possible to achieve good interactive response times by avoiding the situation where a small packet gets stuck between a large packet that may take a second or more to transmit over a slow link while at the same time reducing the header overhead. Although the IETF has later done work on IPv6 header compression, it doesn't look like anyone has bothered to implement these techniques, and the minimum MTU of 1,280 bytes creates significant head-of-line blocking when IPv6 is used over slow links.

Another example where low bandwidth considerations are ignored is the widespread practice of enabling RFC 1323 TCP high-performance extensions for all TCP sessions. RFC 1323 includes two mechanisms: a window scale factor that allows much larger windows in order to attain maximum performance over high-bandwidth links with a long delay, and a timestamp option in the TCP header that allows for much more precise round-trip time estimations. With these options enabled, every TCP segment includes 8 extra bytes with timestamp information. In addition to increasing overhead, the timestamp option introduces an unpredictable value into the TCP header that makes it impossible to use header compression, thereby negating the usefulness of RFC 1144. To add insult to injury, almost no applications allocate enough buffer space to actually use the RFC 1323 mechanisms.

Moral of the story for protocol designers and implementers: spend some time thinking about how your protocol works over slow links. You never know when you'll find yourself behind just such a link.

—Iljitsch van Beijnum
iljitsch@muada.com

Gigabit TCP and MTU Size

I appreciated Geoff Huston's thorough description about the current obstacles and research involving Gigabit TCP (IPJ, Volume 9, No. 3, June 2006). I have already shown the article to many of my colleagues. It appears that Geoff did not address one of the solutions, which is to increase the networkwide *Maximum Transmission Unit* (MTU). In theory that would allow the existing TCP congestion control to handle higher-speed connectivity. Perhaps he did not address the issue because it is infeasible to increase the MTU setting Internetwide, especially with 10-Gigabit Ethernet interfaces sporting a default MTU setting of 1,500 bytes. On the other hand, projects that own their own backbone infrastructure may find increasing the default MTU a feasible approach.

For more information about raising the MTU, please see:
<http://www.psc.edu/~mathis/MTU/>

—Todd Hansen, UCSD/SDSC
tshansen@hpwren.ucsd.edu

The author responds:

Yes, it's true that increasing the size of the packet makes sound sense when the available bandwidth has increased. If the bandwidth increases by one order of magnitude and the packet size is increased by the same amount, then it is theoretically possible to effectively increase the throughput of the system without changing the packet processing load.

Effectively, if you regard the protocol interaction as a time sequence, then a coupling of increased bandwidth and comparably increased packet size preserves the time sequence interaction. Of course, as bandwidth on the network has increased we have not seen a comparable increase in MTU sizes, and today's networks exhibit a wide variety of MTUs and the importance of *Path MTU Discovery*, and coherent transmission of related MTU ICMP [*Internet Control Message Protocol*] messages becomes more critical as a consequence. Although the article concentrated on modifications to the TCP control algorithm, there is no doubting the importance of high-speed TCP senders and receivers using large TCP buffers to maximize the payload throughput potential.

—*Geoff Huston, APNIC*
gih@apnic.net

Drop us a Line!

We welcome any suggestions, comments or questions you may have regarding anything you read in this journal. Send us an e-mail to **ipj@cisco.com**. Also, don't forget to let us know if your delivery address changes. You can use the online subscription system to change your own information by supplying your Subscription ID and e-mail address. The system will then send you an e-mail with a "magic" URL which will allow you to update your database record. If you don't have your Subscription ID or encounter any difficulties, just send us the updated information via e-mail.

—*Ole J. Jacobsen, Editor and Publisher*
ole@cisco.com

Book Review

Internet Measurement

Internet Measurement: Infrastructure, Traffic & Applications, by Mark Crovella, Balachander Krishnamurthy, ISBN 0-470-01461-X, Wiley, 2006.

This book is a comprehensive reference guide to about 900 journal, conference, and workshop papers, and RFCs on the important and rapidly advancing field of Internet measurement. Interest in this growing field arises for three major reasons: commercial, social, and technical. Readers need nothing more than a keen interest in a methodical study of the subject matter from either a practical or research perspective to glean something from this book.

Organization

The book is centered on three architectural pillars relevant to measurement: *infrastructure*, *traffic*, and *applications*. Within each of these pillars, the topics are organized into four sections: *properties*, *challenges*, *tools*, and *state-of-the-art*. In the properties section, the authors review metrics that are important to measure in each area. In the challenges section, they discuss various difficulties and limitations that arise when trying to measure the metrics. The tools section covers some of the popular methods and products used to measure these metrics and work around the challenges mentioned previously. The intent is not to provide “user guides” for these tools. The state-of-the-art section presents the latest measurement results about covered properties and metrics, noting that they are subject to relatively fast obsolescence because of the rapidly evolving Internet.

The first three chapters provide background material. The first chapter provides an obligatory introduction to the Internet architecture, including how the “end-to-end” principle has been used for nearly 20 years to guide many design decisions in the Internet. The second chapter provides the analytic background necessary to study the Internet and cast its measurements in quantitative terms. The third chapter examines the nuts and bolts of Internet measurement, addressing the practical topics to consider in designing and implementing them, including the role of time and its sources.

The second part of this book also consists of three chapters, which cover the three pillars in depth. The first chapter defines metrics of interest for measuring the Internet and describes some of the barriers to their measurement, in particular “middleboxes,” *Network Address Translators* (NATs), firewalls, and proxies that deviate from the end-to-end architecture principle, may block *User Datagram Protocol* (UDP) or *Internet Control Message Protocol* (ICMP) packets, or hinder visibility to endpoint IP addresses. The authors next explore various tools and methods for active and passive measurement, estimation, and inference of these metrics.

Readers may wonder why two important metrics are left out—router reliability and high availability—where *Open Shortest Path First* (OSPF) and the *Border Gateway Protocol* (BGP) “Graceful Restart” would be of interest.

The next chapter focuses on traffic properties that are important to understand, measure, and model. The authors examine the challenges in capturing, processing, storing, and managing large volumes of packets and flows, as well as those related to their statistical characterization. Readers engaged in data modeling and performance analysis will benefit from this chapter. The last chapter in this part of the book examines some popular applications: The *Domain Name System* (DNS), Web, and *Peer-to-Peer* (P2P). The authors discuss the shifts in application mix from the 1980s, when FTP was dominant, to the 1990s, when the *Hypertext Transfer Protocol* (HTTP) became dominant, to today, when by most accounts P2P is the dominant Internet protocol. Next, there is a thorough coverage of the what (properties), why (justification), and how to (tools) facets of measurement of the three popular applications, as well as some coverage of online games and streaming media.

The third part of the book covers material that spans multiple areas. Its first chapter deals with anonymization of collected measurement data, which arises because of the need for data sharing, while preserving identity-related, personal-sensitive, or business-sensitive information for applications previously examined. The second chapter provides a short—but important—coverage of the key areas where Internet measurement has played a role in security enforcement. Various attack types and tools to combat them are discussed. The third chapter examines numerous low-level monitoring tools for high-speed traffic capture, as well as an insightful look at the software architecture of two toolsets, *dss* and *Gigascope*, reflecting the experience of one of the authors at AT&T Labs with them. It also reviews some large-scale measurement platforms at the *Cooperative Association for Internet Data Analysis* (CAIDA), the *Réseaux IP Européens* (RIPE) community, and the *High Energy Physics* (HEP) community. The book concludes with a recap of trends, concerns, and emerging questions in Internet measurement.

Synopsis

The authors have blended their academic research and practical experience in Internet measurement and traffic modeling to provide the reader with a structured view to these vast subjects. I would have liked to see a more extensive coverage of *Voice over IP* (VoIP) and its associated performance measurement protocols, *RTP Control Protocol* (RTCP), *RTCP Extended Report* (XR), and *RAQMOM*, given the gradual but inevitable shift of voice traffic from the *Public Switched Telephone Network* (PSTN) to the Internet with *Session Initiation Protocol* (SIP) peering.

Most probably, this book had already been published when the *Federal Communications Committee* (FCC) issued an order in May 2006 for all VoIP service providers to demonstrate compliance with the *Communications Assistance for Law Enforcement Act* (CALEA) wiretapping requirement within a year. This directive represents a notable departure from data anonymization principles covered in the book.

Overall, I consider this book an excellent reference source for diverse research and practical articles published in the field of Internet measurement. I highly recommend it to network planners, engineers, and managers responsible for instrumentation, traffic modeling, or performance analysis.

—Reza Fardid, Covad Communications
rfardid@covad.com

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. Contact us at **ipj@cisco.com** for more information.

Bob Braden and Joyce K. Reynolds receive the 2006 Postel Service Award

Bob Braden and Joyce K. Reynolds are this year's recipients of the Internet Society's prestigious *Jonathan B. Postel Service Award*. The award was presented "For their stewardship of the RFC (*Request for Comments*) series that enabled countless others to contribute to the development of the Internet." The presentation was made by Internet pioneer Steve Crocker (a member of this year's Postel award committee and author of the very first RFC) during the 67th meeting of the *Internet Engineering Task Force* (IETF) in San Diego, California.

The award is named after Dr. Jonathan B. Postel to commemorate his extraordinary stewardship exercised over the course of a thirty year career in networking. Between 1971 and 1998, Postel managed, nurtured and transformed the RFC series of notes created by Steve Crocker in 1969. Postel was a founding member of the Internet Architecture Board and the first individual member of the Internet Society, where he also served as a trustee.

"It is a pleasure and an honor for the Internet Society to recognize the contribution of Bob and Joyce to the evolution of the Internet," said Crocker. "Since its humble beginnings, the RFC series has developed into a set of documents widely acknowledged and respected as a cornerstone of the Internet standards process. Bob and Joyce have participated in this evolution for a very long time and have been primarily responsible for ensuring the quality and consistency of the RFCs since Jon's death in 1998."

Joyce K. Reynolds worked closely with Postel, and together with Bob Braden she has been co-leader of the RFC Editor function at the University of Southern California's *Information Sciences Institute* (ISI) since 1998. In this role she performed the final quality control function on most RFC publications. Reynolds has also been a member of the IETF since 1988, and she organized and led the User Services area of the IETF from 1988 to 1998. In her User Services role, she was an international keynote speaker and panelist in over 90 conferences around the world, spreading the word on the Internet.

Bob Braden, who has more than 50 years of experience in the computing field, joined the networking research group at ISI in 1986. Since then, he has been supported by NSF for research concerning NSFnet and the DETER security testbed, and by DARPA for protocol research. Braden came to ISI from UCLA, where he had technical responsibility for attaching the first supercomputer (IBM 360/91) to the ARPAnet, beginning in 1970. Braden was active in the ARPAnet Network Working Group, contributing to the design of the FTP protocol in particular. He also edited the Host Requirements RFCs and co-chaired the RSVP working group.

The Jonathan B. Postel Service Award was established by the *Internet Society* (ISOC) to honor those who, like Postel, have made outstanding contributions in service to the data communications community. The award is focused on sustained and substantial technical contributions, service to the community, and leadership. With respect to leadership, the nominating committee places particular emphasis on candidates who have supported and enabled others in addition to their own specific actions.

Previous recipients of the Postel Award include Jon himself (posthumously and accepted by his mother), Scott Bradner, Daniel Karrenberg, Stephen Wolff, Peter Kirstein, Phill Gross and Jun Murai. The award consists of an engraved crystal globe and \$20,000.

ISOC (<http://www.isoc.org>) is a not-for-profit membership organization founded in 1992 to provide leadership in Internet related standards, education, and policy. With offices in Washington, DC, and Geneva, Switzerland, it is dedicated to ensuring the open development, evolution and use of the Internet for the benefit of people throughout the world. ISOC is the organizational home of the IETF and other Internet-related bodies who together play a critical role in ensuring that the Internet develops in a stable and open manner. For over 14 years ISOC has run international network training programs for developing countries and these have played a vital role in setting up the Internet connections and networks in virtually every country connecting to the Internet during this time.

First Internet Governance Forum Meeting Concludes

The inaugural meeting of the *Internet Governance Forum* (IGF) took place in Athens, Greece from October 30 – November 2, 2006. For more information see: <http://www.intgovforum.org>

The Government of Brazil will host the next IGF meeting. It will take place in Rio de Janeiro November 12 – 15, 2007.

ARIN to Provide 4-Byte AS Numbers

On August 30, 2006, the *American Registry for Internet Numbers* (ARIN) Board of Trustees, based on the recommendation of the Advisory Council and noting that the Internet Resource Policy Evaluation Process had been followed, adopted the following policy proposal: “2005-9: 4-Byte AS Number.”

Per the implementation schedule contained in the policy (*Number Resource Policy Manual* [NRPM] Section 5.1), commencing January 1, 2007, ARIN will process applications that specifically request 32-bit AS Numbers.

For more information see: <http://www.arin.net/registration>

[Ed. See also: “Exploring Autonomous System Numbers,” by Geoff Huston in *The Internet Protocol Journal*, Volume 9, No. 1, March 2006.]

Celebrating the 25th Anniversary of the TCP/IP Internet Standards

Two of the core protocols that define how data is transported over the Internet are now 25 years old. The *Internet Protocol* (IP) and the *Transmission Control Protocol* (TCP), together known as “TCP/IP,” were formally standardized in September 1981 by the publication of RFC 791 and RFC 793.

Vint Cerf and Robert Kahn are widely credited with the design of TCP/IP, and many others involved in the ARPANET project made significant contributions. The core of the documents was RFC 675, published in December 1974 by Cerf together with co-authors Carl Sunshine and Yogen Dalal. The subsequent sequence of documents leading up to RFC 791 and 793 benefited from the participation of many people including Dave Clark, Jon Postel, Bob Braden, Ray Tomlinson, Bill Plummer, and Jim Mathis, as well as other unnamed contributors to the definition and implementation of what became the Internet’s core protocols.

“We can’t yet say that the Internet is mature,” says Brian Carpenter, chair of the IETF, “but it’s a great tribute to its pioneers that the two most basic specifications that were published a quarter of a century ago are still largely valid today. I hope the IP version 6 standard will do as well.”

The *Request For Comments* (RFC) series, which was launched in 1969 by Steve Crocker at UCLA (and edited for many years by the late Jon Postel), continues today as the public archive of the Internet’s fundamental technology. Since 1977 it has been hosted by The University of Southern California’s *Information Sciences Institute* (ISI). ARPA support ended in 1998, at which time ISOC took over providing funding for the publication of Internet standards. More recently, ISOC extended its support to include other areas critical to the open development of Internet standards.

See also:

<http://www.ietf.org/rfc/rfc0791.txt>

<http://www.ietf.org/rfc/rfc0793.txt>

<http://www.isoc.org/standards/tcpip25years>

<http://www.isoc.org/internet/history/brief.shtml>

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, trouble-shooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2006 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA
--

The Internet Protocol Journal

March 2007

Volume 10, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
AAA—Part One.....	2
DNS Infrastructure	12
Writing RFCs Using XML....	25
Fragments	29
Call for Papers.....	31

FROM THE EDITOR

Every time you dial into a service provider network or connect to a wired or wireless network that offers Internet access, you are most likely using several components of what is referred to as *Authentication, Authorization, and Accounting*, or “AAA” for short. The AAA space is quite complex, so when we asked Sean Convery to give us an overview of these technologies, he decided to divide his survey into two parts. Part One—subtitled “Concepts, Elements, and Approaches”—is included in this issue. Part Two, which discusses protocol details and applications, will follow in our next issue.

The *Domain Name System* (DNS) has been discussed previously in this journal. The most critical part of the DNS is the collection of *Root Servers*. For protocol reasons, there are only 13 “logical” root servers, but a system of more than 100 servers has been deployed using a technique known as *anycast*. Steve Gibbard examines the distribution of the root servers in different parts of the world and discusses operational aspects of the DNS.

If you are tracking any part of the IETF process, you should be aware of several important resources. First, the *IETF Education Team* (<http://edu.ietf.org/>) offers training sessions and educational materials. Second, the *IETF Journal* (<http://www.isoc.org/ietf-journal>) publishes timely reports and updates on the activities of the IETF.

Finally, the *IETF Tools Team* (<http://www.ietf.org/tools.html> and <http://tools.ietf.org>) provides many tools and applications for protocol developers. Marshall Rose and Carl Malamud take a closer look at one of these tools, namely a system for writing Internet Drafts and RFCs using XML.

Please take a moment to renew and update your subscription. You can access your subscription record by clicking on the “Subscriber Services” link at <http://www.cisco.com/ipj>.

—Ole J. Jacobsen, *Editor and Publisher*
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Network Authentication, Authorization, and Accounting

Part One: Concepts, Elements, and Approaches

by Sean Convery, Identity Engines

Network *Authentication, Authorization, and Accounting* (AAA, pronounced “triple-A”) is a technology that has been in use since before the days of the Internet as we know it today. Authentication asks the question, “Who or what are you?” Authorization asks, “What are you allowed to do?” And finally, accounting wants to know, “What did you do?” These fundamental security building blocks are being used in expanded ways today. This article, the first in a two-part series, focuses on the overall concepts of AAA, defines the elements involved in AAA communications, and discusses high-level approaches to achieving specific AAA goals. Part two of the article, to be published in a future issue of IPJ, will discuss the protocols involved, specific AAA applications, and considerations for the future of AAA.

AAA, at its core, is all about enabling mobility and dynamic security. Without AAA, a network must be statically configured to control access; IP addresses must be fixed, systems cannot move, and connectivity options should be well defined. Even the earliest days of dialup access broke this static model, thereby requiring AAA. Today, the proliferation of mobile devices, diverse network consumers, and varied network access methods combine to create an environment that places greater demands on AAA.

AAA has a part to play in almost all the ways we access a network today. Emerging technologies such as *Network Access Control* (NAC) extend AAA even into corporate Ethernet access (historically the “trusted” network that set the benchmark level of security that all other types of access had to match). Today, wireless hotspots need AAA for security, partitioned networks require AAA to enforce segmentation, and remote access of every kind uses AAA to authorize remote users.

It is not clear when the term AAA first gained acceptance, but an examination of academic papers finds “authentication, authorization, and accounting” used as a discrete term (albeit without the AAA acronym) as early as 1983 in an IEEE paper^[1]. Though mired in pre-Internet *Open Systems Interconnection* (OSI)-centric terminology, the ordering of the “A’s” is the same as today’s usage.

For most network administrators, the genesis of AAA coincided with the development of the *Remote Authentication Dial-In User Service* (RADIUS) protocol^[2]. RADIUS was developed by Livingston Enterprises (now part of Alcatel-Lucent) in the early 1990s, became an Internet standard through the IETF in 1997, and today is the most widely accepted AAA protocol.

Another widely adopted AAA protocol, which predates RADIUS as an RFC by four years, is the *Terminal Access Controller Access Control System* (TACACS)^[3]. Though never an Internet standard, TACACS evolved into XTACACS and then TACACS+, the latter of which is the only version of TACACS in use today.

Before we delve into the details of these protocols, it is important to understand the roles played within a AAA system.

Core Components of AAA

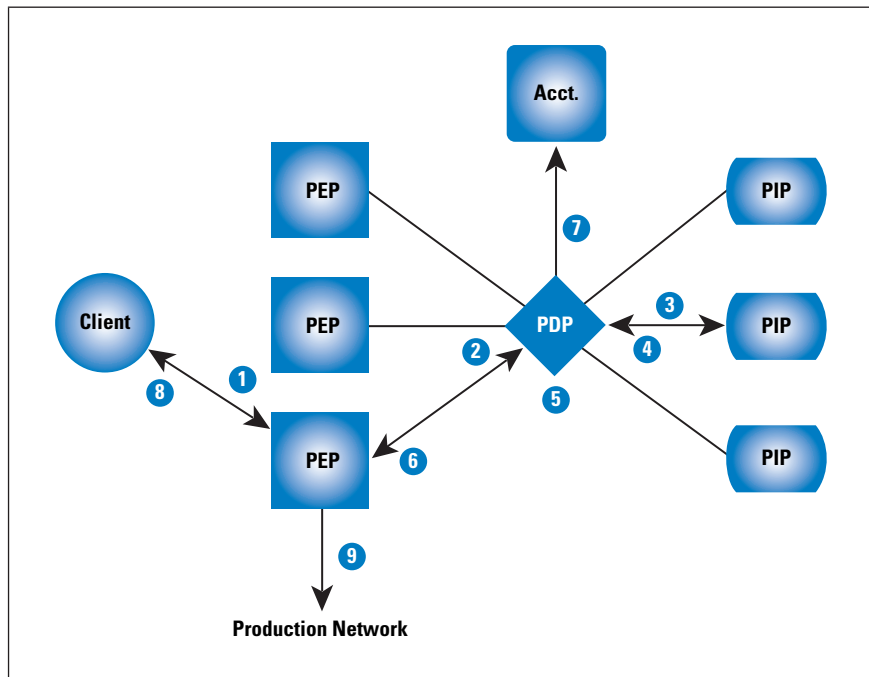
- *Client*: The client is the device attempting to access the network. The client either authenticates itself, or it acts as a proxy to authenticate the user.
- *Policy Enforcement Point (Authenticator)*: The Policy Enforcement Point (PEP) is sometimes called the authenticator or *Network Access Server* (NAS). The PEP is the network device that brokers the access request for the client. The PEP can be a dial-in server, VPN concentrator, firewall, gateway *General Packet Radio Service* (GPRS) support node, Ethernet switch, wireless access point, or an inline security gateway. The PEP is responsible for enforcing the terms of a client's access. This enforcement varies based on the capabilities of the PEP and is discussed later in this article.
- *Policy Information Point*: The Policy Information Point (PIP) is a repository of information to help make the access decision. It could be a database of device IDs, a user directory such as the *Lightweight Directory Access Protocol* (LDAP), a *one-time password* (OTP) token server, or any other system that houses data relevant to a device or user access request.
- *Policy Decision Point (AAA Server)*: The Policy Decision Point (PDP) is the brain of the AAA decision. It collects the access request from the client through the PEP. It also queries any relevant PIPs to gather the information it needs to make the access decision. The PDP, as its name implies, is the entity that makes the final decision around network access. It also can send specific authorizations back to the PEP that apply settings or constraints to the client's network traffic.
- *Accounting and Reporting System*: Whether on a dedicated system or built as part of a PDP, tracking use of the network with accounting is one of the best features of AAA. With all forms of network access now offering controlled access, the AAA service can tell you who got on the network, from where, and what that person was granted access to.

It is important to note that the preceding categories are logical containers of functions and not necessarily dedicated physical devices. Often elements are combined, such as PEP with PDP, and PDP with PIP.

Example AAA Flow

Now that we have examined the components of a AAA solution, walking through a typical use case will help cement our understanding of the role that each entity plays. Figure 1 shows an example of a client attempting to gain access to the network.

Figure 1: A Client Connects to a AAA-Protected Network



1. The client attempts to connect to the network, is challenged for identity information, and sends this information to the PEP. In this example, let's assume the client is a laptop with a worker attempting to access an organization's VPN from a remote location. Additionally, we'll assume this is a valid, permitted use of the network.
2. The PEP sends the collected identity information to the PDP. In some cases (discussed in part two of this article), the PEP cannot see the specific identity information provided but instead relays the information directly to the PDP.
3. The PDP queries any configured PIPs for information about the client and validates that the credential provided by the client is valid. In this example, the PIP is an LDAP directory.
4. The PIP returns a success or failure message from the credential validation step and sends additional information about the client to the PDP for evaluation. This information could include the role of the user, the home location for the user, and so on.
5. The PDP evaluates information learned about the client through the client, PEP, and PIP; the role of the PEP and PIP that serviced the request; and any contextual information (such as time of day) against its configured policies. Based on this information, the PDP makes an authorization decision.

6. The PDP sends the PEP the authentication result and any authorizations specific to the client. These authorizations trigger specific PEP actions to apply to the client. For example, the authorization data might trigger specific *Access Control Lists* (ACLs) or IP pool assignments for the client.
7. The PDP also sends the result of this transaction to the accounting system.
8. The PEP applies the authorization profile learned from the PDP and sends the “authentication successful” message to the client. The PEP can also be configured to send accounting information on this new connection to the accounting and reporting system.
9. The client accesses the production network through the PEP.

Elements of Authentication

When performing authentication, numerous elements can be evaluated before a PDP reaches its access decision. At a high level, these elements can be broken down into three categories: the principal itself (the user, device, or service requesting access), the credential the principal submits (shared key, one-time password, digital certificate, or biometric credential), and the contextual information describing the transaction (location, time of day, software state, and so on).

- *Principal*: The principal is the entity requesting authorization. It is generally some combination of user, device, or service. When concerned with a user, the PIP can provide attributes about the user such as role or group affiliations, job title, e-mail address, physical address, and so on. In specific applications, it can include much more granular information. For example, a higher-education facility might be interested in knowing a student’s class schedule when servicing the student’s authentication request. When the principal is a device, the same thinking applies. The PIP can inform the PDP if the device is a managed asset, what its basic usage parameters are, and so on. User and device authentication can be carried out sequentially for the same transaction, often involving device authentication first and then user authentication. Lastly, a service such as a network management process can authenticate. In this case, the service almost always looks like a user to the AAA infrastructure and is handled accordingly.
- *Credential*: The next element the PDP considers is the credential the user or device submits as proof of identity. There are four main types of credentials: shared key (password), *one-time password* (OTP), digital certificate, and biometric credential. This section examines each of these types. The first and most widely used form of credential is the shared key, typically a user password. AAA deployments that use shared keys can be subdivided based on the protocol the system uses to verify the password, including the *Password Authentication Protocol* (PAP)^[4], *Challenge Handshake Authentication Protocol* (CHAP)^[5], and *Microsoft CHAP Extensions* (MS-CHAP) Versions 1^[6] and 2^[7]. PAP authentication is a plaintext authentication method that is not recommended for use in security-sensitive environments.

However, many newer protocols provide a secure transport for PAP, making its use in AAA still quite common. Some of these methods are discussed in part two of this article. CHAP improves on the security of PAP by not sending the password in the clear but rather a challenge based on a hash of the password. MS-CHAP is a Microsoft extension to CHAP that tunes things a little bit for Microsoft environments. Version 2 of MS-CHAP addresses security weaknesses in Version 1. MS-CHAPv2 is quite common today in Microsoft environments. CHAP in all its forms is vulnerable to dictionary attacks because even if a hash cannot be decrypted, common passwords can be guessed and those hash values can be computed.

A second, also widely used credential type is the OTP. At login time, users refer to their personal token to get the OTP they will type in. The token is generally provided in hardware or software form. Tokens are designed to generate seemingly random passwords that are synchronized with a token server acting as a PIP. The OTP can be sent in the clear because it is used only once; after a configurable time (for example, 30 seconds) a new password is generated. When an OTP is combined with a *Personal Identification Number* (PIN), two-factor authentication is achieved because the client needs to have something (the token) and know something (the PIN).

The third type of credential is the *digital certificate*. Digital certificates can be stored either locally on the client or on some sort of removable device such as a smartcard. A full discussion of asymmetric-key cryptography is outside the scope of this article, but at a high level, certificates work by asserting the identity of their bearer by having the certificate signed by a trusted *Certificate Authority* (CA). CAs can be external entities such as a government or commercial enterprise or they can be internal to a given organization. The certificate itself can be freely distributed, because the only way it can be validated as belonging to the rightful owner is in combination with the private key. Because they reside on the client, certificates are most often used to authenticate a physical entity rather than an individual. However, smartcards are changing this paradigm by enabling users to take their digital certificate (and private keys) with them, thereby disassociating the certificate from the machine itself. Similar to an OTP without a PIN, a digital certificate or smartcard alone does not provide two-factor authentication. Certificate deployments, particularly smartcards, are addressing this problem by requiring a PIN to unlock access to the credential.

The fourth and least widely deployed type of credential is the *biometric credential*. Biometrics^[12] ignores something you *have* and something you *know* and instead focus on something you *are*. Fingerprint scanners, iris scanners, and facial recognition are all forms of biometric authentication. Because biometrics is the newest form of credential, it is currently experiencing heightened anticipation among users regarding potential applications—and also scrutiny for potential weaknesses.

- *Contextual*: The last element the PDP typically considers in its authentication decision is the contextual information associated with the AAA request, including the network and physical location of the request, the type of access provided by the PEP, the time of day, and potentially other elements such as network load, security threat level, and so on. A relatively new entrant into this set of contextual information is client device posture, typically discussed under the rubric of *Network Access Control* (NAC). NAC or posture checks examine the software state of the client before it connects. NAC data allows the PDP to assess the degree of risk posed by the connecting client before granting the client access to the network. For example, if a system is running an out-of-date operating system, has no current security applications running, or otherwise exhibits high-risk behavior, it may not be granted access to the network. NAC will be discussed in more detail in part two of this article.

Authorization Approaches

At its core, authorization means determining what a client is allowed to do on the network. However, the granularity of this authorization is only as good as the sophistication of the PDP and the enforcement capabilities of the PEP. This section examines the authorization options for network AAA, including Layer 2 segmentation, Layer 3 filtering, and Layer 7 entitlements. It closes with an examination of some of the challenges encountered when sending or “provisioning” the authorizations from the PDP to the PEP.

- *Null Authorization (Authentication Only)*: Strangely the most common authorization in AAA is no authorization at all. After the authentication event occurs, the client is immediately granted full access to the network. This characteristic is a holdover from the original goal of remote-access AAA: to perform an authentication check that simply determines whether the client should be trusted as if it were connected to the organization’s home network. Because these home networks employed no segmentation or filtering within them, it was natural that remote-access techniques such as dialup and VPN would likewise employ neither. Today however, authentication is increasingly being used for all forms of network access, with a goal of providing clients with network rights commensurate with their role in the organization. This latter goal requires a strong authorization foundation through the cooperation of the PDP and PEP.
- *Layer 2 Segmentation*: For wireless access points and Ethernet switches, the most common form of authorization enforcement is Layer 2 segmentation, which works by splitting the network into multiple logical segments, isolating certain classes of client from one another. This process is most typically achieved by deploying *Virtual LANs* (VLANs), which separate the members of one VLAN from other VLANs in the same Layer 2 network—even though the VLANs traverse the same physical network infrastructure.

VLANs can be used to restrict access to specific resources by working in coordination with VLAN-specific ACLs on Layer 3 devices upstream from the Layer 2 device. For access points, a given wireless *Service Set Identifier* (SSID) can be associated with a VLAN on the wired side of the access point. *Multiprotocol Label Switching* (MPLS) is more commonly associated as a WAN transport, but there is nothing to prevent labels for traffic based on AAA. More commonly, the client is associated with a VLAN and the VLAN is associated with an MPLS label further into the infrastructure.

- *Layer 3 Filtering:* Layer 3 filtering authorizes access to resources through ACLs configured on Layer 3 devices (routers, Ethernet switches, security gateways, and so on). These ACLs (which generally encompass Layer 4 of the OSI stack as well) can enforce authorizations to a range of hosts, specific hosts, or services on those hosts. As mentioned earlier, Layer 3 filtering can be combined with Layer 2 segmentation to provide aggregate authorizations for an entire VLAN. This filtering is the most common technique on network infrastructure devices, whereas security gateways tend to apply ACLs to specific clients. Additionally, technologies such as *IP Security* (IPsec)^[8] provide a Layer 3 filtering capability by allowing only certain types of traffic to travel through the VPN tunnel.
- *Layer 7 Entitlements:* Increasingly, security gateways are able to go beyond Layer 3 and 4 filtering and are starting to become application-aware, meaning that the authorizations handed from the PDP to the PEP can be very granular, focusing on the specific applications that are needed rather than broader filters based on segments or hosts on the network. Because this technology is still relatively new, there are no standards yet to make this interaction work transparently. As a result, most granular application filters are written on the PEP itself in order to allow the PDP to trigger a preexisting profile on the PEP. These sorts of provisioning challenges are discussed further in the next section.
- *Provisioning Challenges:* In AAA parlance, the term “provisioning” refers to communicating a user’s session rights and constraints to the PEP so that the PEP can grant and enforce these permissions. One of the most difficult aspects of provisioning access rights on a PEP is communicating the decision of the PDP in a format the PEP can understand. This fact is one of the reasons that many PEPs come with a lightweight PDP. This approach solves the narrow problem for that PEP but creates management challenges when coordinating network AAA across a broader enterprise, because the enterprise AAA policies must be implemented individually on each unique type of PEP on the network. Because RADIUS is the most commonly used network AAA protocol, it is natural to communicate the PDP decision using that protocol. RADIUS attributes such as the “filter-id” allow the PDP to trigger a preexisting filter on the PEP.

In addition, many PEP vendors support *Vendor Specific Attributes* (VSAs) in RADIUS to enable the PDP to speak the language of the PEP more specifically. This process works well but creates a significant amount of work on the PDP to enable it to translate the policy result and correctly communicate it to each type of PEP. Another option soon to be sanctioned by the standards bodies is an extension to RADIUS that enables the sending of standard IP ACLs using RADIUS attributes^[9].

One further option for provisioning is through the *Simple Network Management Protocol* (SNMP), which is typically used to assign Layer 2 ports to VLANs or to enable or disable interfaces. This process can work, but remember that the version of SNMP typically in deployment is still SNMPv2c, which is *User Datagram Protocol* (UDP)-based (connectionless) and unencrypted. Therefore, the SNMP traffic is prone to packet loss when links are congested or devices are busy, thereby requiring costly application layer retransmission schemes. It also means the transmissions themselves are vulnerable to inspection or modification. These attributes make SNMP generally a poor choice for security-sensitive tasks. RADIUS also uses UDP, but supports basic retransmission as part of the protocol.

Another provisioning method used today is standard *Secure Shell* (SSH) Protocol or HTTPS-based configuration. This method manages a device through standard administrative interfaces to set enforcement techniques. Although this method gives the PDP full access to the features of the PEP, it is very difficult to coordinate the dynamic aspects of the client AAA event with the static elements of the running configuration of the PEP. Finally, new protocols are emerging to make provisioning easier. NETCONF^[10] is an *Extensible Markup Language* (XML)-based protocol designed as a replacement for network management applications connecting to devices over the *command-line interface* (CLI).

As this section has shown, there are numerous approaches to authorization in AAA. Each PEP has its own capabilities, but the challenge for a diverse network is to consistently authorize clients, regardless of the given PEP they access the network through.

Accounting Techniques

Accounting is an increasingly critical step in the overall AAA process. Regulatory controls are starting to mandate better auditing of network access. The last stage of AAA, accounting simply records which clients accessed the network, what they were granted access to, and when they disconnected from the network. Accounting has always been widely used in the *Internet Service Provider* (ISP) space because auditing network access is the basis for billing ISP customers. Increasingly, accounting is being used as a way to correlate client attribute information (username, IP address, etc.) with actions and events on the network.

This correlation can make other systems that are not user-aware more intelligent in the security decisions that they make. For example, a network *Intrusion Detection System* (IDS) can learn a lot about the behavior of a given IP address. However, when that information is correlated with the user assigned to that IP address—and the permissions that user should have—the relevance of the IDS data increases dramatically.

One of the design considerations of accounting systems is that, given the centralized nature of audit and the decentralized nature of access, they are generally out-of-band with the client's normal communications. This makes them excellent resources to refer to when the network administrator wants to know when the client connected and what the client was granted access to. However, their out-of-band nature makes them poor resources for determining what the client actually did while connected to the network. This information can be learned by the network, as mentioned earlier, by coordinating the AAA accounting information with the rest of the network enforcement and monitoring systems.

Summary and Part Two Teaser

This first part of this article introduced AAA and described many of the foundation concepts necessary to gain a sound understanding of the overall system. After defining the elements involved, a sample flow of a AAA event was described. Additionally, the high-level approaches to authentication, authorization, and accounting were discussed. Part two of this article will discuss the protocols used in AAA, including not just RADIUS, *Extensible Authentication Protocol* (EAP), TACACS+, and Diameter, but many others. Additionally, specific applications of AAA technology will be described, and some conclusions will be drawn as to what the future holds for AAA.

References

- [1] Lagsford et. al., "OSI Management and Job Transfer Services," *Proceedings of the IEEE*, Volume 71, No. 12, December 1983.
- [2] Rigney et. al., "Remote Authentication Dial In User Service (RADIUS)," RFC 2865 (Obsoletes RFC 2138, 2058), June 2000.
- [3] Finseth C., "An Access Control Protocol, Sometimes Called TACACS," RFC 1492, July 1993.
- [4] Lloyd et. al., "PPP Authentication Protocols," RFC 1334, October 1992.
- [5] Simpson W., "PPP Challenge Handshake Authentication Protocol (CHAP)," RFC 1994, August 1996.

- [6] Zorn et. al., “Microsoft PPP CHAP Extensions,” RFC 2433, October 1998.
- [7] Zorn et. al., “Microsoft PPP CHAP Extensions, Version 2,” RFC 2759, January 2000.
- [8] Kent et. al., “Security Architecture for the Internet Protocol,” RFC 2401, November 1998.
- [9] Congdon et. al., “RADIUS Filter Rule Attribute,” Internet Draft, Work in Progress, January 2007,
draft-ietf-radext-filter-08.txt
- [10] Enns et. al., “NETCONF Configuration Protocol,” RFC 4741, December 2006.
- [11] Dory Leifer, “Visitor Networks,” *The Internet Protocol Journal*, Volume 5, No. 3, September 2002.
- [12] Edgar Danielyan, “The Lures of Biometrics,” *The Internet Protocol Journal*, Volume 7, No. 1, March 2004.

SEAN CONVERY is CTO at Identity Engines, a venture-backed startup developing innovative identity management solutions for enterprise networks. Prior to Identity Engines, Sean (CCIE® no. 4232) worked for seven years at Cisco Systems, most recently in the office of the security CTO. Sean is best known as the principal architect of the SAFE Blueprint from Cisco and the author of *Network Security Architectures* (Cisco Press, 2004). Sean has presented to or consulted with thousands of enterprise customers around the world on designing secure networks. Before Cisco, Sean held various positions in IT and security consulting during his 14 years in networking. E-mail: **sconvery@idengines.com**

Geographic Implications of DNS Infrastructure Distribution

by Steve Gibbard, Packet Clearing House

The past several years have seen significant efforts to keep local Internet communications local in places far from the well-connected core of the Internet. Although considerable work remains to be done, Internet traffic now stays local in many places where it once would have traveled to other continents, lowering costs while improving performance and reliability. Data sent directly between users in those areas no longer leaves the region. Applications and services have become more localized as well, not only lowering costs but keeping those services available at times when the region's connectivity to the outside world has been disrupted. I discussed the need for localization in a previous paper, "Internet Mini-Cores: Local connectivity in the Internet's spur regions."^[1] What follows here is a more specific look at a particular application, the *Domain Name System* (DNS).

Most Internet applications depend on the DNS, which maps human-readable domain names to the *Internet Protocol* (IP) addresses computers understand. Two Internet hosts may have connectivity to each other but be unable to communicate because no DNS server can be reached. This article examines the placement of DNS servers for root and top-level domains and the implications of that placement on the reliability of the services these servers provide in different parts of the world. It is not a "how-to" guide to the construction of DNS infrastructure and does not contain recommendations on DNS policy; it is rather a look at the placement of DNS infrastructure as currently constructed.

Although it is possible to access Internet resources without the DNS by entering numeric IP addresses directly, this type of access is not generally done. IP addresses, such as **209.131.36.158**, are difficult to remember, are generally unpublished outside the DNS, and often change without notice. Local caching of DNS information can mask temporary problems with DNS data for commonly accessed domain names, but caches are emptied when caching resolvers are restarted, data in caches expires, and nothing is cached until the first time it is accessed by a local user.

It should be noted that information about DNS deployment is changing rapidly. Several organizations are working on new DNS deployment. Information in this article can be considered current, to the best of my knowledge, as of May 2006.

DNS Hierarchy

The DNS is a hierarchy of domains within domains. The levels of the hierarchy are separated by dots. At the top of the hierarchy is the *root*, usually invisible but sometimes represented as a trailing dot. Using **www.yahoo.com** as an example, the **com** domain is contained within the root. **com** contains **yahoo**, and **yahoo** contains **www**. Domains in the position **com** takes in this example are known as *Top-Level Domains*, or TLDs; they are the first level in the root domain. Domains in the position of **yahoo** are known as *Second-Level Domains*. In this example, **www** occupies the third level, and so forth.

The information that makes up the Domain Name System is stored on DNS *servers*. That information is divided into *zones*, which for our purposes are synonymous with domains. Each zone is stored on a set of *authoritative servers*, which are queried when users or applications attempt to access a service on the Internet. In the simplest case, a domain name query works like the following:

A *caching resolver* (so named because it caches information it receives) that has not yet cached any DNS zone data receives a query for **www.yahoo.com**. Because its cache is empty, it uses the *hints* distributed with the DNS software to contact one of the root servers and asks, “Where is **www.yahoo.com**?” The root server replies with a list of servers for **.com**. The caching server then asks one of the **.com** servers, “Where is **www.yahoo.com**?” and gets a response that directs it to servers for **yahoo.com**. It asks the same question of those servers and finally gets an answer to the question it was asking.

Generally several servers can answer questions about any domain, but if all the servers for any single level are broken or unreachable, the query fails and the service the user is looking for is unreachable. It is therefore important that the DNS be reliable, and that the servers for each zone throughout the hierarchy be reachable from anywhere the servers they point at are being used.

Root Servers

Without root servers, none of the DNS works. As of this writing, 117 root servers exist worldwide, operated by 12 different organizations.^[2] Root servers are added frequently, so the number may be significantly greater by the time this article is in circulation.

Because of protocol limitations, the root servers can use only 13 IP addresses. Each root-server operator is responsible for one or two of those addresses. Using a technique called *anycast*, which allows servers in separate locations to share a *single* IP address, six of those operators operate multiple servers using the same IP address^[3], meaning that only 13 of them are visible at the same time from any single location, but those 13 should in most cases include the topologically closest one.^[4]

The distribution of root servers is rather uneven. North America and Europe have similar numbers: 38 in North America and 35 in Europe. The 35 in Europe are distributed fairly evenly, with the largest concentrations (four each) in London and Amsterdam, Europe's two largest Internet hubs. North America has 8 in Washington, D.C., 8 in the San Francisco Bay Area, and 5 in Los Angeles. In the United States, all cities that host root servers are on the coasts except Atlanta and Chicago. All seven of the remaining IP addresses represented by only a single server, known as *unicast roots*, are in the Washington, D.C., San Francisco, and Los Angeles areas.

Australia has two root servers in Brisbane, one in Perth, and one in Sydney. New Zealand has two, one in Wellington and one in Auckland. Singapore and the wealthier parts of East Asia are well-covered, and there are two root servers in Jakarta and one each in Bangkok and Kuala Lumpur. A year ago, there were none in the vast expanse between Bangkok and Dubai, but three have recently been added in India, along with others in Dhaka and Karachi. Mainland China and the former USSR each have two. There are three in Africa: two in Johannesburg and one in Nairobi. Another will be installed in Nairobi shortly, but most of the rest of Africa lacks direct connectivity to Johannesburg or Nairobi and must cross satellite or intercontinental fiber links to reach the nearest root servers. All four of the root servers in South America are in Brazil and Chile, with two in Sao Paulo and one each in Brasilia and Santiago de Chile.

With some exceptions, root-server density tends to correlate strongly with per-capita income. This fact is not surprising—it is true for other forms of infrastructure as well—but it means that those with the greatest dependence on external infrastructure are those least able to pay for external connectivity.

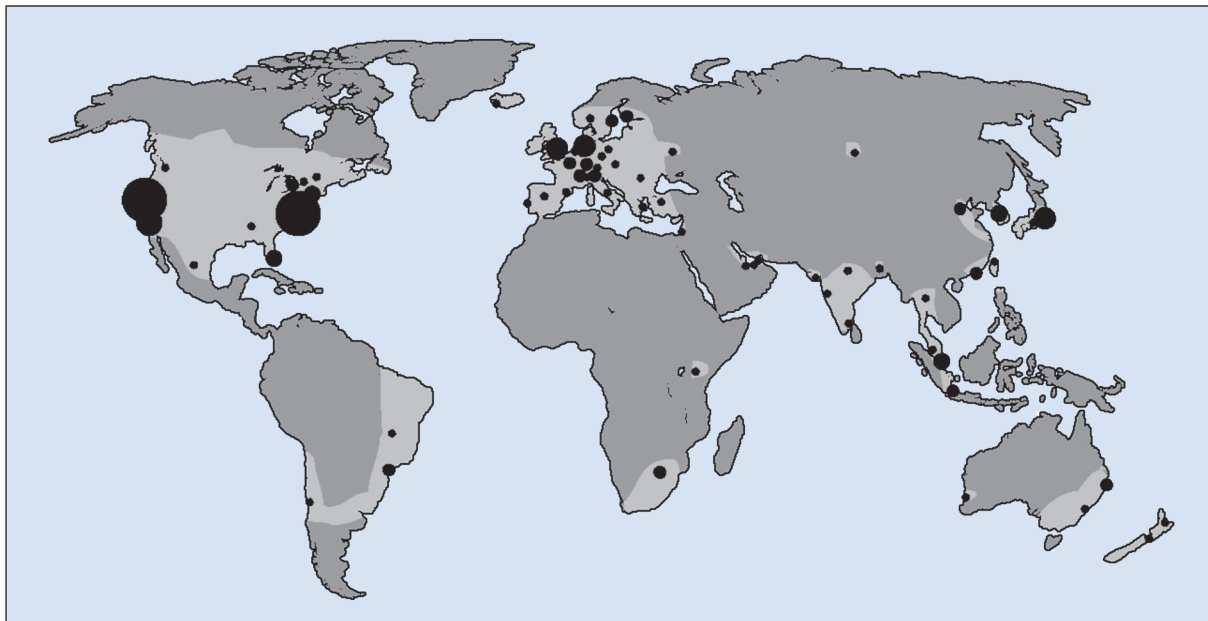
Root-Server Placement

In areas that have local root servers, finding the name servers for a top-level domain should be fairly reliable. In areas without local root servers, the ability to query the root servers is dependent on other long-distance infrastructure. In some places this infrastructure is well-developed, so this problem is not a significant one. Elsewhere long-distance infrastructure is slow, expensive, and unreliable, consisting of satellite links or a single fiber connection that may take several days to fix if it breaks.

Sri Lanka, for example, is connected to the rest of the world by a single fiber connection, which was cut in 2004 by a ship that dragged anchor in the Colombo harbor.^[5] Although Sri Lanka has an exchange point that should have allowed connectivity to local Internet services, news reports said that Sri Lanka's "Internet and long-distance phone service" had been cut off. I have not received a good account of what Internet connectivity looked like from anyone in Sri Lanka at the time, but it is likely that even local Internet connections would not have worked without a local root server.

Sri Lanka is not an isolated case. The dots in Figure 1 show the locations of all root servers. The light grey areas are regions in which multiple fiber paths are available to root servers. The remainder of the world can reach root servers only by a single fiber path or by satellite. Large areas of the world are poorly covered.

Figure 1: Root Server Locations and Areas of Redundant Connectivity



Root-Server Expansion

Four of the 12 root-server operators are presently working to install root servers in areas that lack them. Although the 117 root servers currently in operation are a big improvement over the 13 that were in operation three years ago, many regions still do not have any. Those root-server operators are installing servers wherever they can get the funding to do so.

Funding is generally provided either by grants, especially from the *Asia Pacific Network Information Centre* (APNIC) in the Asia-Pacific region, by local governments or *Internet Service Provider* (ISP) associations. Because the addition of new anycast copies of root servers is relatively easy given sufficient funding, the main limitation preventing the installation of root servers in new locations is lack of funding.

One question probably best addressed in a more central manner is whether it makes sense to have many copies of one or two root-server IP addresses in some regions or whether it would be better to have more of a mix of root-server IP addresses. Currently, only 6 of the 13 root-server addresses are anycasted, only 4 are anycasted in large numbers, and 2 of those focus on specific regions, meaning that in many of the more remote parts of the world the only nearby root servers are *Internet Systems Consortium* (ISC)'s "F" and Autonomica's "I" roots, and some places have several of one of those closer than the next one of the other.

Because some DNS resolvers have their own mechanisms for finding the closest server and for handling failures of types that do not include route withdrawals, having multiple IP addresses nearby seems like a good thing. A more complex question is whether it would be worthwhile to anycast all 13 of them widely, or if there is some smaller number that would be sufficient to have nearby. Previous research on this topic has assumed a limit of 13 root servers, producing conclusions that are not applicable to the modern Internet.^[6]

This article should not be seen as a criticism of the places with large numbers of root servers. Although the U.S. distribution looks strange, with the San Francisco Bay and D.C. area clusters perhaps excessive, it comes close to following the Internet topology in the United States. Indeed, the U.S. concentration may be appropriate to handle server load. Western Europe's dense but relatively even distribution of root servers through the region appears to be an optimal distribution, because most populated areas have multiple root servers nearby. Likewise, Jakarta is one of the very few cities in the developing world to have more than one, and that provides local redundancy that much of the developing world lacks. If root-server deployment were funded from a single global budget, the distribution across the world's regions would look very unfair. But because Internet infrastructure is mostly funded locally, Jakarta and Western Europe are examples other regions could emulate.

TLD Distribution

Use of the DNS also requires access to TLD servers. To access something in the **.com** domain, a user's local DNS resolver must be able to reach the **.com** servers. This statement is true for any TLD, whether it is a *generic TLD* (gTLD), such as **.com**, **.net**, and **.org**, or a *country code TLD* (ccTLD). Unlike the root, it is not necessary that all TLDs be reliable from all locations; if a TLD is not used to name local resources in a region, having local access to that TLD will not help if that the region gets cut off from the rest of the world.

gTLD Distribution

Of the gTLDs, **.com** is by far the largest. It is well-connected to the Internet core, the area with well-meshed internal connectivity mainly comprising North America, Western Europe, East Asia, and Singapore. (See Figure 2.) The **.com** servers are located in Australia, Brazil, Japan, South Korea, the Netherlands, Sweden, the United Kingdom, and the U.S. states of California, Florida, Georgia, Virginia, and Washington. The **.com** servers are well-connected to areas well-connected to those regions but poorly connected to Africa, South Asia, and parts of South America.

Figure 2: Server Locations for .com and .net and Areas of Redundant Connectivity



UltraDNS, the operator of **.org**, **.info**, **.mobi**, and **.coop**, among others, is also somewhat well-connected to the Internet core, although not to the extent the **.com** servers are. It has publicly accessible servers in four metropolitan areas in the United States as well as in London and Hong Kong. It has a couple of noncore locations, in Delhi and Johannesburg. UltraDNS also has servers in other locations, accessible only to the resolvers of certain large ISPs. Because those servers are not available to the general public in their regions, they are omitted from discussion here. (See Figure 3.)

Figure 3: Server Locations for .org, .info, and .mobi and Areas of Redundant Connectivity



Other gTLDs do not do considerably better. Table 1 shows the locations of all the gTLDs.

Table 1: Locations of TLD Servers

gTLD	Locations by Country or U.S. State
.aero	Switzerland, Germany, India, Hong Kong, United Kingdom, and the following states in the United States: California, Illinois, and Virginia
.biz	Australia, Hong Kong, Netherlands, New Zealand, Singapore, United Kingdom, and the following states in the United States: California, Florida, Georgia, New York, Virginia, and Washington
.com	Australia, Brazil, Canada, Japan, South Korea, Netherlands, Sweden, Singapore, United Kingdom, and the following states in the United States: California, Florida, Georgia, Virginia, and Washington
.coop	United Kingdom and the following states in the United States: California, Illinois, and Massachusetts
.edu	Netherlands, Singapore, and the following states in the United States: California, Florida, Georgia, and Virginia
.gov	Canada, Germany, and the following states in the United States: California, Florida, New Jersey, Pennsylvania, and Texas
.info	India, Hong Kong, South Africa, United Kingdom, and the following states in the United States: California, Illinois, and Virginia
.int	Netherlands, United Kingdom, and California in the United States
.jobs	Netherlands, Singapore, and the following states in the United States: California, Florida, Georgia, and Virginia
.mil	The following states in the United States: California, Maryland, Virginia, and other unknown locations
.mobi	India, Hong Kong, South Africa, United Kingdom, and the following states in the United States: California, Illinois, and Virginia
.museum	Sweden and California in the United States
.name	Singapore, United Kingdom, and the following states in the United States: California, Florida, Georgia, Virginia, and Washington
.net	Australia, Brazil, Canada, Japan, South Korea, Netherlands, Sweden, Singapore, United Kingdom, and the following states in the United States: California, Florida, Georgia, Virginia, and Washington
.org	India, Hong Kong, South Africa, United Kingdom, and the following states in the United States: California, Illinois, and Virginia
.pro	Canada and the following states in the United States: Illinois and Texas
.travel	Australia, Hong Kong, Netherlands, New Zealand, Singapore, United Kingdom, and the following states in the United States: California, Florida, Georgia, New York, Virginia, and Washington

Although gTLDs are typically marketed for their applicability to specific types of organization, or in the case of **.com** because it is the only domain many people have heard of, geography should also be considered in selecting domains. Most of the gTLDs have reasonable coverage throughout the Internet core region, but there are exceptions. The **.int** and **.museum** domains are hosted only in North America and Europe, and **.pro** is hosted only in North America.

Outside the Internet core there is little gTLD presence. Only **.biz**, **.travel**, **.com**, and **.net** are present in Australia and New Zealand. South Africa and India have **.aero**, **.info**, **.mobi**, and **.org**, making them the only gTLDs hosted in either Africa or the South Asian region. South America hosts only **.com** and **.net**, with servers in two cities in Brazil. Taken together, these are the only Southern Hemisphere gTLD locations as of this writing, and no gTLD has any presence in parts of the world without external fiber-optic connectivity, although that may be changing.

Where gTLDs should be hosted, and with what scope, are somewhat open questions. Should these domains address resources anywhere, or should their scope be local? This question is really one for the *Internet Corporation for Assigned Names and Numbers* (ICANN), or for the gTLD sponsors or registries, and beyond the scope of this article. Verisign, the company that administers **.com** and **.net**, points out that database replication with the amount of changes in the **.com** zone is a significant problem over slow network links.

ccTLD Distribution

Questions about where ccTLDs, the top-level domains assigned to individual countries, ought to work seem more straightforward. Working effectively in their own countries seems like the top priority, with connectivity to the Internet core and to other regions with which people in the country communicate regularly being somewhat lower priorities. Just over two-thirds of ccTLDs are hosted in their own countries; refer to Figure 4 for the bigger countries, and the online appendices for the full list. Although the third of ccTLDs not hosted in their own countries include some marketed more for international use than global use—Cocos Island’s **.cc**, Tonga’s **.to**, Turkmenistan’s **.tm**, and Tuvalu’s **.tv**, among others—those are very much the exception.

Indonesia has local access to the root and to its ccTLD (**.id**). Pakistan has a root server, but no local access to its ccTLD (**.pk**). Let’s compare what happens when someone in Indonesia does a lookup on an **.id** domain name with what happens when someone in Pakistan does a lookup of a name in the **.pk** domain.

In Indonesia, the query goes to a root DNS server at the Indonesian Internet Exchange in Jakarta, where it is answered with the locations of the **.id** servers, several of which are also in Indonesia. The query then goes to the local **.id** server and is answered locally, whereupon the user can start sending traffic to the host he or she was trying to connect to, which is presumably also local. The traffic need not leave Indonesia, and if all the parties involved are in Jakarta it need not leave town.

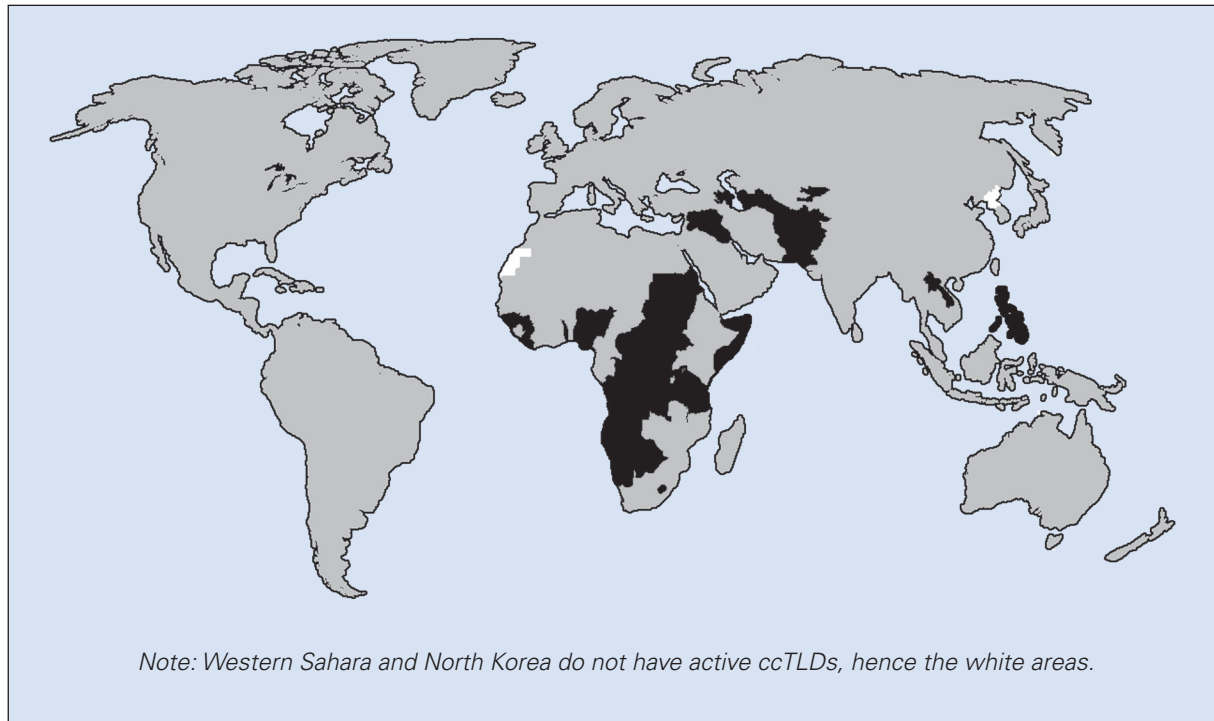
The Pakistani case is quite different. Until early 2006, there were no root servers in Pakistan, nor were there local servers for the **.pk** domain. There is now a root server in Karachi, but lacking servers for any TLDs it is of limited utility. DNS resolvers start out querying the local root server, but the response directs them to servers for the **.pk** domain, all located in the United States, at least 10 time zones away. They then send their lookup packets across the single fiber connection all the way to the United States and wait for the response. At best, this process is slow. If that fiber connection goes down, or if there is any other problem between Pakistan and these U.S. servers, local communications in Pakistan are crippled.

The situation with traditional gTLDs (**.com**, **.net**, and **.org**) in Indonesia and Pakistan is somewhat different. In Indonesia, local root servers provide addresses for the **.com**, **.net**, and **.org** servers. The **.com** and **.net** lookups can be handled in Singapore, 18 milliseconds away. Theoretically, **.org** lookups can be handled in Hong Kong, but *traceroutes* indicate **.org** queries being answered in California instead. Thus, in Indonesia, **.id** is hosted mostly locally, **.com** and **.net** are nearby, and **.org** is considerably farther away. In Pakistan, in contrast, **.pk** queries and **.org** queries are answered from the United States, more than 200 milliseconds away, while **.com** and **.net** are answered from Singapore, 80 milliseconds away. For Pakistani users of all TLDs, there are single points of failure, but **.com** and **.net** do appear to be somewhat better connected than **.pk**.

In Nairobi, Kenya, there are local copies of a root server and the local ccTLD (**.ke**). All external connectivity is by satellite, and most ISPs have only a single satellite link. Two Internet users in Nairobi wanting to communicate can do a lookup on the local root server to find the servers for **.ke** and can do a lookup on a local **.ke** server to find the servers for a subdomain of **.ke**. Assuming the subdomain being used is served locally, they can do a local lookup for a host within that subdomain and then send data across the local exchange point. Thus the two users in the same town can send data back and forth without having to send any data elsewhere.

According to Verisign, Nairobi will soon have servers for **.com** and **.net** as well. In contrast, to use the **.org** domain they can again obtain addresses of the **.org** server from their local root server, but the lookup of the **.org** domain must go over a satellite link to Europe in order to be answered by a server in London. If the satellite link is up, this process adds half a second of latency to the query. If the satellite link is down, whatever local resource they are trying to connect to is out of reach.

Figure 4: Countries that Host Their Own ccTLDs in Grey; Those that Do Not in Black



There is also a concern about ccTLDs not served from the global core; if their region or upstream provider is cut off from the Internet outside their region, the rest of the world is unable to see that ccTLD. (See Table 2). This situation may or may not be of concern; if all Internet resources within that ccTLD become unreachable in the same outage, the DNS portion of the outage may have no additional effect. However, if there is anything in that ccTLD that is not in the ccTLD's region, or if people or systems outside prefer to get a DNS response for an unreachable IP address rather than no DNS response at all, it may be of concern. Indeed, having servers that are well-connected to "the Internet as a whole" is a recommendation of RFC 2182, though the RFC does not consider the case of large portions of the Internet not being well-connected to each other.

Table 2: TLDs Not Served in the “Internet Core” Region

TLD	Country	Location of DNS Servers
BB	Barbados	Barbados
BD	Bangladesh	Bangladesh
BH	Bahrain	Bahrain
CN	China	China
EC	Ecuador	Ecuador
GF	French Guiana	French Guiana and Guadeloupe
JM	Jamaica	Jamaica
KG	Kyrgyzstan	Kazakhstan
KW	Kuwait	Kuwait
MP	Northern Mariana Islands	Guam
MQ	Martinique	Guadeloupe and Martinique
PA	Panama	Brazil, Chile, Costa Rica, and Panama
PF	French Polynesia	French Polynesia
QA	Qatar	Qatar
SR	Suriname	Suriname
TJ	Tajikistan	Tajikistan
ZM	Zambia	South Africa and Zambia

Lack of Exchange Points and Local Peering

In the “Internet Mini-Cores” article^[1], I noted that local hosting of critical infrastructure is moot if there is not either a local exchange point or a monopoly transit provider in the region. If data needed in a poorly connected region must leave the area and return to reach the user requesting it, the communication has double the latency, and possibly double the reliability problems, that it would have if it were hosted somewhere in the core. For the specific examples used in this article, I have mostly chosen areas that do have exchange points. I have not analyzed the underlying local infrastructure in all countries.

Methodology

The addresses of DNS servers for a TLD are available through several means: by looking at the root zone, by doing *digs* for the name servers, and by looking in the *Internet Assigned Numbers Authority* (IANA) *whois* data, among others. I did lookups against an anycast root server on my own network, because that seemed easiest to automate. My script then did a lookup for the address of each name server, stripped off the last octet, and produced a list of TLDs hosted in each /24 subnet.

There are 635 /24s containing name servers for TLDs; 142 of them host multiple TLDs; the rest host just one. I assumed that all DNS servers in a given /24 were likely to be in the same or nearby locations. This situation appears not to be the case for the UUNet name servers, and there are probably a few other exceptions that will show up as errors in my data.

I looked at a few automated geolocation systems to attempt to attach locations to the DNS servers, but none of them appeared to be producing accurate information. Instead, I guessed at the locations of the 600 subnets, using *traceroutes* from a variety of locations, paying attention to DNS, latency, and the results of whois queries for address space along the way. I also asked lots of questions of DNS operators and others and am particularly grateful to several anycast DNS operators, whose locations would not have all been found by my *traceroutes*. Some of my guesses are likely incorrect, and corrections are appreciated.

I may be missing some information about the UltraDNS TLD servers, because UltraDNS has locations it regards as confidential. This information about UltraDNS servers is from Afiliast's *.net* application, *traceroutes* from a variety of locations, and UltraDNS.^[7]

Locations of root servers are easier to find; they are listed at <http://www.root-servers.org>. Some supplemental information about [j.root-servers.net](http://www.j.root-servers.net) was supplied by Verisign. If there are operational root servers not included on www.root-servers.org other than the J-roots, I did not count them.

The full lists of locations of all TLDs and TLD servers are in the appendices to this article, at:

<http://www.pch.net/resources/papers/infrastructure-distribution/dns-distribution-appendices.pdf>.

References

- [1] Steve Gibbard, "Internet Mini-Cores: Local connectivity in the Internet's spur regions" (2005):
<http://www.pch.net/resources/papers/Gibbard-mini-cores.pdf>
- [2] Root Server Technical Operations Association:
<http://www.root-servers.org>
- [3] Joe Abley, "Hierarchical anycast for global service distribution," ISC Tech Notes (2003):
<http://www.isc.org/index.pl?pubs/tn/index.pl?tn=isc-tn-2003-1.html>

- [4] Bradley Huffaker, “Two days in the life of three DNS root servers” (2006):
http://www.caida.org/publications/presentations/2006/brad_wide0611_anycast_analysis
- [5] Tim Richardson, “Ship’s anchor cuts cable to Sri Lanka,” *The Register*, August 24, 2004:
http://www.theregister.co.uk/2004/08/24/sri_lanka_anchor
- [6] Tony Lee, Bradley Huffaker, Marina Fomenkov, and kc claffy, “On the problem of optimization of DNS root servers’ placement” (2003):
<http://www.caida.org/publications/papers/2003/dns-placement/>
- [7] Aflias, “.NET Application Form”:
<http://www.icann.org/tlds/net-rfp/applications/aflias.htm>

STEVE GIBBARD is a Network Architect for the nonprofit organization, Packet Clearing House (**www.pch.net**), based in Berkeley, California. He runs an anycast DNS network that hosts the top-level domains for several countries and several of the “I” root anycast DNS servers, maintains PCH’s network of route collectors and route servers at exchange points around the world, and researches the interconnection of Internet networks. In addition, Steve carries out network architecture and peering work as a consultant for several ISPs in the San Francisco Bay Area and elsewhere. Steve is a former Senior Network Engineer at Cable & Wireless, and has held network engineering positions at Digital Island and World Wide Net. E-mail: **scg@pch.net**

Writing Internet Drafts and RFCs Using XML

by Marshall T. Rose, *Dover Beach Consulting, Inc.* and Carl Malamud, *Public Resource, Inc.*

What is the work product of the *Internet Engineering Task Force* (IETF)? Some cynical observers might suggest “many fine lunches or dinners,” but we argue that those niceties are merely the means to an end. The goal of the IETF is to provide open standards for the Internet community, and those standards are memorialized as written documents called *Request For Comments* (RFCs).

In general, two organizations control the publication of documents as RFCs:

- The *Internet Engineering Steering Group* (IESG) determines which documents are suitable for publication as RFCs—typically by chartering working groups, reviewing their progress (through reading the work-in-progress *Internet Drafts*)—and ultimately approving their documents for publication.
- The RFC Editor strives for “quality, clarity, and consistency of style and format,” and has developed a particular editorial style. The latest RFC that documents this style, RFC 2223^[1], is about a decade old. A somewhat more current version can be found in a text file maintained by the RFC Editor.^[4]

For a more detailed discussion of the interaction between these two organizations, consult RFC 3932^[2].

As an organization, the IETF excels at “eating its own dog food,” including its work product: just as a protocol specification describes interactions on the wire but does not dictate the programming language used for implementation, so too, the IETF has not really cared which document preparation tools are used. The IESG worries about technical quality, and the RFC Editor worries about stylistic consistency (and, to be fair, technical quality as well). This policy works because of the careful choices made by the early Internet community, and in particular the RFC Editor, with respect to the “final form” footprint of the documents. (A discussion of these design decisions is far beyond the scope of this short article—for now, we note that it is hard to argue with success.)

An unfortunate side effect of this focus on stylistic consistency is that, for many years, the RFC Editor has had to recode documents for consistent formatting. Internally, the RFC Editor used *nroff*^[5] for this purpose, and sophisticated authors wishing to minimize RFC Editor “downtime” tended to use the same *nroff* boilerplate. The *nroff* text-formatting program has many strengths, but it can also be fairly viewed as a textual “assembly language,” with the result that authors spent a lot of time dealing with low-level formatting concerns.

In some limited cases, the high degree of formatting-specific expertise is warranted, but for the vast majority of documents, the high entry cost is not.

From Assembly Language to Markup

In early 1999 we were working at a startup company, and we needed a way to organize, search, and retrieve information from documents. We decided to use a markup language for this purpose. We also decided to use the RFC series as one of the testing grounds for the technology, because this series was one we were familiar with. Although today everyone knows what the *Extensible Markup Language* (XML) is, then there were only two widely known markup languages for authoring: SGML and HTML.

The “SG” in SGML is an abbreviation for *Standard Generalized* and not *Simple Generic*. SGML is used for the formatting of a great many books; further, it is used in large projects with long lifetimes. Although truly excellent from an “enumerate every possibility” standpoint, it has a very high cost of entry, making it difficult to use for anything other than specialized applications.

In contrast, the *Hypertext Markup Language* (HTML) embodies elegance of design, but (in the absence of *Cascading Style Sheets* [CSS]), is a presentation language, not unlike *nroff* in many respects. In other words, we needed something with the structural richness of SGML and the elegant simplicity of HTML. The newly invented XML seemed to meet the requirements.

This process led us to develop a language based on XML, which captured high-level RFC constructs (for example, authorship information) and largely ignored presentational concerns. The result is called the 2629 *format*^[3] (also known as the “xml2rfc format,” named after the initial processor for this language).

The Advantages of Markup

To understand the advantages of this approach, let’s look at one example: references. Like most archival series, the RFC Editor has a very rigorous, yet unstructured, syntax for citations. Although this consistency is good for readers of RFCs, achieving consistency of references using tools such as *nroff* was often the hardest part of creating a new document. With the 2629 format, the **<reference>** element contains a small number of subordinate elements that capture all the semantics of the reference. The XML processor takes this information and produces a properly formatted document.

Further, because this information is structured, it is possible to develop automated bibliographic databases for a wide range of data sources. In fact, using the XML “include” mechanism, a document author usually includes just a pointer to the reference, and lets the processor do all the complicated work.

A second advantage is that processors can produce different kinds of output. Some people prefer to view their documents in HTML rather than the canonical textual format. Julian Reschke has written a library of XSLT files that convert to various HTML formats (Strict, Transitional, XHTML, and so on). For example, references are hyperlinked in line, allowing for easy traversal of citations. Still others prefer the *Portable Document Format* (PDF) for printing. By using one of Julian's XSLT scripts and the truly excellent Prince^[6] XML/CSS processor, the result is high-quality, printer-ready output.

However, the primary advantage is that the “high-level” approach allows the author to focus more on content and less on format: a processor can enforce the vast majority of the esoterica associated with the RFC Editorial style, including:

- Inserting required boilerplate (and in particular, the desired revision of the boilerplate)
- Checking for mandatory sections such as “Security Considerations” or “Normative References”
- Generating a specialized table of contents, etc.

To Infinity and Beyond

After publishing RFC 2629, an unexpected result occurred: people outside the IETF started using the 2629 format for their projects. Most credit for this side effect goes to the universality of the canonical textual format. However, some authors are using the 2629 format when writing books (they convert the 2629 format to SGML, which is sent to the publisher), business plans, and software documentation—and even to create a new series of non-IETF technical documents. The constituency here seems to revolve around having a simple yet structured way to author documents.

For the last few years, a large number of XML editing programs have been deployed, and many of these support the 2629 format. These editors offer two advantages: first, they provide a natural paradigm for editing nested content; and, second, sophisticated editors can be integrated into an automated work flow. (Having said that, the authors still use *Emacs* and *vi* for their XML editors.)

A good example of the use of XML editors is a “plug-in” for the *XMLMind* Editor^[7]. This plug-in, written by Bill Fenner, provides a variety of services to the author, such as graphical editing of sections, templates for common constructs, and validation of references.

Over the last 10 years, the 2629 format has evolved in true IETF fashion, based on running code and a rough consensus. Originally created by the authors for our own convenience, we have been more than pleased to see this format used first by an informal community of developers and writers, and more recently by the IETF secretariat, tools team, and administrative entity and by the RFC Editor.

Today, many people use a common high-level markup language for writing RFCs. The next step in this natural evolution will be making the repository of XML-tagged RFCs available to those involved in document distribution, so that RFC repositories will be able to take advantage of the meta-data in the creation of search engine, alternative formats, and any other value-added constructs that would be of use to the community. (At present the RFC Editor prefers input in the 2629 format, but ultimately runs a processor that generates *nroff* for “tweaking”—in the near future, we hope that the `xml2rfc` textual output can be tuned to avoid this final step.)

To find out more, go to the `xml2rfc` Website^[8] or visit the official directory of IETF authoring tools^[9].

References

- [1] Postel, J. and J. Reynolds, “Instructions to RFC Authors,” RFC 2223, October 1997.
- [2] Alvestrand, H., “The IESG and RFC Editor Documents: Procedures,” BCP 92, RFC 3932, October 2004.
- [3] Rose, M., “Writing I-Ds and RFCs using XML,” RFC 2629, June 1999.
- [4] Reynolds, J. and R. Braden, “Instructions to Request for Comments (RFC) Authors,” August 2004.
`ftp://ftp.rfc-editor.org/in-notes/rfc-editor/instructions2authors.txt`
- [5] Ossanna, J., “Nroff/Troff User’s Manual,” UNIX Programmer’s Manual – Volume 2 (Bell Laboratories), 1979.
`http://en.wikipedia.org/wiki/Nroff`
- [6] **`http://princexml.com/`**
- [7] **`http://www.xmlmind.com/xmleditor/`**
- [8] **`http://xml.resource.org`**
- [9] **`http://tools.ietf.org/inventory/author-tools.shtml`**

CARL MALAMUD is the co-founder of Public Resource, a nonprofit public-benefit engineering firm. *Exploring the Internet* was published in 1992 as a book, but today would be called a blog. “Geek of the Week” was published in 1993 as an audio file available for download with FTP, but today would be called a podcast.
E-mail: **`carl@media.org`**

MARSHALL T. ROSE is Principal of Dover Beach Consulting, Inc. He has authored 9 books, 74 RFCs, and 4 patents. With respect to his work on the 2629 format, he claims “self defense.” E-mail: **`mrose@dbc.mtview.ca.us`**

ICANN Board Rejects .xxx Domain Application

On March 30th, 2007 the Board of the *Internet Corporation for Assigned Names and Numbers* (ICANN) voted to reject the *.xxx sponsored Top Level Domain* (sTLD) application from ICM Registry, Inc.

“This decision was the result of very careful scrutiny and consideration of all the arguments. That consideration has led a majority of the Board to believe that the proposal should be rejected,” said Dr Vint Cerf, Chairman of ICANN. “I thank my fellow Board members and the community for their input,” Dr Cerf said.

A copy of the resolution from the Board meeting is available at:

<http://www.icann.org/minutes/minutes/resolutions-30mar07.htm>

A transcript of the Board meeting is also available at:

<http://icann.org/meetings/lisbon/transcript-board-30mar07.htm>

ISOC Fellowship to the IETF

The *Internet Engineering Task Force* (IETF) is the world’s premier Internet standards setting organization. It operates as a large, open international community of network designers, operators, vendor experts, researchers, and other interested technologists. While much of the IETF’s work takes place over mailing lists, the in-person experience promotes a stronger understanding of the standardization process, encourages active involvement in IETF work, and facilitates personal networking with others that have similar technical interests.

Presently, there is limited participation at the IETF by technologists from developing countries. There are, however, many talented individuals in developing regions that have an interest in and follow IETF work and would benefit from the opportunities that attending an IETF meeting presents. As such, the main purposes of the Internet Society (ISOC)’s *IETF Fellowship Program* are to:

- Raise global awareness about the IETF and its work
- Foster greater understanding of and participation in the work of the IETF by technologists from the developing world
- Provide an opportunity for networking with individuals from around the world with similar technical interests
- Identify and foster potential future leaders from developing regions
- Demonstrate the Internet community’s commitment to fostering greater global participation in Internet Forums such as the IETF

ISOC successfully piloted the IETF Fellowship program at the 66th IETF meeting in Montreal in June 2006. Two individuals from Africa participated in this first pilot. Three individuals from the Pacific and Latin America participated in a second pilot phase at the 67th IETF meeting in San Diego in November 2006. All found the experience highly beneficial. Based on the success of the pilots, ISOC decided to formalize the program beginning in 2007.

The ISOC Fellowship pays for the Fellow's IETF meeting registration and social event fees, a round-trip economy class airfare to the meeting, hotel accommodation, and a small stipend to offset incidental expenses.

The program provides fellowships for up to five individuals per IETF meeting. ISOC will be putting out a call for candidates, including through ISOC chapters, at least 3 months before an IETF meeting. A small selection committee comprised of individuals knowledgeable about the IETF will evaluate the applicants against selection criteria and make their fellowship recommendations.

Fellowship recipients will have an obligation to present or otherwise share their experiences at the IETF meeting they attend with their local community and to provide feedback on their experience to ISOC so that the program can be continuously improved. An *ISOC Fellowship Alumni Network* will be established to extend the fellows IETF experience and relationship-building opportunities after the meeting.

For further information on the specifics of the program and how to apply for an ISOC Fellowship see:

<http://www.isoc.org/educpillar/fellowship/application.shtml>

BGP: The Movie

Statistics on Internet resources have been animated to provide a high-level overview of the consumption and use of IPv4 addresses and AS numbers since 1983. The animated video also clearly shows the effect of *Classless Interdomain Routing* (CIDR) and *Regional Internet Registries* (RIR) allocation policies on consumption rates and routing. This animation was developed by *Asia Pacific Network Information Centre* (APNIC) staff members, Geoff Huston and George Michaelson. You can download the 58MB movie from:

<http://www.apnic.net/news/hot-topics/docs/bgp-movie.mpg>

Internet Governance Articles and References

APNIC is also maintaining a collection of articles and references on Internet governance to help the community understand the issues and stay abreast of developments. You can find these at:

<http://www.apnic.net/news/hot-topics/internet-gov/index.html>

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2007 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA
--

The Internet Protocol Journal

June 2007

Volume 10, Number 2

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
AAA—Part Two	2
IPv6 Network Mobility	16
More ROAP	28
Time to Replace SMTP?	34
Fragments	39

Part One of a two-part article on *Authentication, Authorization, and Accounting* (AAA) was published in our previous issue. This time Sean Convery presents Part Two—subtitled “Protocols, Applications, and the Future of AAA.”

Interest in *IP Version 6* (IPv6) is growing in many parts of the Internet technical community; see, for example, the announcement from ARIN on page 39 of this issue. Transition to IPv6 is likely to be one of the greatest technical challenges in the history of the Internet. Several groups are developing parts of the overall solution by creating IPv6-capable versions of protocols such as the *Dynamic Host Configuration Protocol* (DHCP) or including support for IPv6 in the *Domain Name System* (DNS). Although not yet widely deployed, *IP Network Mobility* is expected to play an important part in the Internet of the future. For this reason the IETF is working on IP mobility with an eye toward IPv6. Our second article looks at the *Network Mobility* (NEMO) *Basic Support Protocol*, which is being developed by the NEMO working group in the IETF.

Depletion of IPv4 address space is not the only concern for network operators and developers these days. Questions about the long-term viability of today’s routing protocols and the associated addressing systems center around a basic concern about how we can scale our networks to a size orders of magnitude larger than what we have today. A recently formed *Routing and Addressing Problem Directorate* (ROAP) is tasked to examine these problems in detail. Several ROAP-related sessions took place during the most recent IETF meeting, and Geoff Huston reports on these sessions and gives his analysis and commentary. Incidentally, Geoff was not present in person at this IETF meeting, but the facilities to follow an IETF meeting remotely are now of such a quality that he was able to participate from the other side of the world.

Protocol replacement or enhancement is also the theme in our final article. Dave Crocker asks the question “Is it time to replace SMTP?” Since this is an opinion piece, we invite your feedback or rebuttals.

New on our Website is a linked article index. Visit cisco.com/ipj and click on “Index Files” to explore this feature.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Network Authentication, Authorization, and Accounting

Part Two: Protocols, Applications, and the Future of AAA

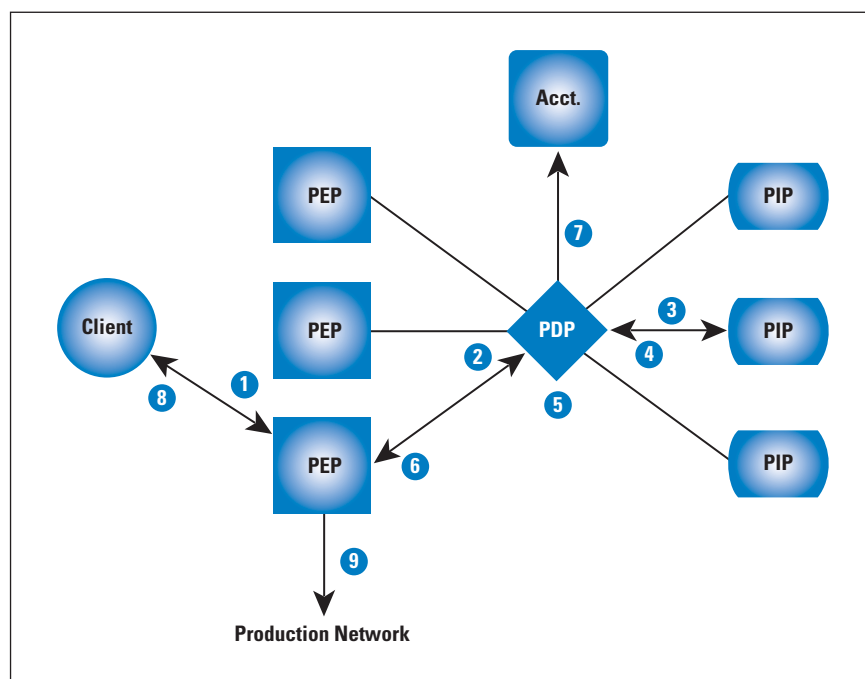
by Sean Convery, Identity Engines

Network Authentication, Authorization, and Accounting has been used since before the days of the Internet as we know it today. Authentication asks the question, “Who or what are you?” Authorization asks, “What are you allowed to do?” And finally, accounting wants to know, “What did you do?” These fundamental security building blocks are being used in expanded ways today. The first part of this two-part series focused on the overall concepts of AAA, the elements involved in AAA communications, and high-level approaches to achieving specific AAA goals. It was published in IPJ Volume 10, No. 1^[0]. This second part of the series discusses the protocols involved, specific applications of AAA, and considerations for the future of AAA.

AAA Protocols

Although AAA is often thought of as the exclusive province of the *Remote Authentication Dial-In User Service* (RADIUS) protocol, in reality a range of protocols is involved at various stages of the AAA conversation. This section introduces these AAA protocols, organized according to the parties involved in the communication. We divide AAA communications into the following categories: Client to *Policy Enforcement Point* (PEP), PEP to *Policy Decision Point* (PDP), Client to PDP, and PDP to *Policy Information Point* (PIP). For easy reference, the AAA flow diagram from Part One of this article is reproduced here. Please refer to Part One^[0] for the explanatory text associated with the diagram.

Figure 1: A Client Connects to a AAA-Protected Network
(from Part One)



Client to PEP

AAA communications between the client and the PEP can travel at Layer 2 of the OSI model, or they can run at higher layers, relying on lower layers as essentially dumb transport. The most common protocols for client-to-PEP communication are the *Point-to-Point Protocol* (PPP)^[1], *PPP over Ethernet* (PPPoE)^[2], IEEE 802.1X^[3], *IP Security* (IPsec), *Secure Sockets Layer* (SSL) VPN, and *Hypertext Transfer Protocol* (HTTP), each of which is discussed in this article.

PPP, the standard protocol for communicating across point-to-point links, includes an optional authentication step—the point at which the AAA element is introduced. During this authentication phase, protocols such as the *Challenge Handshake Authentication Protocol* (CHAP) can be used to identify the client to the PEP. (These protocols were discussed in the credential section of Part One of this article.) PPP is extensively used in dialup access but is otherwise not found in modern AAA. PPPoE, an adaptation of PPP to run over Ethernet, is used by many service providers rolling out broadband services.

PPPoE allows the broadband endpoint to authenticate itself to the service provider's network when making the initial connection. Because many broadband networks use shared Ethernet mediums, PPPoE allows *Internet Service Providers* (ISPs) to maintain the per-user accounting they were familiar with from dialup. The 802.1X protocol is an IEEE standard specifying a way to provide network access control at the port level for wired and wireless networks. The 802.1X standard specifies a way for the client to communicate with the PDP using the *Extensible Authentication Protocol* (EAP)^[4], which is discussed in more detail later in this section. The 802.1X standard requires that the endpoint support 802.1X through a “supplicant” or client sign-on application. This application authenticates the client to the network through the PEP. (See the EAP section later in this article for an explanation showing how EAP and 802.1X can work together.)

For wireless networks, 802.1X has become the standard way of authenticating clients because it supports communicating unique key material to the client to secure its use of the wireless infrastructure. In wired Ethernet networks, 802.1X is rising in popularity as a way to authenticate clients as well. These applications are more fully described in the “AAA Applications” section, later in this article.

At a more generic level, the IPsec protocol has established a standard for securing IP communications, and this approach has become another common method of communicating from a client to a PEP (referred to as a *VPN Gateway* from an IPsec perspective). The initial authentication for IPsec communications uses the *Internet Key Exchange* (IKE) protocol. Version 1^[5] of the IKE protocol had no built-in method for authenticating users with credentials such as passwords, so an extension to IKE called *XAUTH*^[6] was proposed.

XAUTH never became an official standard (though it certainly was a *de facto* one) because the IETF IPsec working group created a second version of IKE^[7] that used EAP as a transport for credentials such as passwords. Finally, in the areas of HTTP and VPN communications, the SSL and *Transport Layer Security* (TLS)^[28] standards are two closely related protocols for securing, among other things, Web communications. SSL/TLS VPNs use these protocols to create a secure session from the client to the PEP (VPN Gateway). Client authentication with SSL and TLS can be done with client-side certificates, but more commonly they use passwords or *One-Time Passwords* (OTPs).

PEP to PDP

The three main protocols for communicating between a PEP and a PDP are TACACS+^[9], RADIUS, and *Diameter*^[10]. First, consider TACACS+: Developed by Cisco, TACACS+ is a proprietary protocol that is used primarily in communicating administrator authorizations for network devices. TACACS+ uses TCP port 49 and features payload encryption for the entire TACACS+ message. Though developed by Cisco, TACACS+ is supported by other companies as well, including Juniper.

Although TACACS+ excels at command-level authorizations and accounting for administrator control, another protocol has become far more common for client AAA: RADIUS. Thanks to nearly ubiquitous support for this protocol in network hardware, RADIUS is the primary protocol for communication between a PEP and a PDP in most environments. RADIUS uses the *User Datagram Protocol* (UDP) port 1812 for authentication and authorization and UDP port 1813 for accounting^[8] (early deployments used ports 1645 and 1646, which are still used sometimes today). RADIUS supports numerous different attributes for communicating information back and forth from the PEP to the PDP, such as client MAC address, username, filter information for enforcement, and so on. It also supports an extensible framework for *Vendor-Specific Attributes* (VSAs), which allow extensions of the functions of RADIUS to support whatever elements a given PEP might need to best serve its role on the network. For example, a PEP manufacturer might support VSAs that allow the assignment of a user to a particular enforcement profile. RADIUS in its default implementation encrypts only the Password field of RADIUS messages, making the RADIUS protocol more prone to leaking information that could be used by an adversary. Both RADIUS and TACACS+ are secured by only a shared secret that is configured on both the PEP and the PDP.

Finally, consider the Diameter protocol. Diameter (the name is a play on words from RADIUS) is the next-generation, *de jure* standard for AAA. It supports stronger security through either IPsec or TLS and greater extensibility than RADIUS. It uses port 3868 for either TCP or the *Stream Control Transmission Protocol* (SCTP)^[11]. The strongest use of Diameter to date is in the carrier space, where it provides AAA for call processing and *third-generation* (3G) mobile networks.

However, the corporate market has been fairly reluctant to embrace Diameter, and that reluctance has translated into a lack of support for Diameter in corporate network infrastructure equipment.

At this point in the discussion, it makes sense to compare RADIUS and Diameter. Although Diameter is an obvious alternative, RADIUS continues to be used in both new and existing deployments, so the IETF has a working group specifically formed to extend RADIUS in the future. The relationship between RADIUS and Diameter is a little like the relationship between IPv4 and IPv6. IPv6 had IPsec as a standard feature, IPv4 integrated IPsec as well, and today, by a large margin, most IPsec deployments are on IPv4 networks. The situation is similar with AAA. RADIUS certainly had limitations, but since Diameter entered the picture, RADIUS has been extended to address some of those shortcomings, particularly with both protocols using EAP as a transport. The result is that RADIUS today does what most people want. Therefore, given the significant added complexity of Diameter, many organizations have elected not to migrate to Diameter. Both RADIUS and Diameter will be around for many years to come.

Client to PDP

Although most of the protocols in this article handle communication from one component to the next component in the AAA chain (that is, client to PEP, PEP to PDP, etc.), there is one protocol that deals with communication from the client to the PDP directly: the *Extensible Authentication Protocol* (EAP). As mentioned earlier, EAP is a flexible mechanism for communicating almost any kind of credential over almost any lower-layer transport. Each technique for authenticating a client is referred to as an *EAP Method*. Originally conceived as an extension to PPP, EAP can now use many transports, including IKEv2 and 802.1X. Cisco's proprietary *Network Admission Control* (NAC) solution offers a deployment option that puts EAP inside UDP. When using 802.1X, for example, EAP uses LAN transport, referred to as *EAPoL* (EAP over LAN). This transport is only for the connection between the client and the PEP though. From the PEP to the PDP, EAP rides inside RADIUS^[12, 13]. The actual conversation, however, takes place between the client and the PDP, with the PEP acting as a relay.

The major benefit of this approach is that the PEP does not need to understand the specifics of the EAP method selected—only the client and the PDP do. The EAP specification in the IETF specifies several different EAP methods, including *EAP Message Digest Algorithm 5* (EAP-MD5, very similar in security to CHAP), *EAP-OTP* (which supports an IETF-defined OTP solution^[14]), and *EAP Generic Token Card* (EAP-GTC). Of the methods explicitly called out in the EAP standard, EAP-GTC is the only one in much use today in production networks. EAP-GTC allows the use of OTP token cards within an EAP context.

Beyond the methods defined in the EAP standard, EAP by its nature can be extended to support additional methods. EAP *Subscriber Identity Module* (EAP-SIM)^[15] specifies a method for authentication using SIM elements in the *Global System for Mobile Communications* (GSM). EAP-SIM was developed by the *Third Generation Partnership Project* (3GPP) as a solution for these second-generation (GSM) mobile networks. EAP-AKA^[16] is the 3GPP's EAP authentication technique for third-generation (*Universal Mobile Telecommunications Service* [UMTS] or *Code Division Multiple Access 2000* [CDMA2000]) mobile networks. Both EAP-SIM and EAP-AKA support authenticating a mobile phone to a Wi-Fi network without using passwords. The problem is that without some sort of user identity federation solution in place, SIM-based authentication can work only with the mobile provider's network that supplied the SIM card. EAP-TLS^[17] specifies a technique for mutual certificate authentication. Although it is widely supported, EAP-TLS is not commonly deployed because of its requirement for client-side certificates.

Though none of the following EAP methods are standards, they—somewhat confusingly—represent the vast majority of EAP deployments. Each of them is referred to as a *Tunneled EAP Method* because it establishes one outer EAP method as a base secure channel and then runs another method (one that may be less secure) over that secure channel. *Protected EAP* (PEAP)^[18], well supported in Microsoft's Windows operating system, has become a de facto standard for EAP methods. Most clients and PDPs support PEAP today. PEAP works by establishing a TLS session authenticated by the server certificate, and then an inner authentication method rides inside that TLS session. The inner method is almost always *Microsoft CHAP Version 2* (MS-CHAPv2), but other methods can be used as well. Another popular tunneled protocol is *EAP Tunneled TLS* (EAP-TTLS)^[19]. This protocol is similar to PEAP except it supports a more arbitrary exchange of information inside the TLS tunnel. For example, one of the primary uses for EAP-TTLS is using the *Password Authentication Protocol* (PAP) as the inner authentication method, allowing an EAP-TTLS-capable PDP to authenticate clients against older password stores (such as those that support only PAP authentication).

Finally, in settings that use primarily Cisco equipment, a common tunneled protocol is *EAP Flexible Authentication via Secure Tunneling* (EAP-FAST)^[20]. This protocol uses TLS to authenticate the PDP, and then a shared key is distributed to allow faster subsequent authentication. An inner EAP method such as MS-CHAPv2 can then be used to authenticate the client to the server. EAP-FAST is used extensively in Cisco products for wireless deployments.

PDP to PIP

The final set of AAA protocols we consider are the ones that govern the communication between the PDP and the PIP. The primary protocol of interest is the *Lightweight Directory Access Protocol* (LDAP)^[21]. From a AAA context, LDAP allows a PDP to query a PIP (typically an X.500 directory^[22]) for information about a client. This information is exposed through a series of group and attribute identifiers, which can include information about a client's home location, organizational role, job title (if referring to a user), and so on. LDAP includes several different authentication options^[23]. This client information learned from the PIP enables the PDP to better make its policy decision. Also useful in the PDP-PIP communications context is the RADIUS protocol. Some large organizations or inter-organization federations use a hierarchy of RADIUS-speaking PDPs where one RADIUS PDP can act as a PIP for another RADIUS PDP further down the AAA hierarchy.

Finally, Microsoft *Active Directory* (AD) uses the LDAP protocol when acting as a PDP but also has its own extension, called *Netlogon*, for validating Microsoft credentials such as MS-CHAPv2. This means that integrating a PDP with Microsoft AD generally involves using LDAP to find information about the client and using Netlogon to validate the client's credential. Other options for PDP-to-PIP interaction—though less often used—include *Structured Query Language* (SQL) databases, *Network Information Service* (NIS), and Kerberos.

AAA Applications

This section surveys the different applications of AAA technology throughout networking. It is divided into three sections covering consumer, enterprise, and carrier applications, with a final section covering emerging applications of AAA technology.

Consumer-Managed Applications

Most consumer network deployments do not perform any advanced AAA beyond a shared key for authentication to a wireless network. In this example, the client is the consumer's host and the wireless access point acts as PEP, PDP, and PIP by validating that any client connecting to the access point presents the correct shared key.

Enterprise-Managed Applications

AAA has numerous enterprise applications, including remote access, wireless security, *Voice over IP* (VoIP), guest access, *Role-Based Access Control* (RBAC), and endpoint posture validation (also known as NAC). This section discusses each of these applications. Remote-access security is the original enterprise AAA application. In the remote-access scenario, remote users connect over a dialup connection or a VPN and authenticate themselves (and optionally their hosts) to the organization's network.

The client's credential is almost always a password, expressed in one of the forms discussed in the credential section of Part One of this article. The main purpose of AAA in the remote-access case is to validate that the client is a valid user of the organization's network.

Wireless security is similar in some respects to remote-access security. The goals of AAA in wireless security are twofold: first it must validate that the wireless client is an authorized user, and second, it must provide the client with a session key for cryptographic protection of the client's traffic. Given these goals, 802.1X using EAP are the ideal protocols to use because they support both client authentication and dynamic keying. Older wireless security approaches relied on an open wireless network and a VPN Gateway separating that network from the rest of the organization's network. In that example, the wireless-security approach mimics the remote-access application just discussed. Other types of networking require different applications of AAA. For example, VoIP deployments have authentication requirements as well. The *Session Initiation Protocol* (SIP)^[24] is used extensively for, well, session initiation in VoIP networks (for example, authenticating the calling parties prior to initiating a new call). Authentication can be handled natively within SIP using HTTP digest authentication, or the same request can be sent to a PDP using RADIUS. AAA for VoIP allows handsets to authenticate themselves to the network and gain access to call-processing services.

Another, very popular application of AAA is guest-access management for networks. This application has grown quickly with the recent growth of wireless networks. Guest access is a method by which guests can be granted temporary access to a network with a full audit trail^[27]. Guest access generally involves creating a distinct PIP, which houses short-term user accounts, and a technique for creating and, after a configurable period of time, automatically deactivating those user accounts. The PIP is often co-resident with the PDP and allows this temporary access without having to provision these users into the organization's more permanent directory. The guest can communicate with the PEP using any of the client-PEP protocols discussed earlier, though HTTP is the most common. The PEP is told by the PDP that the client (because it is a guest) should have restricted access—typically access only to the Internet at large and not any communication with an organization's internal network.

Also growing in popularity as a AAA application is RBAC, an application of AAA that allows customization of the network session based on the role of the client. In fact, guest access is a simple form of RBAC whereby two classes of clients are created: guest and permanent. However, RBAC can be extended to include more levels of delineation, including guest, contractor, and specific classes of permanent users such as sales, human resources, and engineering.

This classification can be done with all forms of AAA-enabled network infrastructure, including wired, wireless, and remote access. Current scalability limitations of VLAN technology and *Access Control Lists* (ACLs) make creating large quantities of roles difficult, but a significant business benefit in audit and regulatory compliance can be realized with usually fewer than five roles.

To implement RBAC, most organizations choose a mix of 802.1X and HTTP authentication for wired and wireless access, combined with VPN technology for remote access. This approach is the most common one to RBAC, though others are used.

Finally, another important AAA application is *Endpoint Posture Validation*, also referred to as *Network Access Control* (NAC). Unfortunately NAC is an inappropriate name because of its almost complete overlap with the more general AAA term—leading to a fair amount of confusion in the market. Endpoint posture validation refers to many different parameters in the industry as it is an emerging technology. These parameters range from very narrow device-centric posture checking to a more identity-centric approach for secure mobile computing. Because this entire article is concerned with the latter, we will consider NAC in its narrow context of endpoint posture checking. With this label, NAC simply acts as another PIP for the PDP to use.

This time, though, instead of checking the client's credential, NAC checks the client's software configuration. This checking generally focuses on security-sensitive configuration details of the endpoint security software and the operating system itself, such as the revision, configuration, and current operating status. This client configuration data is gathered by a host agent on the client and then sent to the PDP or PIP for evaluation. The host agent is either permanent on the client or downloaded dynamically to acquire the information. Some NAC applications rely exclusively on external scanning of the client, although this scanning generally yields far less granular information than an agent would.

The challenge with NAC today is deploying a system built on standards. The IETF and the *Trusted Computing Group* (TCG) are both pursuing standards in this space. Meanwhile Cisco, Microsoft, and a host of smaller companies have offerings not currently based on any standard. Recent announcements from the TCG and Microsoft are changing this. The TCG recently standardized the as-implemented NAC protocol used by Microsoft's NAC approach. Though there is much more work to do, this should allow the beginnings of standards-based interoperability in NAC solutions since a core protocol in Microsoft's NAC is now a standard from the TCG. There is a great base in standards at a low enough layer in all the NAC approaches though, as the emerging standards use the protocols discussed in this article including 802.1X, IPsec, RADIUS, and LDAP.

Carrier-Managed Applications

Some carrier-managed AAA applications are similar to those for the enterprise and others are different. The common distinctions for almost all carrier applications are their large scale and their emphasis on accounting. Carrier applications include dialup, DSL or cable PPPoE, mobile or 3G, wireless hotspot, and metro wireless. Dialup is similar to the remote-access application in the enterprise section, but on a massive scale. *Network Access Servers* (NASs) for a large ISP are geographically dispersed, as are the PDP and PIP systems that support them. Clients communicate with the PEP (NAS) with PPP using one of the password credential techniques discussed in Part One of this article, and the PEP communicates with the PDP using RADIUS or Diameter.

Now consider DSL or cable PPPoE. Though PPPoE-based broadband access seems to be on the decline, many ISPs are still using PPPoE for the enhanced audit trail it provides compared with an unauthenticated connection. In the realm of mobile telephone networks, service providers are increasingly providing data services in mobile phones, and these services require AAA for security and billing. Such data services come in several varieties on both the second- and third-generation mobile networks. Additionally, smartphones are increasingly supporting 802.11-based wireless access as well, creating a complex relationship between the smartphone, mobile voice network, mobile data network, 802.11 data network, and VoIP-based voice services. Previously discussed standards such as EAP-SIM and EAP-AKA are trying to bridge some of these worlds, but there is much work to be done. Ideally, any smartphone should take advantage of the network with the fastest and richest set of services, and callers trying to reach a smartphone user as well as the user himself, should be shielded from this discovery and association process. Business motivators and detractors within the carrier space may affect this convergence.

The next carrier-managed AAA application to discuss is the *wireless hotspot*. Hotspots work much like dialup providers in that regular users get a password-based credential that lets them authenticate to the hotspot. In this context, the 802.1X protocol is less commonly used because the required client software is not yet ubiquitous in the client install base. More common is Web-based authentication much like that used to access broadband in a hotel. A critical security consideration for a hotspot operator is the ability to ensure that a given client is not connected to two hotspots at the same time—a situation that would indicate an account was shared between two or more users. This stipulation places an increasing burden on the accounting aspect of AAA, as with any carrier-based AAA application.

Finally, the last AAA application we examine is the metropolitan wireless network, known as “metro wireless.” In metro wireless, an 802.11 network is deployed throughout a metropolitan area, and access is provided free of charge or for a fee. I live in Mountain View, California, which is home to Google’s headquarters, and is where Google has installed its free, citywide metro wireless network.

Although the service is free, AAA is still required: to sign on to the wireless network, you must authenticate to Google using an ID. This step, much like signing on to a wireless hotspot, allows Google to trace network use to an individual (if necessary) and switch to a fee-based model later on if desired. HTTP authentication is most common in metro wireless environments, and, because of the on/off nature of access, little sophistication in policy decision is required other than validating the client's credential.

Emerging Applications

Several interesting applications of network AAA are emerging. The first is in building just-in-time networks, such as when establishing an on-scene emergency operations center after a disaster. In this situation, emergency workers often need to communicate in a protected environment, and the press that covers the disaster needs network access to send in its reports. The AAA application required here is a cross between wireless security, guest access, and RBAC.

Another emerging application is what we call "granular RBAC." As opposed to RBAC, which associates users into coarse-grained classes of users, granular RBAC knows much more about the users and makes a more sophisticated access decision.

One example of the use of granular RBAC is for classroom control in higher education. Increasingly, classrooms are wireless-enabled as a convenience feature for faculty and students. However, during exam time it is often useful to disable this access to the students taking an exam. Without a granular understanding of which clients are connecting to the network, this setup is very difficult to achieve without physically disabling large portions of the wireless network during exam time. By using AAA, a school could put class schedules inside an LDAP store along with the rest of the students' information. Professors could also register exam times by time and location. AAA could then prevent students from getting on the network inside the classroom during their exam period, while still letting them connect to the network when inside their dorm room.

Finally, the last application we consider is what I call "punitive access restrictions." As networks become more and more an integral part of our lives, it is natural to want as fast a network connection as we can find, creating the situation where denying access to the network based on past behavior (network related or not) can be used as a punitive action. Today, your driver's license can be revoked based on your behavior while on the road. Punitive access restrictions on the network could mirror the same technique (for example, punishing people who propagate a virus by restricting their network access for a time) or could be used even if the infraction is not related to the network. Imagine a university that has trouble getting students to return overdue library books. Fines are one way to get the books back quickly, but if the student's parents are paying the bill, this consequence may not be as effective as the university desires.

However, imagine if the student's account record (in the PIP) had a directory attribute containing a count of the student's overdue library books. The network could then use RBAC or *Quality-of-Service* (QoS) techniques to provide degraded access to the student until the books were returned.

The Future of AAA

AAA as a concept has remained relatively unchanged since its inception. However, as this article has demonstrated, the techniques and applications of AAA continue to evolve. This section discusses some of the ways AAA may change more fundamentally in the future.

Security and Identity Convergence

Today the security and identity services provided by physical building access, network access, and application access are completely distinct. Security can be improved by communicating among these layers. Imagine a user executing a \$10 million purchase order in a financial application. The chance of fraud would be reduced if the application could know that the user was coming from an authorized client with an up-to-date antivirus configuration. The chance of fraud could be further reduced by checking that the same user had accessed the badge access system of the building that day, and that the point of badge-access entry was consistent with the location where the application request originated. Within computer security, the notion of *defense-in-depth* has been around for a long time and is considered a best practice. Security and identity convergence adds new layers to this defense, and can potentially make all the layers more intelligent in their interaction.

User-Centric AAA

In the Web application world, the notion of user-centric identity is gaining ground. Kim Cameron's "Laws of Identity"^[25] makes a compelling case that identity information housed in silos to be used by one organization is problematic. Several circles in the Web and e-commerce communities are beginning to look at identity differently. One change, consistent with the notion of user-centric identity, is that users should own their own identity information and should control how that information is used. The simplest example I can offer is shopping preferences at an online store. Most online stores make suggestions to you based on prior purchases. This data is owned by the online store, though, and not you, the consumer. If you wanted to take your purchasing profile from Amazon.com and transfer it to Barnes and Noble, it would not work. With user-centric identity, this kind of process is possible.

Another example is asserting a user's age. Depending on what you are trying to do on the Internet, you may need to validate that you are above a certain age. To do that, you are often asked to enter your date of birth, but that is more information than the site really needs.

If you could assert, with an identity you control but that is validated by a trusted party, that you are over the required age, it would not be necessary to disclose your date of birth (a process that is sometimes used as an authentication factor when you call places such as your credit card company).

The idea here is that you control your own information and limit what you need to share with others. This is very beneficial for privacy. One user-centric identity approach is included in Windows Vista through an application called “Card Space.” Other approaches include OpenID and the Higgins Project. All of these approaches are somewhat consumer-focused, but if they take hold, it seems natural that there will be pressure for similar identity approaches in the enterprise and carrier space.

Federation

One of the natural evolutions of AAA infrastructure is to start federating access between multiple organizations. Imagine if visiting professors at another university’s campus could access the network as guests using their password from their home location? Federation promises to make this possible, but the most challenging hurdles are political and logistical rather than technological. Protocols such as the *Security Assertion Markup Language* (SAML)^[26] combined with RADIUS and LDAP can overcome this hurdle. The challenge is how to set up the trust relationships between the organizations to make it work. Eduroam based on RADIUS is an early effort delivering federation in Europe today.

Summary

This article, with its companion piece, has explored all aspects of AAA. Part One described the overall approach of AAA, how it works, and the elements that provide authentication, authorization, and accounting. Part Two has explored all the protocols used in the communication between the various AAA elements, the applications of AAA, and some thoughts about the future of AAA. AAA is a giant topic, and each of these sections, protocol descriptions, and applications could be expanded into a paper all by itself. The information in this article, combined with the references provided, should be a good starting point for your own examination of the specific aspects of AAA that are of interest to you.

References

- [0] Convery, S., “Network Authentication, Authorization, and Accounting —Part One: Concepts, Elements, and Approaches,” *The Internet Protocol Journal*, Volume 10, No. 1, March 2007.
- [1] Simpson, W., “The Point-to-Point Protocol (PPP),” RFC 1661, July 1994.

- [2] Mamakos, L., “A Method for Transmitting PPP Over Ethernet (PPPoE),” RFC 2516, February 1999.
- [3] Jeffree et al., “Port-Based Network Access Control,” IEEE Std 802.1X-2004, November 2004.
- [4] Aboba et al., “Extensible Authentication Protocol,” RFC 3748, June 2004.
- [5] Harkins et al., “The Internet Key Exchange (IKE),” RFC 2409, November 1998.
- [6] Beaulieu et al., “Extended Authentication within IKE (XAUTH),” Internet Draft, Work in Progress, October 2001.
draft-beaulieu-ike-xauth-02.txt
- [7] Kaufman C., ed., “Internet Key Exchange (IKEv2) Protocol,” RFC 4306, December 2005.
- [8] Rigney C., “RADIUS Accounting,” RFC 2866, June 2000.
- [9] Carrel et al., “The TACACS+ Protocol Version 1.78,” Internet Draft, Work in Progress, January 1997.
draft-grant-tacacs-02.txt
- [10] Calhoun et al., “Diameter Base Protocol,” RFC 3588, September 2003.
- [11] Stewart et al., “Stream Control Transmission Protocol,” RFC 2960, October 2000.
- [12] Aboba et al., “RADIUS (Remote Authentication Dial In User Service) Support For Extensible Authentication Protocol (EAP),” RFC 3579, September 2003.
- [13] Congdon et al., “IEEE 802.1X Remote Authentication Dial In User Service (RADIUS) Usage Guidelines,” RFC 3580, September 2003.
- [14] Haller et al., “A One-Time Password System,” RFC 2289, February 1998.
- [15] Haverinen et al., “Extensible Authentication Protocol Method for Global System for Mobile Communications (GSM) Subscriber Identity Modules (EAP-SIM),” RFC 4186, January 2006.
- [16] Arkko et al., “Extensible Authentication Protocol Method for 3rd Generation Authentication and Key Agreement (EAP-AKA),” RFC 4187, January 2006.

- [17] Aboba et al., “PPP EAP TLS Authentication Protocol,” RFC 2716, October 1999.
- [18] Palekar et al., “Protected EAP Protocol (PEAP) Version 2,” Internet Draft, Work in Progress, October 2004.
draft-josefsson-pppext-eap-tls-eap-10.txt
- [19] Funk et al., “EAP Tunneled TLS Authentication Protocol Version 1 (EAP-TTLSv1),” Internet Draft, Work in Progress, March 2006. **draft-funk-eap-ttls-v1-01.txt**
- [20] Cam-Winget et al., “The Flexible Authentication via Secure Tunneling Extensible Authentication Protocol Method (EAP-FAST),” Internet Draft, Work in Progress, January 2007.
draft-cam-winget-eap-fast-06.txt
- [21] Zeilenga K., “Lightweight Directory Access Protocol (LDAP): Technical Specification Road Map,” RFC 4510, June 2006.
- [22] Zeilenga K., “Lightweight Directory Access Protocol (LDAP): Directory Information Models,” RFC 4512, June 2006.
- [23] Harrison R., “Lightweight Directory Access Protocol (LDAP): Authentication Methods and Security Mechanisms,” RFC 4513, June 2006.
- [24] Rosenberg et al., “SIP: Session Initiation Protocol,” RFC 3261, June 2002.
- [25] Cameron, “The Laws of Identity,” May 2005.
- [26] OASIS, “Security Assertion Markup Language 2.0,” March 2005.
- [27] Dory Leifer, “Visitor Networks,” *The Internet Protocol Journal*, Volume 5, No. 3, September 2002.
- [28] William Stallings, “SSL: Foundation for Web Security,” *The Internet Protocol Journal*, Volume 1, No. 1, June 1998.

SEAN CONVERY is CTO at Identity Engines, a venture-backed startup developing innovative identity management solutions for enterprise networks. Prior to Identity Engines, Sean (CCIE® no. 4232) worked for seven years at Cisco Systems, most recently in the office of the security CTO. Sean is best known as the principal architect of the SAFE Blueprint from Cisco and the author of *Network Security Architectures* (Cisco Press, 2004). Sean has presented to or consulted with thousands of enterprise customers around the world on designing secure networks. Before Cisco, Sean held various positions in IT and security consulting during his 14 years in networking. E-mail: **sconvery@idengines.com**

IPv6 Network Mobility

by Carlos J. Bernardos, Ignacio Soto, and María Calderón, Universidad Carlos III de Madrid

The *Internet Protocol* (IP) is currently accelerating the integration of voice and data communications. The Mobile IP protocol enables host mobility support, but several scenarios exist today, such as the provision of Internet access from mobile platforms (for example, planes, trains, cars, etc.), making it necessary to also support the mobility of complete networks. In response to this demand, the *Internet Engineering Task Force* (IETF) has developed the *Network Mobility* (NEMO) *Basic Support Protocol*^[1], enabling IPv6 network mobility.

This article explains the Network Mobility Basic Support Protocol, by first providing a general overview and then examining the details.

Why Network Mobility?

Accelerated by the success of cellular technologies, mobility has changed the way people communicate. As Internet access becomes more and more ubiquitous, demands for mobility are not restricted to single terminals anymore. It is also needed to support the movement of a complete network that changes its point of attachment to the fixed infrastructure, maintaining the sessions of every device of the network: what is known as *network mobility* in IP networks. In this scenario, the mobile network has at least a (mobile) router that connects to the fixed infrastructure, and the devices of the mobile network connect to the exterior through this mobile router.

Support of the roaming of networks that move as a whole is required in order to enable the transparent provision of Internet access in mobile platforms, such as the following:

- *Public transportation systems*: These systems would let passengers in trains, planes, ships, etc. access the Internet from terminals onboard (for example, laptops, cellular phones, *Personal Digital Assistants* [PDAs], and so on) through a mobile router located at the transport vehicle that connects to the fixed infrastructure.
- *Personal networks*: Electronic devices carried by people, such as PDAs, photo cameras, etc. would connect through a cellular phone acting as the mobile router of the personal network.
- *Vehicular scenarios*: Future cars will benefit from having Internet connectivity, not only to enhance safety (for example, by using sensors that could control multiple aspects of the vehicle operation, interacting with the environment and communicating with the Internet), but also to provide personal communication, entertainment, and Internet-based services to passengers.

However, IP networks were not designed for mobile environments. In both IPv4^[2] and IPv6^[3, 4], IP addresses play two different roles. On the one hand, they are *locators* that specify, based on a routing system, how to reach the node that is using that address. The routing system keeps information about how to reach different sets of addresses that have a common network prefix. This address aggregation in the routing system satisfies scalability requirements. On the other hand, IP addresses are also part of the *endpoint identifiers* of a communication, and upper layers use the identifiers of the peers of a communication to identify them. For example, the *Transmission Control Protocol* (TCP), which is used to support most of the Internet applications, uses the IP address as part of the TCP connection identifier.

This dual role played by IP addresses imposes some restrictions on mobility, because when a terminal moves from one network (IP subnet) to another, we would like to *maintain* the IP address of the node that moves (associated to one of its network interfaces) in order not to change the identifier that upper layers are using in their ongoing sessions. However, we also would like to *change* the IP address to make it topologically correct in the new location of the terminal, allowing in this way the routing system to reach the terminal.

Protocols such as the *Dynamic Host Configuration Protocol* (DHCP)^[5, 6] facilitated the portability of terminals by enabling the dynamic acquisition of IP configuration information without involving manual intervention. However, this automation is not enough to achieve real and transparent mobility because it requires the restarting of ongoing transport sessions after the point of attachment changes. The IETF has studied the problem of terminal mobility in IP networks for a long time, and IP-layer solutions exist for both IPv4 (Mobile IPv4^[7, 8]) and IPv6 (Mobile IPv6^[9]) that enable the movement of terminals without stopping their ongoing sessions.

If we focus on IPv6^[3] networks, Mobile IPv6 does not support, as it is now defined, the movement of complete networks. One way of achieving the transparent mobility of all the nodes of a network moving together (for example, in a plane) could be enabling host mobility support in all of them, so they independently manage their mobility. However, this approach has the following drawbacks:

- Host mobility support (for example Mobile IP^[7, 8, 9]) is required in *all* the nodes of the network. This support might not be possible, for example, because of the limited capacities of the nodes (such as in sensors or embedded devices) or because it is not possible to update the software in some older devices. By having a single entity (the mobile router) that manages the mobility of the complete network, nodes of the network do not require any special mobility software to benefit from the transparent mobility support provided by the (mobile) router.

- The signaling exchanged because of the roaming of the network is limited to a single node sending only one message (avoiding “storms” of signaling messages every time the network moves).
- Nodes of the network must be able to attach to the access technology available to connect to the Internet. This requirement might mean that all the nodes of the network should have *Universal Mobile Telecommunications Service* (UMTS) or WiMAX interfaces, for example. On the other hand, by putting this requirement on a single node (the mobile router), nodes of the network can gain access to the Internet through the mobile router, using cheaper and widely available access technologies (for example, *wireless LAN* [WLAN] or Bluetooth).

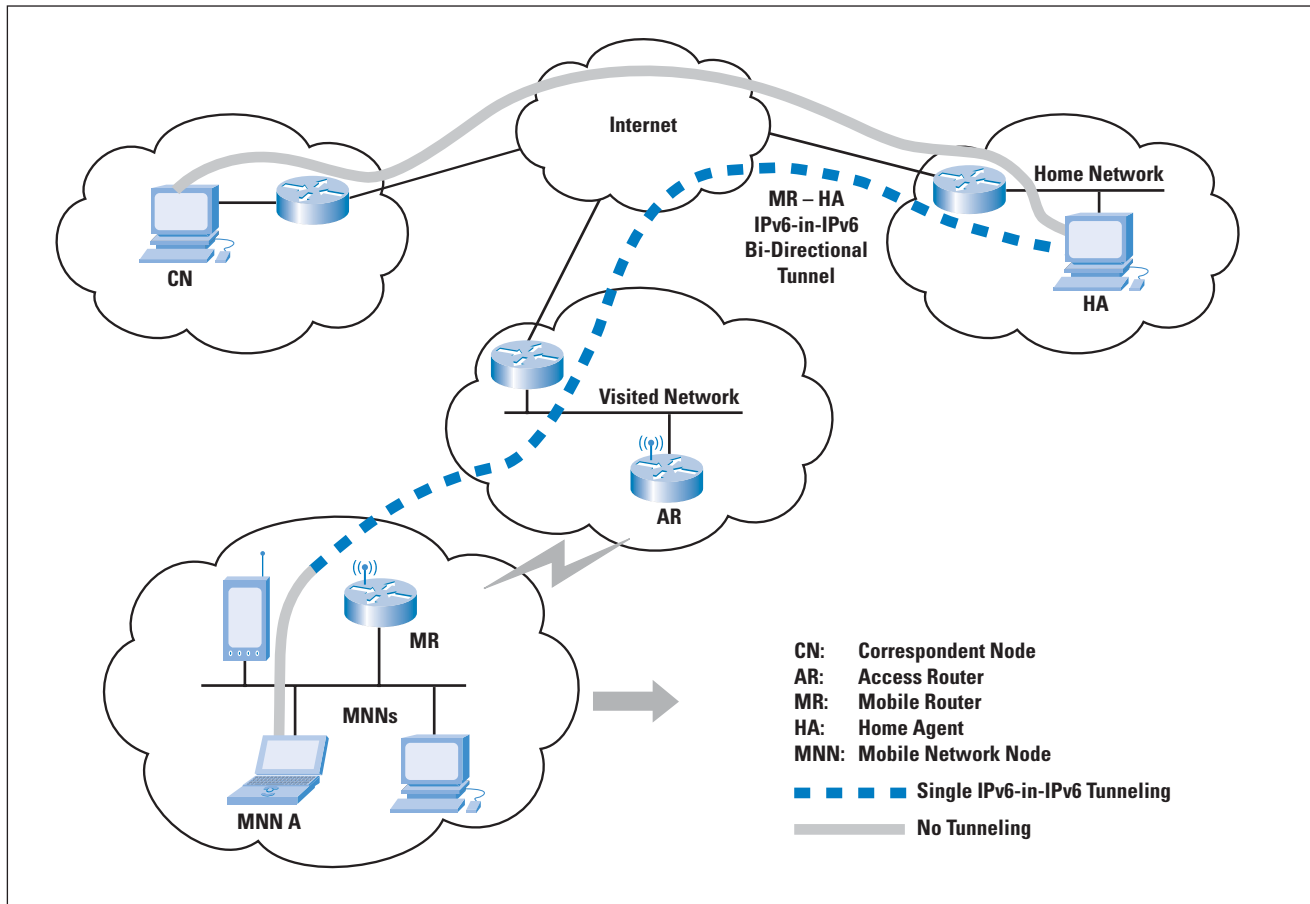
Because of these problems, the IETF *NEMO Working Group* was created to standardize a solution enabling network mobility at the IPv6 layer. The current solution, called the Network Mobility Basic Support Protocol, is defined in RFC 3963^[1].

Operation of the NEMO Basic Support Protocol

A mobile network (known also as a “network that moves,” or *NEMO*) is defined as a network whose attachment point to the Internet varies with time. Figure 1 depicts an example of a network-mobility scenario. The router within the NEMO that connects to the Internet is called the *Mobile Router* (MR). It is assumed that the NEMO is assigned to a particular network, known as its *Home Network*, where it resides when it is not moving. Because the NEMO is part of the home network, the mobile network has configured addresses belonging to one or more address blocks assigned to the home network: the *Mobile Network Prefixes* (MNPs). These addresses remain assigned to the NEMO when it is away from home. Of course, these addresses have topological meaning only when the NEMO is at home. When the NEMO is away from home, packets addressed to the nodes of the NEMO, known as *Mobile Network Nodes* (MNNs), are still routed to the home network. Additionally, when the NEMO is away from home, the mobile router acquires an address from the visited network, called the *Care-of Address* (CoA), where the routing architecture can deliver packets without additional mechanisms.

When any node located at the Internet, known as a *Correspondent Node* (CN), exchanges IP datagrams with a *Mobile Network Node* (MNN; A in Figure 1), the following operations are involved in the communication:

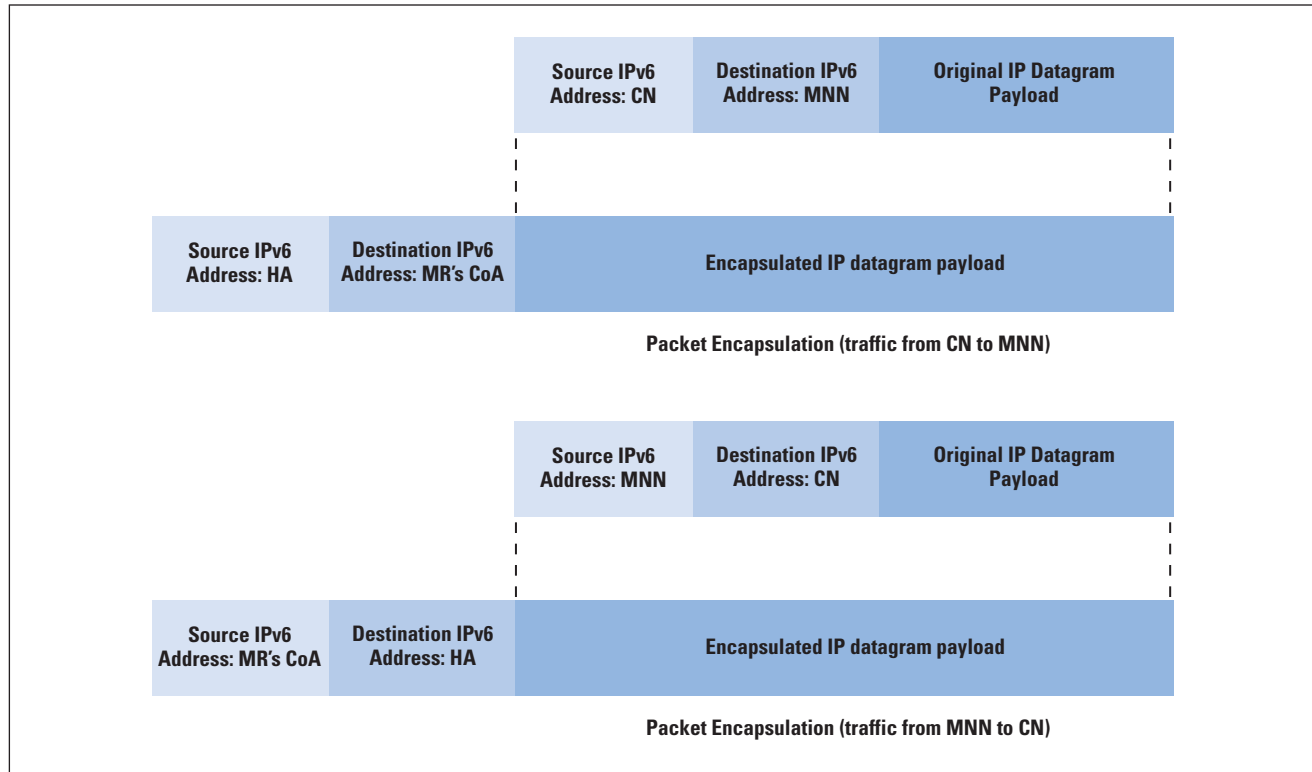
Figure 1: Example of NEMO Basic Support Protocol Operation



1. The correspondent node transmits an IP datagram destined for MNN A. This datagram carries as its destination address the IPv6 address of MNN A, which belongs to the MNP of the NEMO.
2. This IP datagram is routed to the home network of the NEMO, where it is encapsulated inside a new IP datagram by a special node located on the home network of the NEMO, called the *Home Agent* (HA). The new datagram is sent to the CoA of the mobile router, with the IP address of the home agent as source address. This encapsulation (as shown in Figure 2) preserves mobility transparency (that is, neither MNN A nor the correspondent node are aware of the mobility of the NEMO) while maintaining the established Internet connections of the MNN.
3. The mobile router receives the encapsulated IP datagram, removes the outer IPv6 header, and delivers the original datagram to MNN A.

4. In the opposite direction, the operation is analogous. The mobile router encapsulates the IP datagrams sent by MNN A toward its home agent, which then forwards the original datagram toward its destination (that is, the correspondent node). This encapsulation is required to avoid problems with ingress filtering, because many routers implement security policies that do not allow the forwarding of packets that have a source address that appears topologically incorrect.

Figure 2: Overview of NEMO Basic Support Protocol Encapsulation

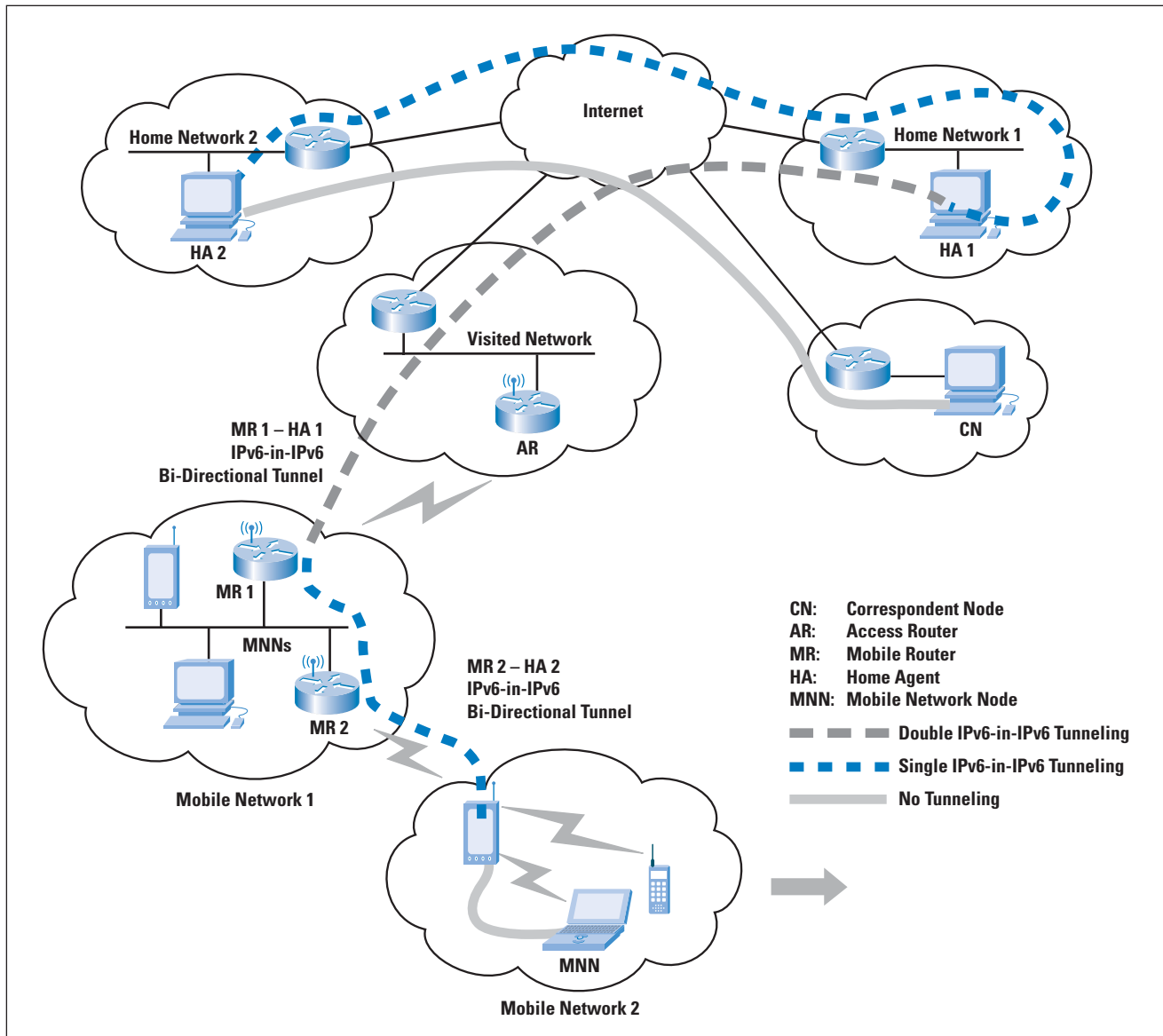


Following are different types of MNNs:

- *Local Fixed Node* (LFN): This node has no mobility-specific software and therefore cannot change its point of attachment while maintaining ongoing sessions. Its IPv6 address is taken from a MNP of the NEMO to which it is attached.
- *Local Mobile Node* (LMN): This node implements the Mobile IPv6 protocol; its home network is located in the mobile network. Its *home address* (HoA) is taken from an MNP.
- *Visiting Mobile Node* (VMN): This node implements the Mobile IP protocol (and therefore, it can change its point of attachment while maintaining ongoing sessions), has its home network outside the mobile network, and it is visiting the mobile network. A VMN that is temporarily attached to a mobile subnet (used as a foreign link) obtains an address on that subnet (that is, its CoA is taken from an MNP).

Additionally, mobile networks can be *nested*. A mobile network is said to be nested when it attaches to another mobile network and obtains connectivity through it (refer to Figure 3). An example is a user who enters a vehicle with his personal area network (mobile network 2) and connects, through a mobile router—like a Wi-Fi enabled PDA—to the network of the car (mobile network 1), which is connected to the fixed infrastructure.

Figure 3: Nested Mobile Network: Operation of the NEMO Basic Support Protocol (multiangular routing)



Protocol Details: NEMO Versus Mobile IPv6

The NEMO Basic Support Protocol is an extension of the solution proposed for host mobility support, *Mobile IPv6* (MIPv6)^[9].

In Mobile IPv6, three mechanisms support the mobility of a host: movement detection, location registration, and traffic tunneling. The NEMO Basic Support Protocol extends some of these mechanisms to support the movement of complete networks. These mechanisms are described next, with those parts that are different from the Mobile IPv6 protocol highlighted.

Movement Detection

In Mobile IPv6, the host needs to discover its own movement, so it can proceed with the required signaling and operations that allow its transparent mobility. Mobile IPv6 defines a generic movement-detection mechanism based on the *Neighbor Discovery Protocol*^[10], which basically consists of listening to *Router Advertisements* (RAs). Routers send these router-advertisement messages, both periodically and in response to a *Router Solicitation* message issued by a host. By looking at the information contained in the router advertisements, a host can determine whether or not it has moved to a new link.

The NEMO Basic Support Protocol does not introduce any change on the movement-detection mechanisms that a mobile router can use.

Location Registration

When a host moves to a new network, it has to configure a new IPv6 address on the visited link (belonging to the IPv6 address space of that visited network): the CoA, and inform the home agent of the movement. In Mobile IPv6, the mobile node (that is, a mobile host) informs its home agent of its current CoA using a mobility message called the *Binding Update* (BU). This message is carried in an IPv6 datagram using a special extension header defined by Mobile IPv6 to encapsulate all messaging related to the creation and management of mobility bindings, called the *mobility header*. The binding-update message contains information required by the home agent to create a mobility binding, such as the home address of the *Mobile Node* (MN) and its CoA, where the home agent should encapsulate all the traffic destined to the mobile node. The home agent replies to the mobile node by returning a *Binding Acknowledgement* (BA) message.

The NEMO Basic Support Protocol extends the binding-update message to convey the following additional information:

- *Mobile Router Flag* (R): The mobile router flag is set to indicate to the home agent that the binding update is from a mobile router. A mobile router can behave as a mobile host: by setting this flag to 0, the home agent does not forward packets destined for the mobile network to the mobile router, but forwards only those packets destined to the home address of the mobile router.

- *Mobile Network Prefix Option:* This option is in the binding update to indicate the prefix information for the mobile network to the home agent. There could be multiple mobile network prefix options if the mobile router has more than one IPv6 prefix in the mobile network and wants the home agent to forward packets for each of these prefixes to the current location of the mobile router.

When the NEMO Basic Support Protocol is used to provide mobility to a complete network, only one binding-update or binding-acknowledgement signaling messages exchange is performed, whereas if the Mobile IP protocol were used by all the nodes of an N -node network, $N \times$ (Binding-update or Binding-acknowledgement) signaling messages synchronized exchanges would be required—usually referred to as a “binding-update signaling storm.”

Mobile IPv6 defines a route-optimization mechanism that enables direct path communication between the mobile node and a correspondent node (avoiding traversal of the home agent). This route optimization is achieved by allowing the mobile node to send binding-update messages also to the correspondent nodes. In this way the correspondent node is also aware of the CoA, where the home address of the mobile node is currently reachable. A special mechanism—called the *Return Routability* (RR) procedure—is defined to prove that the mobile node has been assigned (that is, “owns”) both the home address and the CoA at a particular moment in time^[11], and therefore provides the correspondent node with some security guarantees.

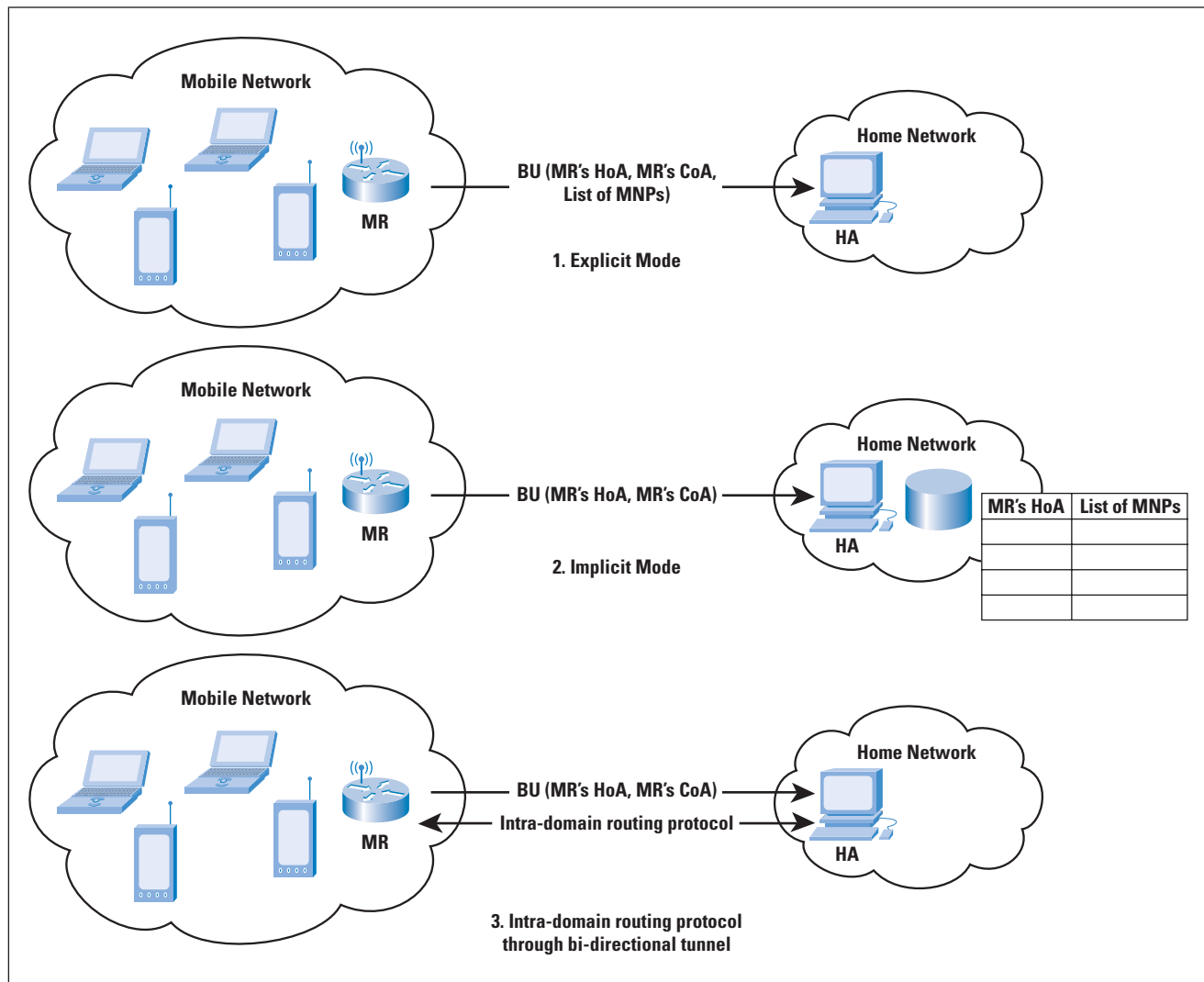
Because of the nature of the network-mobility scenario, the task of providing mobile networks with route-optimization support becomes more complex. The IETF is currently working on this topic^[12, 13, 14].

Traffic Tunneling

In Mobile IPv6, after the mobile node has successfully registered its current location, the home agent starts encapsulating the data traffic destined to the mobile node toward its CoA.

In a NEMO scenario, the home agent forwards not only those IP datagrams arriving at the home network that are destined to the home address of the mobile router, but also all the traffic addressed to any of the mobile-network prefixes managed by the mobile router. The home agent can determine which prefixes belong to the mobile router in three different ways (refer to Figure 4):

Figure 4: NEMO Basic Support Modes of Operation



- *Explicit mode:* The mobile router includes one or more mobile network prefix options in the binding-update message that it sends to the home agent. These options contain information about the mobile-network prefix(es) configured on the mobile network.
- *Implicit mode:* The mobile router does not include prefix information in the binding-update message it sends to the home agent. The home agent determines the mobile-network prefix(es) owned by the mobile router by using any other mechanism (the NEMO Basic Support Protocol does not define any, leaving this prefix determination open to be implementation-specific).

One example would be manual configuration at the home agent mapping the home address of the mobile router to the information required for setting up forwarding for the mobile network.

- *Intradomain Dynamic Routing Protocol through the bidirectional tunnel*: Alternatively to the previous two modes of operation, the home agent and the mobile router can run an intradomain routing protocol (for example, *Routing Information Protocol next generation* [RIPng] or *Open Shortest Path First* [OSPF]) through the bidirectional tunnel. The mobile router can continue running the same routing protocol that it ran when attached to the home link.

Fragmentation may be needed to forward packets through the tunnel between the mobile router and the home agent. In this case, the other end of the tunnel (the home agent of the mobile router) must reassemble the packet before forwarding it to the final destination. This requirement does not contradict the fact that *intermediate* IPv6 routers do not fragment (as opposed to IPv4), because the mobile router and home agent are the actual *ends* of the tunnel.

Performance of the NEMO Basic Support Protocol

The NEMO Basic Support Protocol relies on the creation of a bidirectional tunnel between the mobile router and the home agent to provide transparent mobility support to a complete network. The use of this tunnel causes an additional overhead of 40 bytes per packet, because of the extra IPv6 header added by the encapsulation. The effect of this overhead might be relevant for applications that generate small packets, such as *voice-over-IP* (VoIP) packets, because the 40-byte added overhead may be even bigger than the actual VoIP payload.

The end of the bidirectional tunnel at the side of the mobile router needs to be updated each time the mobile network moves (and also periodically to refresh the binding at the home agent), to reflect the current location of the mobile router. This updating is achieved by the binding-update or binding-acknowledgement signaling exchange between the mobile router and the home agent. As stated previously, only one exchange (two packets, one in each direction) is required per movement, regardless of the number of MNNs that are attached to the mobile router—one of the main advantages of using the NEMO Basic Support Protocol on the mobile router instead of Mobile IPv6 on every node of the mobile network, because the signaling generated by a complete moving network (composed of numerous nodes) is the same as the one generated by a single moving node.

Conclusions

The NEMO Basic Support Protocol^[1] extends the functions of Mobile IPv6 to support the mobility of complete networks. The current specification supports basic mobility, and the IETF is currently working on new enhancements and extensions to provide route-optimization support, multihoming capabilities, and IPv4 support.

Some implementations of the NEMO Basic Support Protocol are already available. For example, the latest Cisco IOS® Software releases provide network mobility support. Open-source implementations also exist, such as the *NEMO Platform for Linux* (NEPL) (<http://www.mobile-ipv6.org/>) and SHISA (<http://www.mobileip.jp/>), for Linux and *Berkeley Software Distribution* (BSD) operating systems, respectively.

References

- [1] Vijay Devarapalli, Ryuji Wakikawa, Alexandru Petrescu, and Pascal Thubert, “Network Mobility (NEMO) Basic Support Protocol,” RFC 3963, January 2005.
- [2] Jon Postel, “Internet Protocol,” RFC 791, September 1981.
- [3] Iljitsch van Beijnum, “IPv6 Internals,” *The Internet Protocol Journal*, Volume 9, No. 3, September 2006.
- [4] Stephen E. Deering and Robert M. Hinden, “Internet Protocol, Version 6 (IPv6) Specification,” RFC 2460, December 1998.
- [5] Ralph Droms, “Dynamic Host Configuration Protocol,” RFC 2131, March 1997.
- [6] Ralph Droms, Jim Bound, Bernie Volz, Ted Lemon, Charles E. Perkins, and Mike Carney, “Dynamic Host Configuration Protocol for IPv6 (DHCPv6),” RFC 3315, July 2003.
- [7] William Stallings, “Mobile IP,” *The Internet Protocol Journal*, Volume 4, No. 2, June 2001.
- [8] Charles E. Perkins, “IP Mobility Support for IPv4,” RFC 3344, August 2002.
- [9] David B. Johnson, Charles E. Perkins, and Jari Arkko, “Mobility Support in IPv6,” RFC 3775, June 2004.
- [10] Thomas Narten, Erik Nordmark, and William A. Simpson, “Neighbor Discovery for IP Version 6 (IPv6),” RFC 2461, December 1998.
- [11] Pekka Nikander, Jari Arkko, Tuomas Aura, Gabriel Montenegro, and Erik Nordmark, “Mobile IP Version 6 Route Optimization Security Design Background,” RFC 4225, December 2005.
- [12] Chan-Wah Ng, Pascal Thubert, Masafumi Watari, and Fan Zhao, “Network Mobility Route Optimization Problem Statement,” Internet Draft, Work in Progress, September 2006.
draft-ietf-nemo-ro-problem-statement-03.txt

- [13] Chan-Wah Ng, Fan Zhao, Masafumi Watari, and Pascal Thubert, “Network Mobility Route Optimization Solution Space Analysis,” Internet Draft, Work in Progress, September 2006. **draft-ietf-nemo-ro-space-analysis-03.txt**
- [14] María Calderón, Carlos J. Bernardos, Marcelo Bagnulo, Ignacio Soto, and Antonio de la Oliva, “Design and Experimental Evaluation of a Route Optimisation Solution for NEMO,” *IEEE Journal on Selected Areas in Communications* (J-SAC), Issue on Mobile Routers and Network Mobility, Volume 24, Number 9, pages 1702–1716, September 2006.

CARLOS J. BERNARDOS received a telecommunication engineering degree in 2003, and a Ph.D. in telematics in 2006, both from University Carlos III of Madrid. His Ph.D. thesis focused on Route Optimisation for Mobile Networks in IPv6 Heterogeneous Environments. He has been working as a research and teaching assistant in Telematics Engineering since 2003. His current work focuses on IP-based mobile communication protocols. E-mail: **cjbc@it.uc3m.es**

IGNACIO SOTO received a telecommunication engineering degree in 1993, and a Ph.D. in telecommunications in 2000, both from the University of Vigo, Spain. He was a research and teaching assistant in telematics engineering at the University of Valladolid from 1993 to 1999. In 1999 he joined University Carlos III of Madrid, where he has been an associate professor since 2001. His research activities focus on mobility support in packet networks and heterogeneous wireless access networks. E-mail: **isoto@it.uc3m.es**

MARÍA CALDERÓN is an associate professor at the Telematics Engineering Department of University Carlos III of Madrid. She received a computer science engineering degree in 1991 and a Ph.D. degree in computer science in 1996, both from the Technical University of Madrid. She has published more than 20 papers in the fields of advanced communications, reliable multicast protocols, programmable networks, and IPv6 mobility. E-mail: **maria@it.uc3m.es**

More ROAP: Routing and Addressing at IETF68

by Geoff Huston, APNIC

Over the past year or so we have seen a heightened level of interest in Internet routing and addressing. Speculation regarding the future role of the Internet raises the possibility of the Internet supporting as many as hundreds of billions of chattering devices. What does such a future imply in terms of the core technologies of the Internet? Consideration of this topic has prompted a critical examination of the architecture of the Internet, including the scaling properties of routing systems, the forms of interdependence between addressing plans and routing, and the roles of addresses within the architecture.

The March 2007 meeting of the IETF, IETF68, saw some further steps in analysing these topics, and many sessions addressed aspects of routing and addressing. This article reports on these sessions, and includes some conjecture as to what lies ahead.

Plenary ROAP – The Plenary Session on Routing and Addressing

The plenary session presented an overview of the topic, looking at the previous initiatives in routing and addressing, as well as providing some perspectives on the current status of work in this area. There are concerns that the technology platform cannot scale by further orders of magnitude without some changes. Also of concern are the scalability of routing, the “transparency” of the network, renumbering questions, provider-based addressing, and service and traffic engineering and routing capabilities—and these concerns are potentially even more relevant and challenging for tomorrow’s Internet.

Our routing technology does not localize the external effects of local configuration choices. Far from being a protocol that damps instability, the *Border Gateway Protocol* (BGP) is a highly effective amplifier of noise components of routing events. So although it is a remarkably useful information-dissemination protocol, the properties of BGP in an ever-more connected world with ever-finer granularity of information raise some questions about its scaling properties. Will the imposed “noise” of the behaviour of the protocol completely swamp the underlying information content? Will we need to deploy disproportionately larger routers to support a larger network? The prospect here is that routing may become far less efficient because as we simultaneously increase the degree of interconnection and the information load, the inability to effectively localize information creates a far greater load on network routing.

In addition to these observations about routing, there is the continuing suspicion that the semantic load of addresses in the Internet architecture, where an address simultaneously conveys the concepts of “who,” “where,” and “how,” contributes to routing load.

To what extent the semantic intent of endpoint identity (or “id”) can be separated from the semantic intent of network location and forwarding lookup token (or “loc”) is a question of considerable interest. Although the current IP address semantics removes the need to support an explicit mapping operation between identity and location, the cost lies in the inability to support an address plan that is cleanly aligned to network topology, and the inability to cleanly support functions associated with device or network mobility. In the end it is the routing system that carries the consequent load. The questions in this area include an evaluation of the extent to which identity can be separated from location, and the effect of such a measure on the operation of applications. How much of today’s Internet architecture would be affected by such a change, and what would be the resultant benefits if this measure were deployed? Are we necessarily looking at a single model of such an id/loc split, or should we think about this scenario in a more general manner with numerous potential id/loc splits?

Obviously this study of routing and addressing, and the related aspects of name space attributes and mapping and binding properties, has a very broad scope. The larger question posed here is whether we can defer resolution of this problem to a comfortably distant future, or whether its effect on the present network is imminent. Are we accelerating toward some form of near-term technical limit that will cause a significant disruptive event within the deployed Internet, and will volume-based networks economics hold or will bigger networks start to experience disproportionate cost bloat—or worse? Is it time to be alarmed?

The unallocated IPv4 address pool will certainly be exhausted in the coming years, but this sense of alarm over routing and addressing is more about whether there are real limits in the near future in the capability to continue to route the Internet within the deployed platform, using the current technologies, and working within current cost-performance relationships irrespective of whether the addresses in the packet headers are 32 or 128 bits in size. There was a strong sense of “Don’t panic!” in the plenary presentation, with the relatively confident expectation that BGP will be able to carry the routing load of the Internet over the next 3 to 5 years without the need for major protocol “surgery,” and that Moore’s Law will continue to ensure that the capacity and speed of hardware will track the anticipated growth rates. Expectations are that the current technologies and cost-performance parameters will continue to prevail in this time frame.

The *Internet Engineering Steering Group* (IESG) has followed the *Internet Architecture Board’s* (IAB’s) initiative and has begun working with a focus group, the *Routing and Addressing Problem Directorate* (ROAP), to refine the broad space into many more specific work areas, and has assumed a role of coordination and communication across the related IETF activities.

In addition, because a relatively significant research agenda is posed by such long-term questions, the *Routing Research Group* of the *Internet Research Task Force* (IRTF) has been rechartered and, judging by the participation at its most recent meeting, effectively reinvigorated to investigate various approaches to routing that take us well beyond tweaking the existing routing toolset.

Internet ROAP – The Internet Area Meeting

The Internet Area meeting concentrated on aspects of this approach of supporting an identifier/locator split within the architecture of the Internet, and gathering some understanding as to whether this approach would assist with routing scaling. One of the important considerations in this area is working through what could be called boundary conditions of the study. For example, is this matter purely one for protocol stacks within an endpoint, or should distributed approaches that have active elements within the network also be considered? To what extent should a study consider mobility, traffic engineering, *Network Address Translation* (NAT), and *Maximum Transmission Unit* (MTU) behaviour? What appears to be clear at the outset is that this network is not a “clean-slate” network, and any approach should be deployable on the existing infrastructure, should use capability negotiation to trigger behaviours so that deployment can be incremental and piecemeal, should allow existing applications and their identity referential models to operate with no changes, and, hopefully, should have a direct benefit to those parties who decide to deploy the technology.

From the routing perspective, the overall desire is to reduce the growth rates of the interdomain routing space. The desired intent is to reduce the amount of information associated with locators so that locators reflect primarily network topology in such a way that the locators can be efficiently aggregated within the routing system that attempts to maintain a highly stable view of the network topology.

More detailed consideration of the implications of disambiguating aspects of identity from those of network location involves many dimensions—including the structure of the spaces—the mapping functions, and the practicalities of any form of deployment of such a technology.

A critical topic appears to be how an identity-mapping function relates to the forwarding-mapping function. Assuming that the existing name spaces remain unaltered, then the resultant framework appears to require distinct “name-to-identifier” and “identifier-to-locator” mappings and a “locator-to-forwarding” mapping. Where these mapping functions should be performed, who should perform them, when they should be performed, the duration of the validity of the outcomes, whether the mapping function outcomes are relative or universal, the scope and level of granularity in time and space of the map elements, the security of these mapping functions, and whether there is a simple operation in each mapping function or multiple operations all remain undefined at this point.

Other questions include whether the mapping is explicit or implicit, what evidence of a previous mapping operation is held in a packet in a visible manner, and what is occluded from further inspection after the mapping operation has been performed. In addition, what level of state is required in each host, and is there true end-to-end transparency—at what level?

It is likely, at least at this stage of the study, that such a split can have a variety of approaches, both in the intended roles of identifier and location tokens and in their binding. The expectation at this stage of the study is that further ideas will surface, and such ideas will be helpful rather than distracting. It is unclear if a single solution can emerge from this activity, or whether different actors have a sufficiently different set of relative priorities that multiple approaches—each of which expresses different prioritization of functions—are viable longer-term outcomes.

The critical consideration here is that it is unlikely that scaling routing over the longer term to a much larger network is simply a matter of just changing the operation of the routing system itself. Real improvement in this area appears to also require an understanding of the meaning of the objects, or “addresses,” that are being passed within the routing system. The motivation for opening up the identifier or locator space within the Internet area appears to be strongly tied to the notion that if you can unburden some of the roles of the addresses used in routing, and treat these routed tokens as unadorned network locality tokens, then you can gain some additional capability in routing.

Routing ROAP – The Routing Area Meeting

The first part of the Routing ROAP session looked at the trends in the routing system over 2005 and 2006. The overall trend appears to be a system that is increasingly densely interconnected, carrying more information elements, each of which expresses finer levels of granularity in reachability. There appears to be two forms of dynamic BGP load: the BGP “supernova” that burst with an intense BGP update load over some weeks and then disappear, and “background radiation” generators that appear to be unstable at a steady update rate for months or even the entire year.

In looking at scaling the BGP routing environment, one response is that of behavioural changes in local instances of BGP that reduce the potential for unnecessary updates to be propagated beyond a “need-to-know-now” radius. Another response is to consider changes to BGP in terms of additional attributes to BGP updates—such as a “withdrawal-at-origin” flag, or selective advertisement of “next best path”—both of which are intended to limit the span of advertised intermediate transitions while the BGP distance vector algorithm converges to a stable state.

It appears that we could improve our understanding of the operational profile of the routing space, looking particularly at the various forms of pathological routing behaviours and comparing these behaviours against the observations of known control points. Such a study may also lead to some more effective models of projections of the size of the routing space in the near- and medium-term future, and allow some level of quantification as to what “scaling of the routing space” actually implies.

The second part of the Routing ROAP session considered the current status of the routing world, updating some of the observations made at the IAB Routing Workshop and outlining some further perspectives on this space. One critical perspective on BGP is the behaviour of BGP under load. It was noted that most BGP implementations use adaptive responses to peer load, so that BGP attempts to ensure that its peer receives only the most current state information when the peer signals that it is not keeping pace with the update rate.

Another critical factor is the nature of “convergence” in BGP. The claim was made that this problem was the biggest, yet least important, problem with BGP. Convergence delays can be mitigated by *Graceful Restart*, *Nonstop Routing*, and *Fast Reroute*. One of the measures that exacerbates convergence is the use of *Route Reflectors*. The model of information hiding or Route Reflectors is intended to reduce the number of BGP peer sessions and the update load, but the benefits they do achieve are at the cost of slower convergence with a higher message rate during the intermediate-state transitions. Perhaps it is appropriate to consider small-scale changes to BGP behaviour to mitigate the transient BGP update bursts caused by path hunting, including those already mentioned of “withdrawal-at-origin” notification and propagation of backup paths.

The approach advocated here is based on the perspective that BGP is not in danger of imminent collapse, and there is still considerable “headroom” for BGP operation in today’s Internet.

More ROAP?

The routing space is a classic example of the commons, where each party can use routing to solve a multitude of business problems. This includes, for example, using routing to perform load balancing of traffic over a set of transit providers, using a “spot market” in Internet transit services, creating differentiated transit offerings using more specific routes and selective advertisements. The ultimate cost of these local efforts in optimising local business outcomes lies in the increasing bloat in the routing system and the consequent escalation in costs across the entire network in supporting the routing system. There is no way to impose administrative controls on the global routing system, nor have we been able to devise an economic model of routing where the incremental costs of local routing decisions are visible to the originator as true economic costs for the business, and the benefit of a conservative and prudent use of the routing system reaps economic dividends in terms of relatively lower costs for the business.

Like the commons, there are no effective feedback mechanisms to impose constraint on actors in the routing space. Also, like the commons, there is the distinct risk that the cumulative effect of local actions in routing creates a situation that pushes the routing system, either as a whole or in various locales, into a nonfunctioning state.

Whether it needs a sense of urgency to motivate the work, or a sense that there can and should be a better way to plan a future than crude crisis management, the underlying observation is that the routing and address world is fundamental to tomorrow's Internet. Unless we make a concerted effort to understand the various interdependencies and feedback systems that exist in the current environment, and understand the interdependences that exist between network behaviours and routing and addressing models, then I'm afraid that the true potential of the Internet will always lie within our vision—but frustratingly just beyond our grasp.

Further Reading

Following are references to further material on this topic, as presented at IETF68:

- <http://tools.ietf.org/html/draft-iab-raws-report-01>
- http://submission.apricot.net/chatter07/slides/future_of_routing/apia-future-routing-john-scudder.pdf
- http://submission.apricot.net/chatter07/slides/future_of_routing/apia-future-routing-jari-arkko.pdf
- <http://www3.ietf.org/proceedings/07mar/slides/plenaryw-3.pdf>
- <http://www3.ietf.org/proceedings/07mar/agenda/intarea.txt>
- <http://www3.ietf.org/proceedings/07mar/agenda/rtgarea.txt>
- <http://www1.tools.ietf.org/group/irtf/trac/wiki/RRG>
- <http://www.ietf.org/IESG/content/radir.html>

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. The author of numerous Internet-related books, he is currently the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He was a member of the Internet Architecture Board from 1999 until 2005, and served on the Board of the Internet Society from 1992 until 2001. E-mail: gih@apnic.net

Opinion: Is It Time to Replace SMTP?

by Dave Crocker, Brandenburg InternetWorking

The first Internet (ARPANET) e-mail, sent 35 years ago, was remarkably similar to a basic text e-mail of today: From, To, CC, Subject, Date, followed by lines of text, and the familiar @-sign in addresses. The right side of the address changed from a simple string into the multilevel domain name that we now use. The body can now be a set of multimedia attachments rather than just lines of text, but it can still be in its original, simpler form. The means of moving mail was the *File Transfer Protocol* (FTP) in the early 1970s. The current mechanism, the *Simple Mail Transfer Protocol* (SMTP)^[1a, 1b], was not created until 10 years later, but a mere 25 years of use is not bad, either.

All of the technical specifications for e-mail have undergone many changes over the years, but a core requirement has been to protect the installed base of users and operators by incrementally adding features as options, rather than by performing wholesale replacement of any infrastructure service component. E-mail has changed the way we communicate, yet it is also now viewed as having a serious problem: As the Internet grew, it acquired the full mixture of participants, some of whom do not make nice neighbors.

Frustration with the effect of abusive users is often expressed as a belief that the solution lies in replacing some or all of the core technology of the e-mail service, or even by moving to an entirely different paradigm, such as querying Webpages using *Really Simple Syndication* (RSS)^[2]. Although different paradigms make sense for some forms of human communications, what is forgotten in these pleas for massive change is the power of the classic mail model, whether by paper or by electrons: Spontaneous or occasional communication requires the ability to “push” the message to the recipient, without prior arrangement. This ability is, of course, also what leaves the door open for abuse—anyone may walk in, uninvited and unwanted.

The alternative proposals might work well enough for ongoing, regular communication among people who already know each other. And for most of us, that is probably 80 percent of our exchanges, or more. Unfortunately, as soon as anyone starts worrying about the remaining 20 percent, these alternative approaches require cascading hacks, producing a design that looks no better than what we have today, except that it is based strictly on theory rather than decades of practice. It is easy for a paper proposal to beat a deployed system; making it work as promised is, of course, more difficult.

Mantra

I have developed a simple mantra, in response to calls for replacing today's Internet mail:

0. The basic problems we are experiencing with e-mail are really based on undesirable social behaviors, long popular outside the Internet. The Internet enables broader reach, to more victims, and in much shorter time spans, but the core misbehaviors have existed for all of recorded human history. We should not assume that there are technical solutions to social problems.
1. The beginning of changing a human service is to gain community consensus about the change that is needed, because a mechanism will not be successful unless it is perceived as needed. Only then can the engineers work on designing the change.
2. When there is community consensus about the way that e-mail needs to be changed, the folks who are currently contributing to its 35-year evolution need to try to find a way to add the desired features to the existing service. Given the record of accomplishment of e-mail, the odds seem favorable that any new requirement can also be satisfied without disrupting the installed base.
3. When that effort fails, it will be time to create a replacement infrastructure.

Alas, as those who track e-mail abuse technical discussions are well aware, we have not completed Step 1. As soon as we try to formulate community consensus about basic messaging communication policies, discussion devolves into cacophony or marginalized community fragments. It is certain that there will eventually be a change required for e-mail, which we cannot fit into the current service, but we do not yet have any evidence that e-mail abuse is going to produce that requirement.

Trust Models

One hopeful sign is that we do have a solid set of efforts to evolve e-mail to support mechanisms that are based on trust. This evolution begins with the ability to associate a validated identity to a message and then requires assessing the "safety" of that identity's owner. Until recently, only the IP address of the last-hop sending SMTP server could be used as an identifier. Using addresses as identifiers sounds reasonable at first glance, but turns out to have long-term scaling and administrative problems. As a result, there has been a broad effort to find ways to use domain names, which are more stable, and they align better with organizational boundaries. This process is well under way, with the recent IETF standardization of the *Domain Keys Identified Mail* (DKIM)^[3] message-signing specification, as well as path-based registration schemes, such as Sender-ID^[4] and SPF^[5].

That took about 5 years. And now comes the hard part: developing a range of *assessment mechanisms*—sometimes generically called *reputation services*—that satisfy requirements for quality, strength, convenience, and stability. Assessment services tell recipients whether the author of the message, or the service that sent it, can be trusted. Some mechanisms need to work for small groups, others need to work for mass-market business-to-consumer mailings, and others need to work among business partners. A few startup companies have recently joined the few, surviving volunteer services, to satisfy this need. It is too early to tell whether they will suffice, or whether additional services will be needed. What is important is that these services are generally regarded as producing good results.

For the long term it seems likely that this capability will result in an Internet mail service that is logically split into two types of traffic. One has substantial trust associated with its messages, so that they can be delivered with a reasonable degree of comfort. The other is the current, open-to-all service that requires heavy filtering and the use of various heuristics, to reduce the effect of abuse mail. If the first traffic flow is sufficiently successful, filters for the second can become much more stringent. The aggregate effects of these changes will be that wanted mail is likely to be received and identified much more reliably, and unwanted mail is more likely to be rejected.^[8, 9]

So the current Internet mail technical infrastructure is safe, right? Well, maybe.

Enhancements?

What gets less attention, but perhaps should worry us more, is the general lack of user-level functional enhancement for e-mail. What users can do with e-mail, today, is pretty much the same as they could do 25 years ago. The evolution of Internet mail has been primarily in support of performance, reliability, and scaling. Although important, they have not produced functional changes that are apparent to end users. Human communication is a very rich space, yet most e-mail is limited to a narrow range of styles: person-to-person informal communications, and informal, unstructured group communications. Toss in some very basic, one-way “transactional” mail, such as order confirmations from businesses to their customers, and that about covers it.

Instead, new functions for human collaboration have tended to appear in new services. *Instant Messaging* (IM), blogging, and wikis are the most popular examples. In each case, they rely on a centralized service, rather than the highly distributed model that e-mail uses. Users must all go to a single, centralized address to obtain a given service. Most of the IM world does not even know that there are two (!) Internet standards for distributed IM—*Extensible Messaging and Presence Protocol* (XMPP)^[6] and SIMPLE^[7]. Even for these standards, most of their production use tends to be within noninteroperable, centralized services.

Is there something about e-mail that is a barrier to functional enhancements for end users?

For these new services, the interservice relaying that is at the core of e-mail is absent. Indeed, centralized services are easier to create and operate than are distributed services, but they also carry scaling, administration, and control challenges. So the issue is not so much what is easier, but who will do the work—and when? With a centralized service, all the interesting work is done by the single provider. For a distributed model, like e-mail, the work is shared across participating organizations. The Internet was designed to avoid single points of failure (and failure), so it is ironic that these new services risk exactly these problems.

For a distributed model, like e-mail, to add end-user functions, useful adoption is required by all user software that participates, and possibly by all the intermediate, relaying services. The adoption is in three parts: agreeing on the enhancement, modifying existing software, and making it available to users. These are daunting barriers, so the appeal of centralized services is clear: a single organization decides what to change, changes it, and makes it available to end users with, at most, a natural software upgrade.

Interorganization partnerships provide the best argument for adoption of distributed services, because they do not naturally permit agreement on a central point of control. The counterforce is, again, the simplification (for the partners) that comes from agreeing to use independent third-party services. The scaling problem here is with end users having to juggle a large number of independent services. Note the emergence of IM clients that support a variety of independent IM services.

Perhaps the real danger to e-mail is not its wholesale and traumatic replacement, stemming from frustration about abuses, but a gradual attrition, as portions of its traffic move to services that evolve more quickly, but leave end users with a complicated array of narrow, specialized, and noninteroperable venues.

References

- [1a] Postel, J. B., “Simple Mail Transfer Protocol,” RFC 821, August 1982.
- [1b] Klensin, J., “Simple Mail Transfer Protocol,” RFC 2821, April 2001.
- [2] Really Simple Syndication Specifications,
<http://www.rss-specifications.com/rss-specifications.htm>

- [3] Allman, E., et al., “DomainKeys Identified Mail (DKIM) Signatures,” February 2007. (*RFC publication pending.*)
<http://dkim.org/specs/draft-ietf-dkim-base-10.html>
- [4] Lyon, J. and Wong, M., “Sender ID: Authenticating E-Mail,” RFC 4406, April 2006.
- [5] Wong, M. and Schlitt, W., “Sender Policy Framework (SPF) for Authorizing Use of Domains in E-Mail, Version 1,” RFC 4408, April 2006.
- [6] Saint-Andre, P. (ed.), “Extensible Messaging and Presence Protocol (XMPP): Instant Messaging and Presence,” RFC 3921, October 2004.
- [7] Campbell, B. (ed.), Rosenberg, J., Schulzrinne, H., Huitema, C., and Gurle, D., “Session Initiation Protocol (SIP) Extension for Instant Messaging,” RFC 3428, December 2002.
- [8] Crocker, D., “Challenges in Anti-Spam Efforts,” *The Internet Protocol Journal*, Volume 8, No. 4, December 2005.
- [9] Klensin, J., “Taking Another Look at the Spam Problem,” *The Internet Protocol Journal*, Volume 8, No. 4, December 2005.

DAVE CROCKER is a principal with Brandenburg InternetWorking. He has authored or contributed to most Internet mail standards, and an assortment of e-mail products and businesses, as well as working on facsimile, security, e-commerce, and EDI. He received the 2004 *IEEE Internet Award* for his work on e-mail. Dave is a contributor to the development efforts for DKIM, CSV, and BATV, motivated by a strong desire to protect more than 30 years of professional investment that is being threatened by spamming. E-mail: [**dcrocker@bbiw.net**](mailto:dcrocker@bbiw.net)

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

ARIN Board Advises Internet Community on Migration to IPv6

The *American Registry for Internet Numbers* (ARIN) and the other *Regional Internet Registries* (RIRs) have distributed Internet Protocol version 6, IPv6, alongside IPv4 since 1999. To date, ARIN has issued both protocol versions in tandem and has not advocated one over the other. ARIN has closely monitored trends in demand and distribution for both protocol versions with the understanding that the IPv4 available resource pool would continue to diminish.

The available IPv4 resource pool has now been reduced to the point that ARIN is compelled to advise the Internet community that migration to IPv6 is necessary for any applications that require ongoing availability from ARIN of contiguous IP number resources. On 7 May 2007, the ARIN Board of Trustees passed the following resolution:

“Whereas, community access to *Internet Protocol* (IP) numbering resources has proved essential to the successful growth of the Internet; and,

Whereas, ongoing community access to *Internet Protocol version 4* (IPv4) numbering resources can not be assured indefinitely; and,

Whereas, *Internet Protocol version 6* (IPv6) numbering resources are available and suitable for many Internet applications,

Be it Resolved, that this Board of Trustees hereby advises the Internet community that migration to IPv6 numbering resources is necessary for any applications which require ongoing availability from ARIN of contiguous IP numbering resources; and,

Be it Ordered, that this Board of Trustees hereby directs ARIN staff to take any and all measures necessary to assure veracity of applications to ARIN for IPv4 numbering resources; and,

Be it Resolved, that this Board of Trustees hereby requests the ARIN Advisory Council to consider Internet Numbering Resource Policy changes advisable to encourage migration to IPv6 numbering resources where possible.”

Implementation of this resolution will include both internal and external components. Internally, ARIN will review its resource request procedures and continue to provide policy experience reports to the Advisory Council. Externally, ARIN will send progress announcements to the ARIN community as well as the wider technical audience, government agencies, and media outlets. ARIN will produce new documentation, from basic introductory fact sheets to FAQs on how this resolution will affect users in the region. ARIN will focus on IPv6 in many of its general outreach activities, such as speaking engagements, trade shows, and technical community meetings. For more information, visit ARIN’s IPv6 Information Center at:

<http://www.arin.net/v6/v6-info.html>

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2007 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA
--

The Internet Protocol Journal

September 2007

Volume 10, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

FROM THE EDITOR

In This Issue

From the Editor	1
Secure Multivendor Networks.....	2
IPv4 Address Depletion	18
IPv4 Address Consumption ..	22
Awkward /8 Assignments	29
Book Review	32
Call for Papers.....	35

For the last 10 or so years I have been involved with the organization of APRICOT, the *Asia Pacific Regional Internet Conference on Operational Technologies*. APRICOT has at its core a set of workshops featuring expert instructors with years of operational network experience. A recent addition to the APRICOT workshop program is a course focusing on Internet security in a multivendor environment. Our first article, written by Kunjal Trivedi from Cisco Systems, Inc., and Damien Holloway from Juniper Networks, is based on this workshop. It's not every day that you see an article co-authored by instructors from competing companies, but this is exactly the type of cooperation that is needed in order to deploy security in a multivendor network.

The rest of this issue is mostly devoted to IPv4 depletion and the transition to IPv6. The first article, by Geoff Huston, summarizes many of the concerns related to IPv4 depletion and IPv6 transition, and gives numerous pointers to further articles and documents of interest. Our second addressing-related article, by Iljitsch van Beijnum, looks more closely at the numbers relating to address allocation by the *Regional Internet Registries* (RIRs). The final article concerns some address blocks that are currently unassigned but actually in use. Leo Vegoda explains the potential problems that may arise when these blocks eventually become part of the RIR assignment pool.

We are pleased to announce a new online addition to this journal. *The Internet Protocol Forum* (IPF) available at www.ipjforum.org is designed to allow discussion of any article published in the printed edition of IPJ. In addition to article discussions, the forum will be used to provide updates and corrections, downloads, expanded versions of some articles, configuration and programming examples, and news and analysis that does not fall into our quarterly publication schedule. The IPF's editor and moderator is Geoff Huston, long-time contributor to this journal and chief scientist at APNIC. I am confident that IPF will become an important addition to IPJ, and I hope you will take the time to participate in the online discussions. Of course, you can always contact us at the usual e-mail address: ipj@cisco.com

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

A Standards-Based Approach for Offering a Managed Security Service in a Multivendor Network Environment

By Kunjal Trivedi, Cisco Systems and Damien Holloway, Juniper Networks

As transport becomes a commodity, service providers are seeking new revenue sources and new ways to differentiate themselves. Managed security services address a growing market because business customers are struggling to comply with regulatory requirements such as the *Payment Card Industry-Data Storage Standards* (PCI-DSS), the *Sarbanes-Oxley Act*, the *Gramm-Leach Bliley Act*, *Health Insurance Portability and Accountability Act* (HIPAA), *Directive 2002/58/EC*, and the *Asia-Pacific Economic Cooperation-Organization for Economic Cooperation and Development* (APEC-OECD) initiative on regulatory reform. Increasingly, business customers recognize that outsourcing network security is less costly than staffing with highly specialized security personnel who can provide 24-hour incident detection and response. Another incentive for outsourcing is to free existing IT resources to focus on the core business.

A standards-based approach helps service providers take best advantage of the managed security service opportunity because it increases the potential breadth and depth of the service offering. Multivendor solutions are becoming the norm when deploying services on an integrated backbone. Therefore, standards simplify deployment and management, helping control operational costs and accelerating time to market.

Service providers are experiencing a growing need for skilled engineers who understand multivendor environments—the motivation for conducting a multivendor security workshop at the 2006 *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT 2006)^[15], held in Perth, Australia, in February 2006. During the workshop (which was repeated again at APRICOT 2007 in Bali), participants successfully deployed and tested a multivendor service environment using *IP Security* (IPsec)-based Layer 3 *Virtual Private Networks* (VPNs)^[1, 2, 3] over a *Border Gateway Protocol/Multiprotocol Label Switching* (BGP/MPLS) core^[4].

Technical Challenges

To offer managed security services, service providers need the following:

- A secure network infrastructure, including tools and techniques for risk mitigation
- Technical solutions for the customer's business needs, such as VPNs based on BGP/MPLS, IPsec, or both

- Web-based reporting tools that business customers can use to monitor the security service in accordance with *Service-Level Agreements* (SLAs). Service providers can scale cost-effectively by offering customers a secure, Web-based portal that shows open trouble tickets, security incident-handling detail, SLAs, and access reports that customers need to comply with regulations.

An effective managed security service requires tools and techniques to address the following challenges:

- *More sophisticated threats, and less time between vulnerability and exploitation:* In addition to worms and viruses, the service provider needs to protect its own and its customers' networks against *Denial-of-Service* (DoS) attacks. Today's botnets can launch thousands or even a million bots that carry out outbound DoS attacks. New varieties of worms have side effects similar to those of DoS attacks. These threats can take down the service provider infrastructure, thereby violating SLAs and eroding revenue.
- *A need for proactive rather than reactive threat response:* Many service provider security groups are stuck in reactive mode. Every network device and security system produces voluminous event logs every day, and vendors use different formats. Therefore, identifying security incidents in order to react to them can take hours or days—or not happen at all. The connection between two separate events in different parts of the network can easily escape human detection, especially when the clues are buried among tens of thousands of harmless events that took place around the same time.
- *Multivendor networks:* Network security and reporting are easier to achieve in single-vendor networks. Realistically, however, many service providers and business customers have multivendor networks, sometimes because of mergers and acquisitions. Even if the service provider itself has a single-vendor network, some of its customers will use other vendors' equipment.
- *Slow progress toward adopting IP Next-Generation Networks (IP NGNs):* When service providers complete the migration to IP NGN, they will achieve greater control, visibility, and operational efficiency. Until then, service providers will incur higher costs and labor requirements for support and migration.
- *A need to comply with industry standards from IETF and ITU:* Standards facilitate security in multivendor networks. MPLS helps ensure infrastructure security, whereas IPsec provides secure connectivity among the customer's branches and remote offices. By using industry standards, the service provider can select best-of-class products based on performance, features, or cost.

- *Scalability challenges:* The security operations center for a managed services provider cannot cost-effectively scale to process several million events for each customer. However, it can scale to process a few security-incident trouble tickets. Scalability hinges on the ability to minimize false positives. Products such as Cisco *Security Monitoring, Analysis and Response System* (MARS), IBM Micromuse, and NetIQ provide analysis and correlation of events from multiple elements in the IT infrastructure. They process events using consolidation, filtering, normalization, enrichment, correlation, and analysis techniques, and also notify IT staff about critical events.

Infrastructure Security in Multivendor Environments

Securing the service provider infrastructure requires the following common best practices:

- Point protection
- Edge protection
- Remote-triggered black-hole protection
- Source-address validation on all customer traffic
- Control-plane protection
- Total visibility into network activity

Point Protection

Before offering a managed security service, providers need to protect the backbone; security operations center or network operations center; *Authentication, Authorization, and Accounting* (AAA)^[10, 11] server; and remote-access networks. Securing individual network devices requires enforcing AAA, controlling the type of packets destined to network devices, and performing regular configuration audits to ensure that no unauthorized changes have been made. Best common practices include:

- *Protect the backbone by locking down the vty and console ports:* This protection helps prevent unauthorized access to network devices.
- *Encrypt management commands that staff send to devices:* Use of the *Secure Shell* (SSH) protocol helps prevent hackers from obtaining passwords that they could later use to compromise the network. Service providers that use out-of-band management for device configuration should also encrypt this management traffic and restrict access to authorized personnel.
- *Deploy a AAA server:* Using a AAA server is preferable to relying on local authorization on the devices themselves because it enables centralized policy control. The AAA server controls a user's access to the device, or even the specific commands that the user is authorized to execute.

It is strongly recommended that service providers use TACACS+^[13] authentication rather than *Remote Authentication Dial-In User Service* (RADIUS)^[12] authentication. With RADIUS, traffic is sent in the clear between the AAA servers and network devices using the *User Datagram Protocol* (UDP), which defeats the use of SSH to encrypt logins and passwords. Open-source implementations of TACACS+ are available.

- *Use one-time passwords (OTPs)*: To distribute one-time passwords, service providers can provide authorized users with a token card, soft token, or soft key. One-time passwords ensure that the user was authorized at the time of login, and was not an attacker who used a packet-sniffer program to intercept a password.
- *Protect the AAA infrastructure from DoS attacks*: Some service providers set up local accounts on routers and switches so that staff can log in if the AAA infrastructure is down, creating vulnerability. If the service provider does not secure management-plane access to the device, hackers can use SSH or Telnet and attempt a brute-force attack to crack the local account. The local account is often not as secure as an OTP because it is changed only once every 30 days, providing a longer window of opportunity for hackers to gain device access. It is strongly advised to not use default or easy-to-guess passwords. To prevent attacks against the AAA infrastructure, service providers should harden the infrastructure and consider placing the server behind a firewall with stateful inspection. Use *Access Control Lists* (ACLs), which are packet filters, on the firewall to restrict traffic between the AAA server and network devices only. Also be sure to distribute the AAA servers so that they do not create a single point of failure.
- *Regularly audit device configurations*: Frequently, the first indications of an attack, often unnoticed, are unauthorized commands executed on routers that change the configuration. An easy way to monitor configurations is using RANCID (*Really Awesome New Cisco Config Differ*)^[14], a UNIX or Linux freeware tool that logs into each of the devices in the device table file, runs various show commands, processes the output, and sends e-mail messages reporting any differences from the previous collection to staff. RANCID works with routers from Cisco and other vendors. Another tool for auditing device configurations, the *Router Audit Toolkit* (RAT) assigns security scores to ACLs and other security best practices to show the relative security of routers.

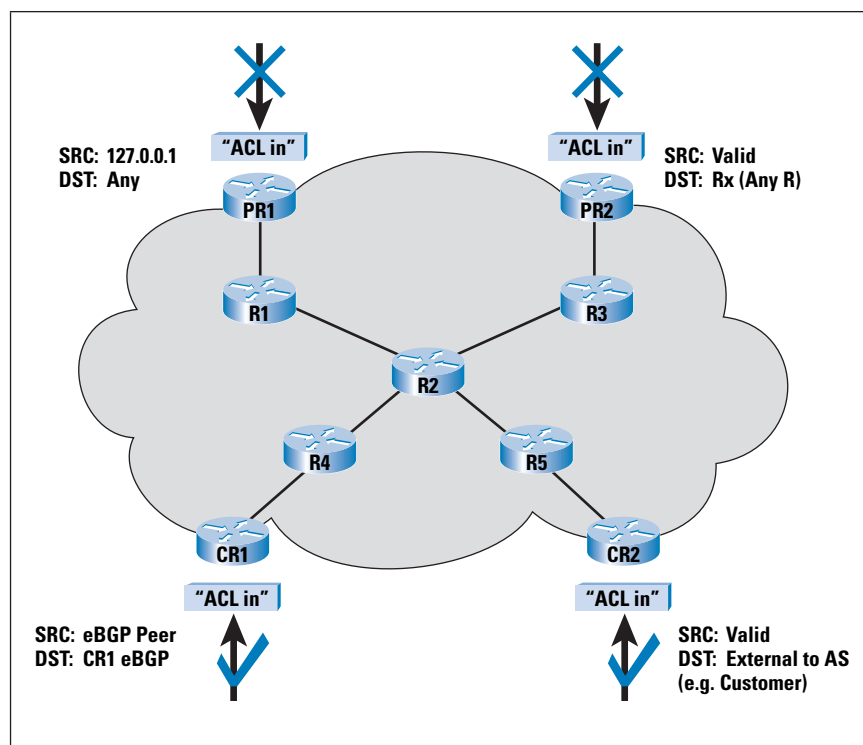
Traditionally, service providers enforced policy at the process level, using vty ACLs, *Simple Network Management Protocol* (SNMP) ACLs, and others. Some service providers used ingress ACLs when possible. Today, it is far preferable to stop DoS traffic at ingress points: the peer edge, downstream and upstream routers, colocated network devices, and the customer access edge, enabling central policy enforcement and more granular protection schemes.

In addition, many network devices at the network edge have hardware acceleration, which provides far more robust resistance to attack than the process level.

Edge Protection

In many service provider networks, each core router is individually secured but still accessible to outsiders using SNMP or Telnet. Now service providers can supplement individual router protection with infrastructure protection that prevents undesired traffic from ever touching the infrastructure.

Figure 1: Protecting the Network Edge



The following steps help protect the network edge (Figure 1):

1. Classify the required protocols that are sourced from outside the *Autonomous System* (AS) access core routers, such as *external BGP* (eBGP) peering, *Generic Routing Encapsulation* (GRE)^[5], and IPsec. (Examples of nonrequired protocols are SNMP and Telnet.) Classification can be performed using a classification packet filter or Cisco *NetFlow* telemetry. The classification packet filter comprises a series of permit statements that provide insight into required protocols. Gradually narrow down the list, keeping in mind that very few protocols need access to infrastructure equipment, and even fewer are sourced from outside the autonomous system. Summarize the IP address space as much as possible, for simpler and shorter ACLs. Be cautious: just because certain types of traffic appear in a classification packet filter or NetFlow telemetry data does not mean they should be permitted to pass through to the routers.

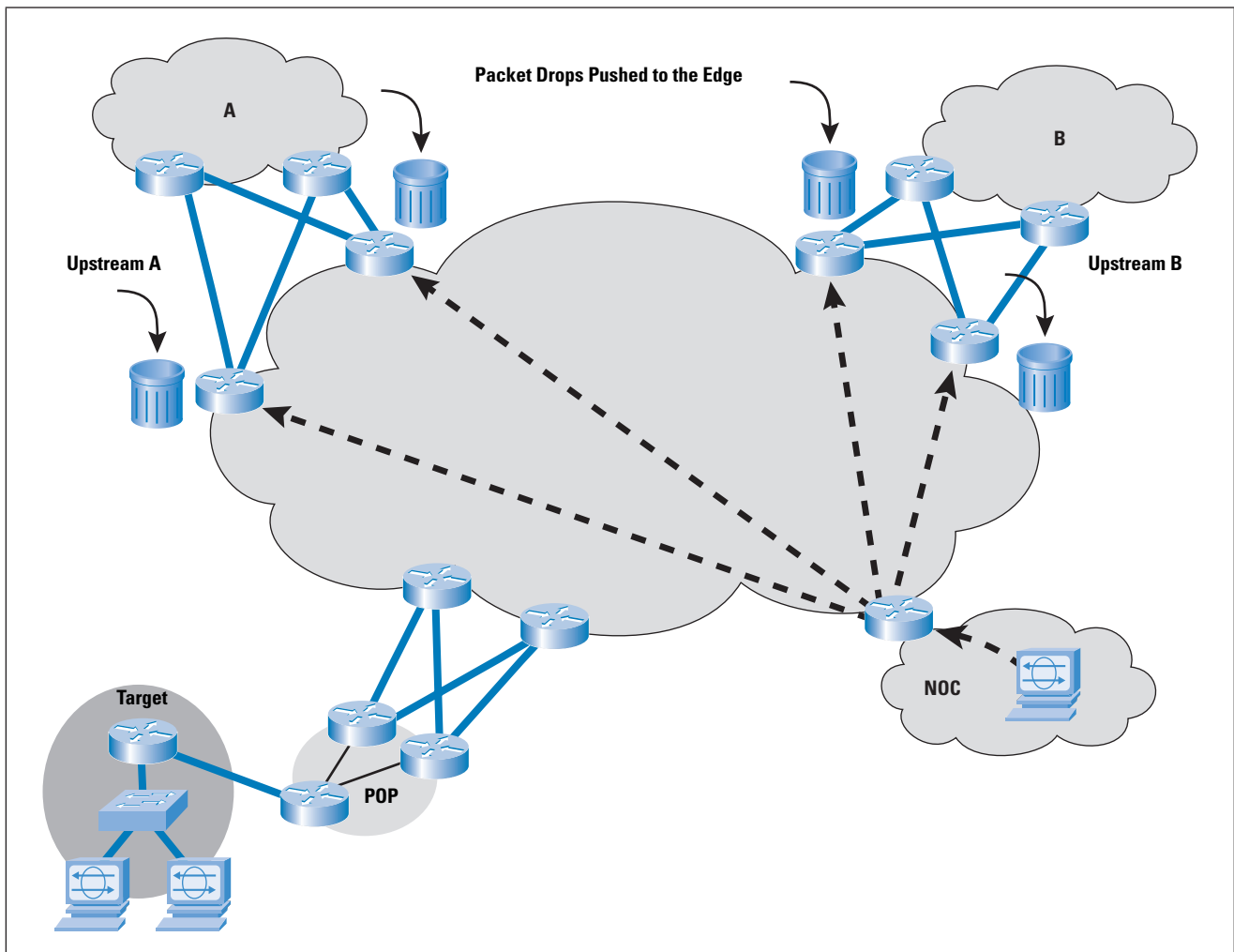
2. Begin filtering. Use an infrastructure packet filter to permit only the required protocols to access infrastructure-only address blocks, denying all other protocols. It is important to monitor the packet filter entry counters, because a high volume of hits, whether or not a protocol has been identified as required, might signal an attack. To permit transit traffic, use the following as the final line of the *Infrastructure ACL* (iACL): **permit ip any any**, protecting the core network with a basic iACL that admits only the required protocols. Note that iACLs also provide antispoof filtering by denying access to the space from external sources, denying the RFC 1918 space^[6], and denying multicast source addresses. RFC 3330^[7] defines special-use IPv4 addressing.
3. Further protect the core by identifying legitimate source addresses for the required protocols, such as external BGP peers and tunnel endpoints.
4. Deploy destination filters when possible.

Infrastructure packet filters at the edge of the network protecting the infrastructure are an effective first layer of defense. Service providers need additional forms of infrastructure protection for their older routers that do not support infrastructure packet filters and for packets that cannot be filtered with infrastructure packet filters.

Remote Triggered Black Hole Filtering

Remote Triggered Black Hole Filtering (RTBH) is among the most effective reaction and mitigation tools for DoS, *Distributed DoS* (DDoS), and backscatter tracebacks. It enables service providers to quickly drop DoS traffic at the network edge (Figure 2). Rather than sending commands to every router to drop DoS or other problem traffic, the service provider can deploy a trigger router that uses BGP to signal all other routers—just as fast as iBGP can update the network. In destination-based RTBH, all traffic headed to the destination under attack is dropped—the good traffic as well as the bad. In source-based RTBH, traffic from all or certain sources are blocked. The advantage of sourced-based RTBH is that service providers can whitelist certain addresses, such as the *Network Operations Center* (NOC) or route-name servers, so that they can continue providing services.

Figure 2: DoS Packets Dropped at the Network Edge



Source Address Validation on all Customer Traffic

Source address validation, defined in *Best Current Practices* (BCP) 38^[8], prevents service provider customers from spoofing traffic—that is, sending IP packets out to the Internet with a source address other than the address allocated to them by the service provider. Best practices from BCP 38 are to filter as close to the edge as possible, filter precisely, and filter both the source and destination address when possible.

Every access technology has antispoofing mechanisms derived from BCP 38:

- Packet filters
- Dynamic packet filters that are provisioned to be AAA profiles; when a customer signs in with RADIUS, a packet filter is set up for the customer
- *Unicast Reverse Path Forwarding* (URPF)
- Cable-Source Verify and packet cable multimedia (cable)
- IP Source Verify and DHCP Snooping (Metro Ethernet)

To gain operational confidence in BCP 38, service providers can take a phased approach—for example, implementing it first on one port, then on a line card, then on an entire router, and then on multiple routers.

Control-Plane Protection

Protecting the infrastructure control plane helps prevent an attacker from taking down a BGP session and thereby causing denial of service. The exploits a service provider needs to prevent include saturating the receive-path queues so that BGP times out, saturating the link so that the link protocols time out, dropping the *Transmission Control Protocol* (TCP) session, and dropping the *Interior Gateway Protocol* (IGP), which causes a recursive loop-up failure.

Following are techniques for control-plane protection.

- *Generalized Time-to-Live (TTL) Security Mechanism (GTSM)*: This technique protects BGP peers from multihop attacks. Routers are configured to transmit their packets with a TTL of 255, and to reject all packets with a TTL lower than 254 or 253. Therefore, a device that is not connected between the routers cannot generate packets that either router will accept.
- *Configuring routing authentication*: The *Message Digest Algorithm 5* (MD5) peer authentication feature instructs the router to certify the authenticity of its neighbors and the integrity of route updates. MD5 peer authentication can also prevent malformed packets from tearing down a peering session, and unauthorized devices from transmitting routing information. Be aware that MD5 peer authentication does not protect the router if an attacker compromises the router and begins generating bogus routing updates. Although it is not a panacea, MD5 peer authentication does raise the level of protection.
- *Customer ingress prefix filtering*: Prefix hijacking is an exploit in which a service provider customer announces an address space that belongs to another customer. The remedy is customer ingress prefix filtering, which enables service providers to accept only those customer prefixes that have been assigned or allocated to their downstream customers. For example, if a downstream customer has a **220.50.0.0/20** block, customers can announce this block only to their peers, and upstream peers accept this prefix only. Service providers can apply ingress prefix filtering to and from customers, peers, and upstream routers.

Visibility into Network Activity

To gain visibility into the network for early detection of security incidents, service providers can use open-source tools to analyze flow-based telemetry data, which is retrieved from routers and switches. Open-source tools for visibility into security incidents include RRDTool, FlowScan, Stager, and NTOP *Remote Monitoring* (RMON).

These tools provide information such as packets per second, bits per second, and traffic types. For example, RRDTool shows the number of *Domain Name System* (DNS) queries per second, according to record type. A spike in *Mail Exchange* (MX) Record queries might indicate that a customer's router has been compromised and is being used as a spam proxy. Similarly, a sharp increase in round-trip-time latency might indicate a DoS attack.

MPLS Security in a Multivendor Environment

In addition to securing the infrastructure, managed security service providers need to secure packets as they travel from one customer-edge router to another—regardless of the equipment the customer uses at the edge. Layer 3 VPNs meet this need. RFC 4364, which replaced RFC 2547bis, defines a BGP/MPLS IP VPN that creates multiple virtual routers on a single physical router: one virtual router for each customer.

In BGP/MPLS VPNs, *Customer Edge* (CE) routers send their routes to the *Service Provider Edge* (PE) routers. Customer edge routers at different sites do not peer with each other, and the customer's routing algorithms are not aware of the overlay. Data packets are tunneled through the backbone so that the core routers do not need to know the VPN routes. BGP/MPLS IP VPNs support either full mesh or partial mesh, although full mesh is more cost-effective.

A unique advantage of BGP/MPLS VPNs is that two service provider customers with overlapping IP addresses can connect across the service provider backbone. The router distinguishes between traffic from different companies by examining the label at the beginning of the packet, and then instantly forwards the traffic based on the *Label Switching Path* (LSP) that has been established for each customer's VPN. Eliminating the need to look at the packet in depth enables faster forwarding. That is, the service provider core does not impose any latency as packets pass between the provider edge routers.

IPsec Security in a Multivendor Environment

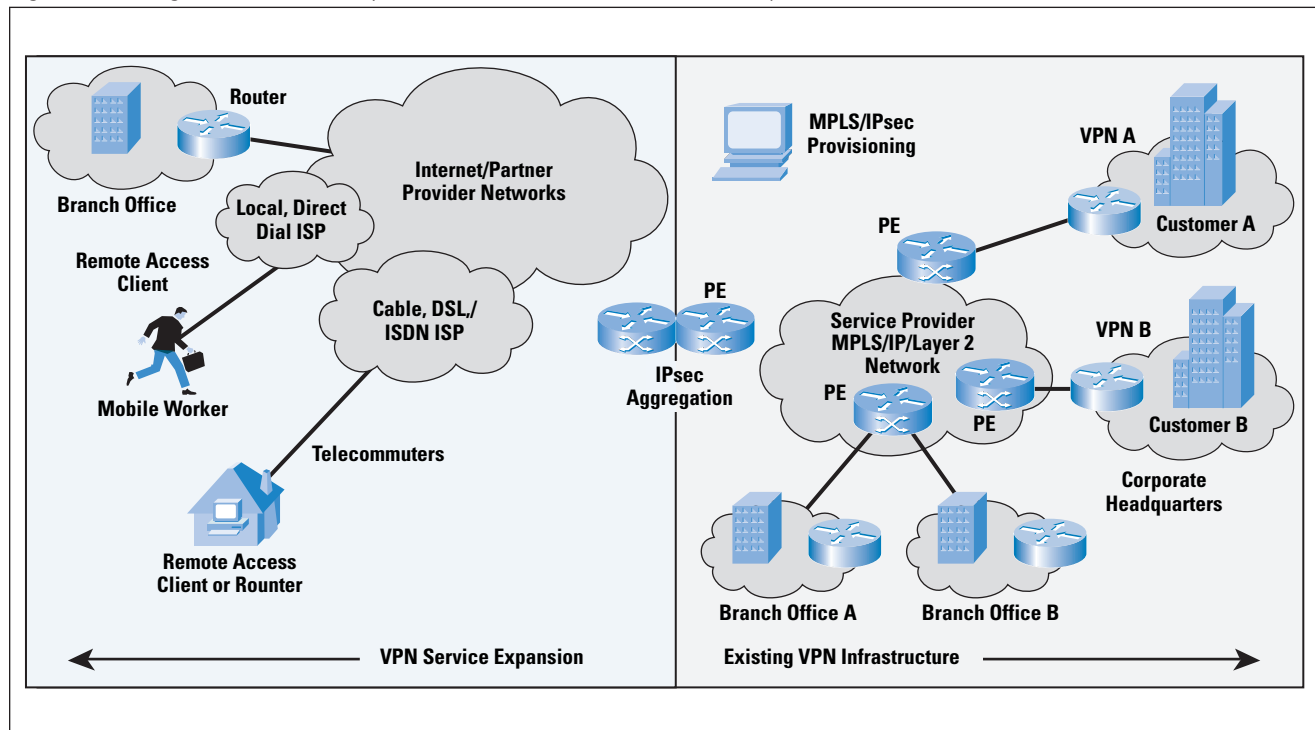
In addition to or instead of deploying a BGP/MPLS IP VPN, the service provider can extend its service to other partner provider networks using IPsec. The options are to use MPLS alone, IPsec alone, or a combination (Figure 3 on page 12). A retail customer that needs to comply with PCI-DSS, for example, needs IPsec or *Secure Sockets Layer* (SSL) encryption for payment card transaction data as part of its managed security service.

Table 1 summarizes the process based on the option the service provider selects. In the table, VPNA refers to one customer's VPN on a router that hosts VPNs for multiple customers.

Table 1: Comparing Packet Flow in IPsec VPNs, BGP/MPLS VPNs, and Combination VPNs

IPsec	BGP/MPLS VPN	BGP/MPLS VPN and IPsec
<ol style="list-style-type: none"> Host A in site 1 of VPNA sends packets to host B in site 2 of VPNA. Routers A and B negotiate an Internet Key Exchange (IKE) [9] phase-one session in aggressive or main mode to establish a secure and authenticated channel between peers. Routers A and B negotiate an IKE phase-two session to establish security associations on behalf of IPsec services. Information is exchanged securely through an IPsec tunnel. The tunnel is terminated. 	<ol style="list-style-type: none"> Host A in site 1 of VPNA sends packets to host B in site 2 of VPNA. Packet arrives on a VPN Route-Forwarding (VRF) VPNA interface on the PE1 router. The PE1 router performs an IP lookup, determines the label stack and the outgoing core-facing interface, and forwards the packet to the MPLS core. The packet is label-switched at each hop in the core until it reaches the penultimate hop router. At this point, the top label is popped before the packet is forwarded to the egress provider edge router. The egress PE2 router performs a MPLS lookup and determines that it should remove the label before forwarding the packet to host B in site 2. Router B in site 2 receives a regular IP packet and forwards it to host B. 	<ol style="list-style-type: none"> Router A in site 1 and the associated PE1 router negotiate an IKE phase-one session in aggressive or main mode to negotiate a secure and authenticated channel between peers. Router A and the PE1 router negotiate an IKE phase-two session to establish security associations on behalf of IPsec services so that information is exchanged securely through an IPsec tunnel. Host A in site 1 of VPNA sends packets to host B in site 2 of VPNA. The PE1 router, which is enabled with VRF-aware IPsec, creates a direct association through the IPsec tunnel that connects site 1 and the corresponding VRF ID (VPNA) on the provider edge router over the Internet. Encrypted traffic arrives on an Internet-facing interface on the provider edge router A, which terminates the IPsec tunnel, decrypts the incoming packet, and forwards the plaintext packet to the VRF VPNA for further processing. The PE1 router performs an IP lookup, determines the label stack and the outgoing core-facing interface, and forwards the packet to the MPLS core. The packet is label-switched at each hop in the core until it reaches the penultimate hop router. At this point, the top label is popped before the packet is forwarded to the egress provider edge router. The egress PE2 router performs a MPLS lookup and determines that it should remove the label before forwarding the packet to host B in site 2. Router B in site 2 receives a regular IP packet and forwards it to host B. If site 2 is also reachable over the Internet and the egress PE2 router is enabled with VRF-aware IPsec, the packet is encrypted and sent to site 2 across the Internet over an IPsec tunnel. Router B in site 2 terminates the IPsec tunnel, performs a regular IP lookup, and forwards the packet to host B.

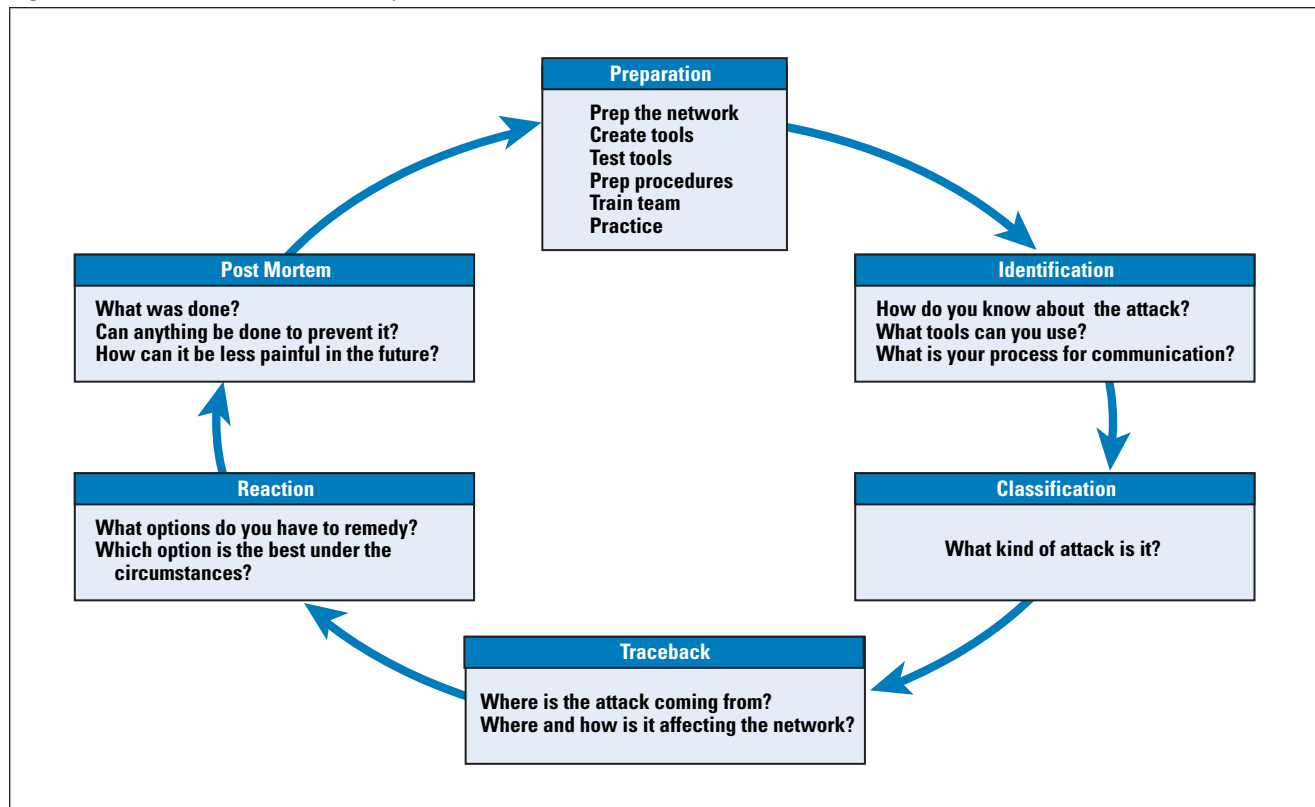
Figure 3: Managed IP VPN Security Services: IP MPLS, IPsec, or MPLS plus IPsec



Six-Step Methodology

Service providers can detect and mitigate attacks on the infrastructure using a six-step incident-response methodology (Figure 4).

Figure 4: Six Phases of Incident Response



- *Preparation:* The service provider needs to prepare the network, acquire the needed tools, develop and document a security plan, implement security procedures, and train NOC staff to use tools and procedures. *It is vital that security be a practice; the first time that the NOC staff follows its incident-response procedures should not be during an actual attack.*
- *Identification:* Unfortunately, service providers sometimes learn about a security incident from their customers. It is far better to be able to identify the threat before it becomes a problem, using NetFlow telemetry data and analysis tools, for example.
- *Classification:* The service provider needs to be able to quickly assess the nature of the threat and its scope: single customer, multiple customers, or entire infrastructure.
- *Traceback:* After classifying the threat, the IT staff needs to identify the point of ingress: peer, upstream server, downstream server, or compromised network device in the data center.
- *Reaction:* Following classification and traceback, the IT team applies the tools and processes needed to mitigate the attack. Success requires visibility into the network and well-defined procedures. Adherence to standard operating procedures helps prevent the service provider from inadvertently making the problem worse.
- *Post-mortem:* After the incident, the security team should analyze the root causes and integrate new insights into the security incident-handling procedures for use during the next incident.

Real-Life Observations About Interoperability from the APRICOT Workshops

Cisco and Juniper conducted a multivendor security workshop at APRICOT 2006 in Perth, Australia, and again at APRICOT 2007 in Bali, Indonesia. The workshops were offered in response to the fact that service providers often deploy a multivendor network for reasons ranging from financial to political.

Hands-on workshops were conducted in a lab using 12 routers running the Cisco IOS Software and another 12 running JUNOS software. Topics included:

- Password protection
- Packet filtering at the network edge
- Protecting the control plane
- Securing routing protocols
- Network monitoring techniques: NetFlow, syslog, SNMP, and *Network Time Protocol* (NTP)
- BGP MPLS Layer 3 VPNs
- IPsec VPNs

The goal of the workshops was to achieve a working configuration that interoperated with JUNOS and the Cisco IOS Software, resulting in consistent technology implementation, as well as common security policy enforcement. The workshops underscored the fact that interoperability is not automatic—even among standards-based network products. The reason is that standards bodies such as IETF, ITU, IEEE, and others define some aspects of protocols but leave others to vendor discretion. Standards do define *protocol format*, which is a syntactical structure identifying bit-field definition, length, and more. They also define *protocol behavior*, which specifies when actions occur, such as sending Hello and Keepalive timer probes and handling retransmission and reset packets. For purposes of analogy, a spoken language such as English is like a protocol format, and polite conversation conventions, such as beginning with a greeting and concluding with goodbye, is like a protocol behavior.

What standards do *not* cover are vendor-specific internal implementations, such as software coding techniques, hardware acceleration for performance, *command-line interface* (CLI) structure, and so on. Therefore, even though the APRICOT workshops involved deploying standards-based technology such as BGP-based MPLS VPNs and IPsec, vendor-specific differences had to be accounted for in the workshop materials and were noticed by participants. Following are examples noted at the APRICOT workshop:

- *Label Distribution Protocol*: With BGP MPLS VPN, JUNOS and Cisco IOS Software did not interoperate in their default configurations. However, routers from the same vendor did establish *Label Distribution Protocol* (LDP) sessions. The explanation, which participants found by troubleshooting with debug commands and referring to the manual, is that Cisco IOS Software uses the *Tag Distribution Protocol* (TDP) by default, whereas JUNOS uses LDP. After the Cisco IOS Software was changed to use LDP, the BGP-based MPLS VPN configuration succeeded.
- *IPsec tunnel establishment*: To simplify IPsec configuration, the workshop employed a *Graphical User Interface* (GUI) that prompted the user to choose source and destination IP addresses for the tunnel endpoints, a shared key, and the prefixes that defined the “interesting” traffic that was to use the IPsec tunnel. On the first attempt, the IPsec tunnel was not established. Workshop participants used the CLI to determine the problem, which was that the default encryption being negotiated was incompatible. The root cause for this mismatched encryption standard was that some routers were using an export version of software and needed an upgrade to support a higher encryption standard. Furthermore, even with common encryption capabilities, the two operating systems used different criteria to identify the interesting traffic that would be encrypted. Using the GUI, JUNOS defined interesting traffic as sourced from “ANY” network and destined to **192.168.1.0/24**.

In contrast, the Cisco IOS Software defined interesting traffic as sourced from **10.1.1.0/24** and destined to **192.168.1.0/24**. Following a discussion about whether the JUNOS default was too permissive or the Cisco IOS Software default was too restrictive, workshop participants agreed to disallow traffic that did not require encryption in the IPsec tunnel. The consensus was that the customer's security policy would provide a more conclusive answer to how permissive the policy should be, and that it was reasonable to require use of the CLI to tweak the configuration because the GUI performed most of the more difficult parts of the configuration on both platforms.

- *Loopback interface cost with Open Shortest Path First (OSPF):* During the OSPF deployment, participants noticed that the OSPF cost associated with interfaces was the same for each vendor. The OSPF cost is based upon a reference bandwidth of 100 Mbps. However, the loopback interfaces had different values: a default OSPF cost of 1 for the Cisco IOS Software and 0 for JUNOS. It is advisable to change one of the defaults to make them the same.

Although these subtle differences in protocols are documented by the vendors, service provider operational teams often have little time to research them. Therefore, it can be valuable for them to participate in multivendor hands-on workshops. Anecdotal evidence suggests that operators who are comfortable with multiple vendors understand the protocols, helping them design networks that can support new, revenue-generating services.

It is hoped that events such as the APRICOT workshops will help build a community of professionals who can add value for their employers, each other, and the broader Internet community. The result will be a secure and trusted networking environment that people and industry can rely on and use to connect in new and innovative ways.

Summary

Managed security services represent a growing revenue opportunity for service providers. Most service providers operate in a multivendor environment, either because of mergers and acquisitions or because their customers use other vendors' equipment. Therefore, a standards-based approach positions providers to capitalize on the managed security service opportunity. Providers can secure their infrastructure in a multivendor environment by following best practices for point protection, edge protection, RTBH protection, source-address validation, control-plane protection, and total visibility into network activity.

References

- [1] S. Kent and R. Atkinson, "Security Architecture for the Internet Protocol," RFC 2401, November 1998.
- [2] S. Kent and R. Atkinson, "IP Authentication Header," RFC 2402, November 1998.
- [3] S. Kent and R. Atkinson, "IP Encapsulating Security Payload (ESP)," RFC 2406, November 1998.
- [4] E. Rosen and Y. Rekhter, "BGP/MPLS VPNs," RFC 2547, March 1999.
- [5] S. Hanks, T. Li, D. Farinacci, and P. Traina, "Generic Routing Encapsulation (GRE)," RFC 1701, October 1994.
- [6] Y. Rekhter, B. Moskowitz, D. Karrenberg, G. J. de Groot, and E. Lear, "Address Allocation for Private Internets," RFC 1918, February 1996.
- [7] Internet Assigned Numbers Authority (IANA), "Special-Use IPv4 Addresses," RFC 3300, September 2002.
- [8] F. Baker and P. Savola, "Ingress Filtering for Multihomed Networks," RFC 3704, March 2004.
- [9] D. Harkins and D. Carrel, "The Internet Key Exchange (IKE)," RFC 2409, November 1998.
- [10] Convery, S., "Network Authentication, Authorization, and Accounting – Part One: Concepts, Elements, and Approaches," *The Internet Protocol Journal*, Volume 10, No. 1, March 2007.
- [11] Convery, S., "Network Authentication, Authorization, and Accounting – Part Two: Protocols, Applications, and the Future of AAA," *The Internet Protocol Journal*, Volume 10, No. 2, June 2007.
- [12] Rigney et. al., "Remote Authentication Dial- In User Service (RADIUS)," RFC 2865 (Obsoletes RFC 2138, and RFC 2058), June 2000.
- [13] Carrel et al., "The TACACS+ Protocol Version 1.78," Internet Draft, Work in Progress, **draft-grant-tacacs-02.txt**, January 1997.
- [14] <http://www.shrubbery.net/rancid/>
- [15] <http://www.apricot.net>

KUNJAL TRIVEDI joined Cisco in 1999 as a consulting engineer initially and then worked in product management covering Cisco IOS Software infrastructure security. Currently, he is helping Cisco shape a Managed Security Services marketing vision and strategy. A widely respected networking security expert, Kunjal presents infrastructure security, IP Security, and Managed Security topics at Cisco Networkers events as well as at conferences such as APRICOT. Kunjal has a Bachelor of Engineering degree with honors in electrical and electronics engineering from University of Wales, College of Cardiff, and a Master of Science degree in Artificial Intelligence from Cranfield Institute of Technology, UK. He holds CISSP and CCIE designations in routing and switching as well as security. Recently, he published a book titled *[Read Me First]: Building or Buying VPNs*; Kunjal has been awarded Chartered Engineer status by Institute of Engineering and Technology. He can be reached at kunjal@cisco.com

DAMIEN HOLLOWAY joined Juniper Networks in 2004 as an Instructing Engineer. He contributes to the development of the Juniper Technical Certification Program and custom delivery of training in the Asia Pacific region. Previously he was a consulting engineer and provided design, installation, and training to providers in Australia and the United States. Damien has presented a wide variety of topics relevant to customers, including backbone design, application acceleration, and *Broadband Remote Access Server* edge design, to audiences, including APRICOT and SANOG. Damien has a Bachelor of Electrical Engineering and Bachelor of Science from University of Sydney, Australia. He is a CCIE expert in routing and switching and JNCIE-M, JNCIP-E, and CISSP. He can be reached at holloway@juniper.net

Kunjal Trivedi (left) and Damien Holloway (center) share a joke with workshop students at APRICOT 2007



Kunjal with APRICOT 2007 workshop attendees



IPv4 Address Depletion and Transition to IPv6

by Geoff Huston, APNIC

At the recent APNIC meeting in New Delhi, the subject of IPv4, IPv6, and transition mechanisms was highlighted in the plenary session^[1]. This article briefly summarizes that session and the underlying parameters in IPv4 address depletion and the transition to IPv6.

IPv4 Status

As of September 2007 we have some 18 percent of the unallocated IPv4 address pool remaining with the *Internet Assigned Numbers Authority* (IANA), and 68 percent has already been allocated to the *Regional Internet Registries* (RIRs) and through the RIRs to *Internet Service Providers* (ISPs) and end users. The remaining 14 percent of the IPv4 address space is reserved for private use, multicast, and special purposes. Another way of looking at this situation is that we have exhausted four-fifths of the unallocated address pool in IPv4, and one-fifth remains for future use. It has taken more than two decades of Internet growth to expend this initial four-fifths of the address space, so why shouldn't it take a further decade to consume what remains?

At this point the various predictive models come into play, because the history of the Internet has not been a uniformly steady model. The Internet began in the 1980s very quietly; the first round of explosive growth in demand was in the early 1990s as the Internet was adopted by the academic and research sector. At the time, the address architecture used a model where class A networks (or a /8) were extremely large, the class B networks (/16) were also too large, and the class C networks (/24) were too small for most campuses. The general use of class B address blocks was an uncomfortable compromise between consuming too much address space and consuming too many routing slots through address fragmentation. The subsequent shift to a classless address architecture in the early 1990s significantly reduced the levels of IPv4 address consumption for the next decade. However, over the past five years the demand levels for addresses have been accelerating again. Extensive mass-market broadband deployment, the demand for public non-*Network Address Translation* (NAT) addresses for applications such as *Voice over IP* (VoIP), and continuing real cost reductions in technology that has now brought the Internet to large populations in developing economies all contribute to an accelerating IPv4 address consumption rate.

Various approaches to modeling this address consumption predict that the IANA unallocated address pool will be fully depleted sometime in 2010 or 2011^[2, 3, 4, 5].

Transitioning to IPv6

The obvious question is “What then?”, and the commonly assumed answer to that question is one that the *Internet Engineering Task Force* (IETF) started developing almost 15 years ago, namely a shift to use a new version of the Internet Protocol: what we now know as IP Version 6, or IPv6. But if IPv6 really is the answer to this problem of IPv4 unallocated address-pool depletion, then we appear to be leaving the transition process quite late. The uptake of IPv6 in the public Internet remains extremely small as compared to IPv4^[6]. If we really have to have IPv6 universally deployed by the time we fully exhaust the unallocated IPv4 address pools, then this objective appears to be unattainable during the 24 months we have to complete this work. The more likely scenario we face is that we will not have IPv6 fully deployed in the remaining time, implying a need to be more inventive about IPv4 in the coming years, as well as inspecting more closely the reason why IPv6 has failed to excite much reaction on the part of the industry to date.

We need to consider both IPv4 and IPv6 when looking at these problems with transition because of an underlying limitation in technology: *IPv6 is not “backward-compatible” with IPv4*. An IPv6 host cannot directly communicate with an IPv4 host. The IETF worked on ways to achieve this through intermediaries, such as a protocol to translate NATs^[7], but this approach has recently been declared “historic” because of technical and operational difficulties^[8]. That decision leaves few alternatives. If a host wants to talk to the IPv4 world, it cannot rely on clever protocol translating intermediaries somewhere, and it needs to have a local IPv4 protocol stack, a local IPv4 address, and a local IPv4 network and IPv4 transit. And to speak to IPv6 hosts, IPv6 has the same set of prerequisites as IPv4. This approach to transition through replication of the entire network protocol infrastructure is termed “Dual Stack.” The corollary of Dual Stack is continued demand for IPv4 addresses to address the entire Internet for as long as this transition takes. The apparent contradiction here is that we do not appear to have sufficient IPv4 addresses in the unallocated address pools to sustain this Dual Stack approach to transition for the extended time periods that we anticipate this process to take.

What Can We Expect?

So we can expect that IPv4 addresses will continue to be in demand well beyond any anticipated date of exhaustion of the unallocated address pool, because in the Dual Stack transition environment all new and expanding network deployments need IPv4 service access and addresses. But the address distribution process will no longer be directly managed through address allocation policies after the allocation pool is exhausted.

Ideas that have been aired in address policy forums include encouraging NAT deployment in IPv4, expanding the private use of IPv4 address space to include the last remaining “reserved-for-future-use” address block, various policies relating to rationing the remaining IPv4 address space, increased efforts of address reclamation, the recognition of address transfers, and the use of markets to support address distribution.

Of course the questions here are about how long we need to continue to rely on IPv4, how such new forms of address distribution would affect existing notions of fairness and efficiency of use, and whether this effect would imply escalation of cost or some large-scale effect on the routing system.

On the other hand, is IPv6 really ready to assume the role of the underpinning of the global Internet? One view is that although the transition to a universal deployment of IPv6 is inevitable, numerous immediate concerns have impeded IPv6 adoption, including the lack of backward compatibility and the absence of simple, useful, and scalable translation or transition mechanisms^[9]. So far the business case for IPv6 has not been compelling, and it appears to be far easier for ISPs and their customers to continue along the path of IPv4 and NATs.

When we contemplate this transition, we also need to be mindful of what we need to preserve across this transition, including the functions and integrity of the Internet as a service platform, the functions of existing applications, the viability of routing, the capability to sustain continued growth, and the integrity of the network infrastructure.

It appears that what could be useful right now is clear and coherent information about the situation and current choices, and analyzing the implications of various options. When looking at such concerns of significant change, we need to appreciate both the limitations and the strengths of the Internet as a global deregulated industry and we need, above all else, to preserve a single coherent networked outcome. Perhaps this topic is far broader than purely technical, and when we examine it from a perspective that embraces economic considerations, business imperatives, and public policy objectives, we need to understand the broader context in which these processes of change are progressing^[10].

It is likely that some disruptive aspects of this transition will affect the entire industry, and this transition will probably be neither transparent nor costless.

References

- [1] APNIC 24 Plenary Session: “The Future of IPv4,” September 2007. <http://www.apnic.net/meetings/24/program/plenaries/apnic/>
- [2] Geoff Huston, “The IPv4 Report.” <http://ipv4.potaroo.net>
- [3] Tony Hain, “IPv4 Address Pool.” <http://www.tndh.net/~tony/ietf/ipv4-pool-combined-view.pdf>
- [4] Tony Hain, “A Pragmatic Report on IPv4 Address Space Consumption,” *The Internet Protocol Journal*, Vol. 8, No. 3, September 2005. http://www.cisco.com/web/about/ac123/ac147/archived_issues/ipj_8-3/ipv4.html
- [5] K.C. Claffy, CAIDA, “ ‘Apocalypse Then’: IPv4 Address Space Depletion,” Presentation to ARIN XVI, October 2005. http://www.arin.net/meetings/minutes/ARIN_XVI/PDF/wednesday/claffy_ipv4_roundtable.pdf
- [6] Geoff Huston, “IPv6 / IPv4 Comparison Metrics.” <http://bgp.potaroo.net/v6/v6rpt.html>
- [7] G. Tsirtsis and P. Srisuresh, “Network Address Translation – Protocol Translation (NAT-PT),” RFC 2766, February 2000.
- [8] C. Aoun and E. Davies, “Reasons to Move the Network Address Translator – Protocol Translator (NAT-PT) to Historic Status.” RFC 4966, July 2007.
- [9] Randy Bush, “IPv6 Operational Reality,” APNIC 24 Plenary Presentation, September 2007. <http://www.apnic.net/meetings/24/program/plenaries/apnic/presentations/bush-ipv6-op-reality.pdf>
- [10] Geoff Huston, “IPv4 Exhaustion,” APNIC 24 Plenary Presentation, September 2007. <http://www.apnic.net/meetings/24/program/plenaries/apnic/presentations/huston-ipv4-exhaustion.pdf>

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. The author of numerous Internet-related books, he is currently the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He was a member of the Internet Architecture Board from 1999 until 2005, and served on the Board of the Internet Society from 1992 until 2001. E-mail: gih@apnic.net

IPv4 Address Space: 2.46 Billion Down, 1.25 Billion to Go

by Iljitsch van Beijnum

In September 2005, *The Internet Protocol Journal* published an article about the IPv4 address space consumption^[1]. At that time, projections done by Geoff Huston and Tony Hain varied widely, because the number of /8 address blocks in use had gone up sharply in early 2005. So what has happened since then, and what can we expect for the not-too-distant future?

Address Assignment and Allocation

The *Internet Assigned Numbers Authority* (IANA, part of the *Internet Corporation for Assigned Names and Numbers* [ICANN]) has authority over the IPv4 address space. In the past, IANA gave out address blocks directly to end users, but now IANA distributes address space in the form of /8 blocks, each holding 24 bits worth of address space, or 16,777,216 addresses, to five *Regional Internet Registries* (RIRs). There are a few exceptions, but AfriNIC^[2] gives out address space in Africa; APNIC^[3] in the Asia-Pacific region; ARIN^[4] in North America; LACNIC^[5] in Latin America and the Caribbean; and RIPE NCC^[6] in Europe, the former Soviet Union, and the Middle East. These RIRs sometimes assign address space to end users, but mostly allocate it to *Internet Service Providers* (ISPs), who then assign it to their customers, meaning that there are two pools of available address space: the global pool of /8 blocks that IANA has not delegated to anyone^[7], and the address space held by the RIRs that they have not given out yet. The article in the September 2005 issue of *The Internet Protocol Journal*^[1] looked at the depletion of the IANA global pool, whereas this article mostly looks at the amounts of address space given out by the RIRs, providing a more granular view. The RIRs publish daily reports of their address assignments and allocations on their respective FTP servers. According to these reports as downloaded on January 1, 2007, the amounts of address space shown in Table 1 were given out over the past seven years.

Table 1: Address Space Allocated 2000–2006 [January 2007 data]

	2000	2001	2002	2003	2004	2005	2006
AfriNIC	0.56	0.39	0.26	0.22	0.51	1.03	2.72
APNIC	20.94	28.83	27.03	33.05	42.89	53.86	51.78
ARIN	30.83	28.55	21.08	22.32	34.26	47.57	38.94
LACNIC	0.88	1.61	0.65	2.62	3.77	10.97	11.50
RIPE NCC	24.79	25.36	19.84	29.61	47.49	62.09	56.53
Total	78.00	84.73	68.87	87.82	128.92	175.52	161.48

However, if we compare these totals to the totals seen on January 1, 2006, we see some differences (Table 2).

Table 2: Address Space Allocated 2000–2006 [January 2006 data]

	2000	2001	2002	2003	2004	2005	2006
Total	78.35	88.95	68.93	87.77	128.45	165.45	–

For the years 2000 to 2002, the number of addresses registered as given out is slightly lower, as seen in the January 1, 2007 data compared to the January 1, 2006 data—a result that is to be expected because address space given out in that year that is no longer used is returned. However, for the later years, and especially for 2005, there is a retroactive *increase* in the number of addresses given out. The reason: When ARIN suspects an address space user may come back for more space relatively soon, it takes a larger block than requested, and then fulfills the request from part of that block and keeps the rest in reserve. So an organization requesting a /16 may get the first half of a /15. When that organization then requests another /16 one or two years later, ARIN gives the organization the second half of the /15. ARIN subsequently records this as a /15 given out on the date when the original /16 was requested.

For instance, ARIN’s January 1, 2006, data shows that a block of 12.6 million addresses was given out within **73.0.0.0/8** block:

arin|US|ipv4|73.0.0.0|12582912|20050419|allocated

In the January 1, 2007, data, this number had changed to 13.6 million addresses:

arin|US|ipv4|73.0.0.0|13631488|20050419|allocated

This change means that simply looking at the registration date does not provide very good information. It also does not account for address space given out in earlier years that is returned. An alternative approach is to count the amount of address space given out based on the RIR records published on a certain date (Table 3).

Table 3: RIR Records for Address Space Allocation

	IANA (/8)		RIRs (millions)			Total	
	Delegated	Free	Received	Delegated	Free	Free	Delta
Jan. 1, 2004	133	88	1509.95	1245.63	264.32	1740.71	
Jan. 1, 2005	142	79	1660.95	1351.66	309.30	1634.69	106.02
Jan. 1, 2006	155	66	1879.05	1517.74	361.31	1468.61	166.08
Jan. 1, 2007	166	55	2063.60	1685.69	377.90	1300.65	167.96
May 1, 2007	172	49	2181.04	1754.68	426.36	1248.44	52.21

(Note that block **7.0.0.0/8** shows up as unused in the IANA global pool and is counted as available in the table, but this block is in fact used by the U.S. Department of Defense.)

The jump in address consumption between 2004 (106 million) and 2005 (166 million) is even more dramatic in this light, while consumption numbers of 2005 and 2006 (168 million) are now almost identical. The figure for the first four months of 2007 seems rather modest at 52 million addresses, but the reason lies in the fact that Bolt, Beranek and Newman returned **46.0.0.0/8** to IANA in April. So the number of addresses given out from January to April was 69 million, a rate that puts the RIRs on track to give out more than 200 million addresses in 2007.

The size of address blocks given has been increasing steadily. Table 4 shows the number of requests for a certain range of block sizes: equal or higher than the first, lower than the second value (2005 and earlier values from the January 1, 2006 data, 2006 values from the January 1, 2007 data).

Table 4: Number of Requests for Ranges of Block Sizes

	2000	2001	2002	2003	2004	2005	2006
< 1,000	326	474	547	745	1022	1309	1526
1,000 – 8,000	652	1176	897	1009	1516	1891	2338
8,000 – 64k	1440	868	822	1014	1100	1039	1133
64k – 500k	354	262	163	215	404	309	409
500k – 2M	19	39	29	46	61	60	56
> 2M	3	5	5	6	7	18	13

The number of blocks in the two smallest categories has increased rapidly, but not as fast as the number of blocks in the largest category, in relative numbers. However, the increase in large blocks has a very dramatic effect whereas the small blocks are insignificant, when looking at the millions of addresses involved (Table 5).

Table 5: Millions of Addresses Given Out

	2000	2001	2002	2003	2004	2005	2006
< 1,000	0.10	0.16	0.18	0.25	0.35	0.44	0.52
1,000 – 8,000	2.42	4.47	3.23	3.45	4.49	5.07	6.10
8,000 – 64k	18.79	12.81	11.35	14.00	15.99	15.46	17.17
64k – 500k	35.98	32.19	20.28	25.51	42.01	34.23	49.64
500k – 2M	12.68	24.64	21.30	31.98	44.63	41.63	46.64
> 2M	8.39	14.68	12.58	12.58	20.97	68.62	41.42

The increase in the 2M+ blocks was solely responsible for the high number of addresses given out in 2005. In 2006, there was growth in all categories except the 2M+ one (even the 500k – 2M category increased in number of addresses if not in number of blocks). When the 2M+ blocks are taken out of the equation, 2005 had a total of 96.83 million addresses (January 1, 2006) and 2006 had 119.06 million given out, even without fully correcting for the ARIN reporting particularities. Apparently there is still an underlying upward trend.

Figure 1 shows the amounts of address space given out by IANA and by the RIRs every year from 1994 to 2006.

Figure 1: IPv4 Address Space Given Out from 1994 to 2006

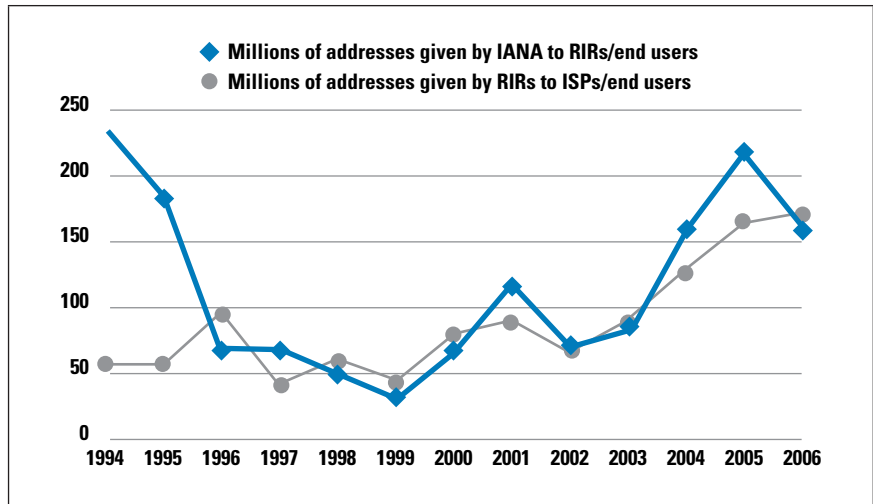
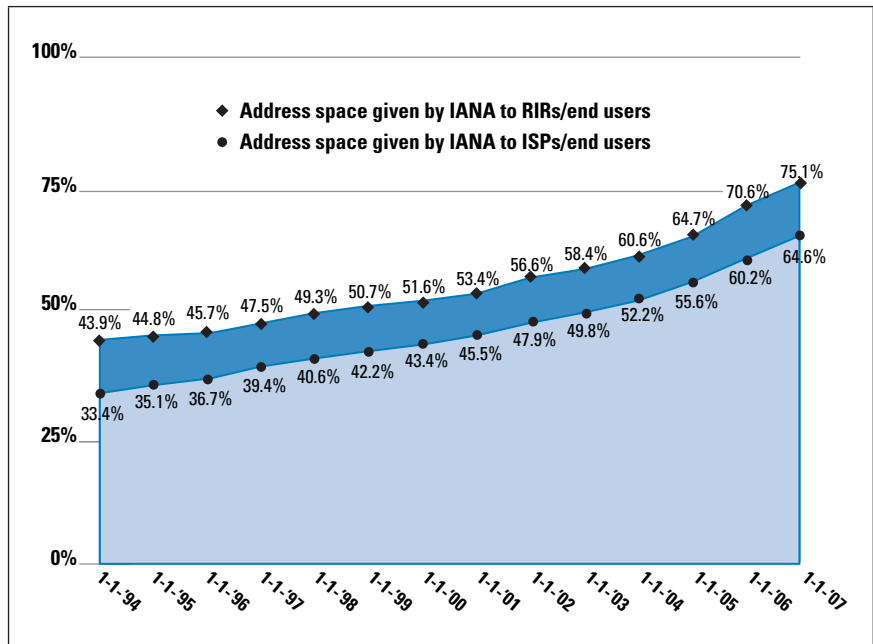


Figure 2 shows the amounts of address space marked as “in use” by IANA and by the RIRs. The difference between the two numbers is what the RIRs hold in order to satisfy day-to-day address space requests. This amount is usually two years’ worth of address space.

Figure 2: IPv4 Address Space in Use from 1994 to 2006



Depletion

The exact moment when the IPv4 address space will be depleted depends on numerous factors. Since 1997, the three-year period with the largest growth in yearly address use was from 2003 to 2005 relative to the 2002 figure: a factor 2.4, or 34 percent per year. If this growth repeats itself in the next three years, we will be out of IPv4 addresses in the second half of 2010.

Interestingly, the period with the lowest growth also includes the year 2003: from 2001 to 2003 relative to 2000. In 2003, 12 percent more addresses were given out than in 2000, for an average increase in yearly use of 4 percent. If this is the new trend for the coming years, we can expect to run out of IPv4 addresses in mid-2013. There is of course no reason to assume that future IP address use will conform to patterns seen in earlier years, but we really have nothing else to base our projections upon.

So anyone expecting to obtain new IPv4 address space more than three years from now is taking a big risk. With the IPv4 reserves visibly diminishing each year, the question is: What can we, as a community, do to make the IPv4 address depletion as painless as possible? IPv4 addresses are useful only if the people who need them can obtain them, meaning that using up addresses unnecessarily fast or locking up the still-available reserves are both suboptimal solutions. It has been suggested that turning IPv4 address space into a tradable commodity would allow a free market to form, aiding the efficient distribution of address space from those who have it to those who need it.

This scenario has several problems. First, when supply is limited and demand is high, prices rise and hoarding becomes lucrative. So the effect of making address space tradable could be a reduction of available address space rather than an increase. And certainly, as trading IPv4 space becomes more likely, holders of large address blocks will be less inclined to return them. Finally, more than half of the IPv4 address space in use is held by organizations in the United States, whereas the developing world has comparatively little address space. The prospect of having to buy address space from American companies that got the space for free is not likely to be popular in the rest of the world.

Address Reclamation a Solution?

There are two large classes of potentially reclaimable address space: the class E reserved space (**240.0.0.0 – 255.255.255.255**) and the class A blocks given out directly to end users by IANA. The class E space has 268 million addresses and would give us in the order of 18 months worth of IPv4 address use. However, many TCP/IP stacks, such as the one in Windows, do not accept addresses from class E space and will not even communicate with correspondents holding those addresses. It is probably too late now to change this behavior on the installed base before the address space would be needed. There are currently 42 class A blocks and another two /8s from class C space listed as given out to end users—738 million addresses. The U.S. government uses about 10 of those blocks; 21 of them are not present in the *Border Gateway Protocol* (BGP) routing table.

Although harsh judgments about the need for so much address space are easily made from the outside without having all the pertinent information, it seems reasonable to try to reclaim some of this space. I would consider getting back half of this space a big success, but that would give us only 2 years worth of additional address space. There are also 645 million addresses of older class B assignments, but reclaiming those will be extremely difficult because nearly 8,000 individual assignments are involved. Reclaiming a class B block is probably not much easier than reclaiming a class A block, but the amount of address space returned is less than half a percent.

Planning for the End Game

So what should we do? In my opinion: promote predictability. The situation where we run out of IPv4 address space much faster than expected would be very harmful as organizations struggle to adjust to the new circumstances. On the other hand, if the IPv4 space unexpectedly lasts longer, people may be disinclined to believe space is really running out and then would be unprepared when it does. Artificially delaying running out of IPv4 address space also prolongs the situation in which it is difficult to get IPv4 space, but not enough people feel the pain to initiate IPv6 deployment. One solution worthy of consideration would be to impose a worldwide moratorium on the change of IPv4 address allocation and assignment policies after a certain date to aid this predictability. If some kind of encouraged or forced reclamation of older class A blocks is desired, this process should be instigated sooner rather than later, both for the sake of predictability and because it gives the address holders involved time to reorganize their networks. Another small but useful step would be to limit the size of address blocks given out. This scenario would be like the agreement between the RIRs and IANA that the RIRs will receive two /8s at a time in the future. The situation where a single /9 or /8 allocation constitutes 5 or even 10 percent of the address space given out in that year makes adequate predictions extremely difficult, and also runs the risk that a good part of the address block in question will never be used as circumstances change. Limiting individual allocations to a /11 or /12 would be better, even if it requires the requesting organization to come back for more address space several times per year.

Finally, it seems prudent for all organizations using public IPv4 address space to start planning for the moment that they themselves, or third parties that they communicate with over the public Internet, can no longer obtain additional IPv4 address space.

References

- [1] Tony Hain, "A Pragmatic Report on IPv4 Address Space Consumption," *The Internet Protocol Journal*, Volume 8, No. 3, September 2005.
- [2] AfriNIC, <http://www.afrinic.net>
- [3] APNIC, <http://www.apnic.net>
- [4] ARIN, <http://www.arin.net>
- [5] LACNIC, <http://www.lacnic.net>
- [6] RIPE NCC, <http://ripe.net>
- [7] IANA Internet Protocol v4 Address Space
<http://www.iana.org/assignments/ipv4-address-space>

ILJITSCH VAN BEIJNUM holds a Bachelor of Information and Communication Technology degree from the Haagse Hogeschool in The Hague, Netherlands. In 1995, he found himself in the emerging Internet Service Provider business. There he learned about system administration, IP networking, and especially routing. After first starting a small ISP with four others and working as a senior network engineer for UUNET Netherlands, he became a freelance consultant in 2000. Not long after that, he started contributing to the IETF Multihoming in IPv6 working group. He wrote the book *BGP: Building Reliable Networks with the Border Gateway Protocol*, ISBN 0-596-00254-8, published by O'Reilly in 2002, and *Running IPv6*, ISBN 1590595270, published by Apress in 2005. E-mail: iljitsch@muada.com

Used but Unallocated: Potentially Awkward /8 Assignments

by Leo Vegoda, ICANN

IPv4 has proven to be exceedingly popular, so it should be no surprise that the time is rapidly approaching when the last /8 block will be allocated and the *Internet Assigned Numbers Authority's* (IANA's) free pool of address space will be empty. At the time of writing, Geoff Huston of the *Asia Pacific Network Information Centre* (APNIC) is projecting^[1] the IANA free pool will run out in mid-2010. Unfortunately, it is possible that some of these remaining /8s may cause problems for enterprise and *Internet Service Provider* (ISP) network operators when they are put back into use. These blocks are not the /8s that have been returned to IANA by the original registrants; they are previously unassigned address blocks.

Concerns

There are many concerns about the IANA free pool depletion, but one of them seems particularly straightforward to identify and fix. Many organizations have chosen to use unregistered IPv4 addresses in their internal networks and, in some cases, network equipment or software providers have chosen to use unregistered IPv4 addresses in their products or services. In many cases the choice to use these addresses was made because the network operators did not want the administrative burden of requesting a registered block of addresses from a *Regional Internet Registry* (RIR)^[2, 11]. In other cases they may not have realized that RFC 1918^[3] set aside three blocks of address space for private networks, so they just picked what they believed to be an unused block, or their needs exceeded the RFC 1918 set-aside blocks. Other organizations used the default address range suggested by their equipment vendor, or supplied in example documentation, when configuring *Network Address Translation* (NAT) devices. Regardless of the reason, these uses of unregistered addresses will conflict with routed addresses when the /8s in question are eventually assigned to ISPs or enterprise users.

A few examples of /8s where problems are likely to occur follow:

- | | |
|-------------------|--|
| 1.0.0.0/8 | Widely used as private address space in large organizations whose needs exceed those provided for by RFC 1918 ^[4] |
| 5.0.0.0/8 | Used by one of numerous zero-configuration Internet applications (including the Hamachi VPN service ^[5, 6]) |
| 42.0.0.0/8 | Default range used in the NAT configuration of at least one Internet appliance (the HP Procurve 700w ^[7]) |

Organizations using these address ranges in products or services may experience problems when more specific Internet routes attract traffic that was meant for internal hosts, or alternatively find themselves unable to reach the legitimate users of those addresses because those addresses are being used internally. The users of unregistered networks may also find problems with reverse *Domain Name System* (DNS) resolution, depending on how their DNS servers are configured. These problems are likely to result in additional calls to helpdesks and security desks at both enterprises and ISPs, with unexpected behavior for end users that might be hard to diagnose. Users of unregistered address space may also experience problems with unexpected traffic being received at their site if they leak internal routes to the public Internet. Many ISPs have already had experience with this type of routing inconsistency as recent /8 allocations reach routing tables and bogon filters are updated.

Alternatives

There are several alternatives to using unregistered IPv4 address space:

- Use RFC 1918 IPv4 address space (no need to obtain this space from an RIR)
- Use IPv4 address space registered with an RIR
- Use IPv6 address space registered with an RIR
- Use IPv6 Unique Local Address^[8] space (no need to obtain this space from an RIR)

Obviously, all of these efforts will involve renumbering networks, a sometimes painful and time-consuming process. Those using unregistered unique IPv4 address space should look at renumbering their networks or services before the previously unallocated /8s are allocated to avoid address clashes and routing difficulties.

Additionally, vendors and documentation writers can clean up their configurations to ensure they use RFC 1918 addresses, or make it clear to their users that they must use registered addresses to avoid routing conflicts.

All RIRs provide free telephone helpdesks that can advise you on obtaining unique IPv4 or IPv6 address space. But if you want to continue using unregistered space and can transition to IPv6, the prefix selection mechanism described in RFC 4193 makes the probability of a clash a mere 1 in 550 billion. Ultimately, transitioning to IPv6 is most likely the best solution, and this approach offers an opportunity for those having to renumber parts of their network to avoid a subsequent renumbering later into IPv6.

About IANA and ICANN

IANA allocates address space to RIRs according to the global IPv4 [9] and IPv6^[10] policies. Enterprise and ISP networks need to obtain IP addresses from their upstream provider or from the appropriate RIR.

The *Internet Corporation for Assigned Names and Numbers* (ICANN) is an internationally organized, nonprofit corporation that has responsibility for *Internet Protocol* (IP) address space allocation, protocol identifier assignment, *generic* (gTLD) and *country code* (ccTLD) *Top-Level Domain* name system management, and root server system management functions. These services were originally performed under U.S. government contract by IANA and other entities. ICANN now performs the IANA function.

References

- [1] <http://www.potaroo.net/tools/ipv4/>
- [2] RFC 1174, para 2.2 states in part, “The term Internet Registry (IR) refers to the organization which has the responsibility for gathering and registering information about networks to which identifiers (network numbers, autonomous system numbers) have been assigned by the IR. An RIR does this function for its service area.”
- [3] <http://www.ietf.org/rfc/rfc1918.txt>
- [4] <http://tools.ietf.org/id/draft-hain-1918bis-01.txt>
- [5] <https://secure.logmein.com/products/hamachi/howitworks.asp>
- [6] <http://en.wikipedia.org/wiki/Hamachi>
- [7] <http://www.hp.com/rnd/support/faqs/700w1.htm>
- [8] <http://www.ietf.org/rfc/rfc4193.txt>
- [9] <http://www.icann.org/general/allocation-IPv4-rirs.html>
- [10] <http://www.icann.org/general/allocation-IPv6-rirs.htm>
- [11] Daniel Karrenberg, Gerard Ross, Paul Wilson, and Leslie Nobile, “Development of the Regional Internet Registry System,” *The Internet Protocol Journal*, Volume 4, No. 4, December 2001.

LEO VEGODA holds a BA (Hons) from the University of Central England. He joined ICANN in 2006 and is the Manager, Number Resources - IANA. He has previously worked for the RIPE NCC, where he ran the Registration Services department. He can be reached at: leo.vegoda@icann.org

Book Review

Uncommon Sense

Uncommon Sense: Out of the Box Thinking for An In the Box World, By Peter Cochrane, ISBN 1-84112-477-x, Published by Capstone, 2004, <http://www.wileyurope.com>

A series of articles published in **silicon.com** form the basis for this book, which looks at the effect that new technology has on business and its implications for society. In many ways it attacks conventional wisdom and forces a reevaluation of the effect of technology, often exposing flaws in the business logic that lead to many investments and decisions.

The book is aimed at technologists, managers, and professionals who are interested in change and progress, offering them a glimpse of the future. It is easy to read, with liberal use of figures and tables to aid understanding.

Organisation

Cochrane begins by looking at the communication of ideas, particularly fairly complex and novel concepts. He notes the lack of agreement on the major concerns of the future and bemoans the handling of complex business and political topics—and the lack of engineering type rigour applied to their assessment. He suggests a much more rigorous modeling of complex business problems is required, especially of business processes, which are typically complex and inter-related, so treating them as isolated “stovepipes” is inappropriate and error-prone. Cochrane emphasises the need for nonlinear thinking.

Cochrane’s analysis continues with an assessment of technology markets, not surprisingly beginning with the forces behind the dot-com bubble, with particular reference to the effect that the so-called new and old economies have had on each other. He suggests that short-term approaches, with their tendency to hit high-visibility symptoms and not the underlying commercial factors, are a barrier to progress. Cochrane reflects that whilst the dot-com boom is over, it is now clear that the online world has been very successful and has dragged the old world along in its wake.

The book then looks at change: considering the adoption of new technology and the impact effect of the Internet, comparing this new technology with the adoption of television. Cochrane spends a significant amount of time on both entertainment and learning. He examines topics as varied as security, the ease of movement of information across borders, and the role of specialist and general devices.

His assessment of security considers the range and rate of spread of threats and some advanced countermeasures such as biometrics. He considers the nature of change programmes and the harmful ways insensitive micromanagement can affect their progress.

Cochrane explores the role of the consumer in deciding which technical innovations survive, as exemplified by the growth of the American cable TV (CATV) market. He notes that most consumers have a fixed level of disposable income and new innovations allow them to redirect rather than increase their level of spending. Cochrane argues that this truth is reflected in the saturation within the mobile handset market and the dynamics seen between the media companies and new innovators such as Napster.

The penultimate collection of essays considers the speed of innovation. Cochrane notes that many consumers are suffering from “technology fatigue” and many products are suffering from “feature death.” Here he discusses stagnation within the mobile market and disillusionment with the *Wireless Application Protocol* (WAP), *General Packet Radio Service* (GPRS), and Bluetooth. He notes that the adoption of technology is linked to the willingness of customers to pay.

Cochrane concludes by looking at leading-edge variables, including reliability, noting that this variable goes hand-in-hand with maturity, with the *Public Switched Telephone Network* (PSTN) delivering extremely high levels of reliability and most modern IT solutions delivering considerably less. He makes this comparison a critical test of the five-nines availability claims of many new technology solutions. Cochrane looks at some more less-conventional ideas such as the replication of ant logic in IT systems and the possible future use of plasma screens and voice recognition as convenient input/output devices. He notes the increasing intelligence of devices, but also acknowledges that rapid communications and minimal hierarchy can triumph over better organised structures as demonstrated by protesters in France in 2000 and 2001.

Synopsis

Cochrane takes the reader through many contemporary technology developments and concerns and in the process invites his readers to form their own views. His mission is to “communicate the implications of what we have done, are doing and are about to do.” In 50 short articles, delivered in 233 pages, it is possible for the author to cover only a small portion of a rapidly growing field, providing sufficient detail to appeal to the technologist without losing the bigger picture. He examines the implications of new technology for society and notes that the progress we are seeing means that we have to take on the new, changing the way we manage, operate, and govern our businesses as a result.

The Author

Peter Cochrane is the ex-BT Chief Technologist, who with a group of ex-Apple Computer technologists founded Concept Labs, where he advises a range of companies across the world. He has published widely, holds B.Sc., M.Sc., Ph.D., and D.Sc. degrees from Nottingham (Trent) and Essex Universities, is an Apple Master, and is a visiting professor at London, Essex, and Southampton Universities. He is best known for his incisive and often provocative views on the United Kingdom and world telecommunications industries.

—Edward Smith, BT, UK

edward.a.smith@btinternet.com

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. Contact us at **ipj@cisco.com** for more information.

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2007 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA
--

The Internet Protocol Journal

December 2007

Volume 10, Number 4

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
IP Spoofing	2
Security Standards	10
Looking Toward the Future	23
Remembering Itojun	32
Book Review	35
Fragments	39
Call for Papers	43

Identity theft is a widely reported problem in today's world. Criminals can use numerous ways to obtain private information such as Social Security Numbers, credit card details, and other information that makes it possible for the perpetrator to successfully "pretend to be" someone else. A similar concept, albeit less personal, is so-called *IP Spoofing*, wherein fake IP datagrams can be generated and sent across the network in order to compromise remote systems in a variety of ways. Farha Ali gives an overview of IP Spoofing and explains ways in which the problem can be mitigated.

Our second article looks at numerous standards for information security management being developed by organizations such as the *International Organization for Standardization* (ISO), the *National Institute of Standards and Technology* (NIST), and others. The author of the article is William Stallings.

On November 2, 2007, Vint Cerf ended his term as chairman of the *Internet Corporation for Assigned Names and Numbers* (ICANN). At the same time he released a document entitled "Looking Toward the Future," which details ICANN's history, as well as outlining its challenges ahead. We've included the document in this issue and added some pointers for those readers who may not be familiar with the workings of ICANN.

In late October, the Internet technical community received the sad news that Dr. Junichiro Hagino, universally known as "Itojun" had passed away. Itojun played a very important role in the development of IPv6 and had many friends across the world. We asked one of them, Bob Hinden, to reflect on Itojun's life and compile some comments from those who knew him well.

We would like to remind you about our online adjunct to this journal. *The Internet Protocol Forum* (IPF) available at <http://www.ipjforum.org/> is a resource you can use to discuss articles and read additional material. Please take a moment to explore IPF.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

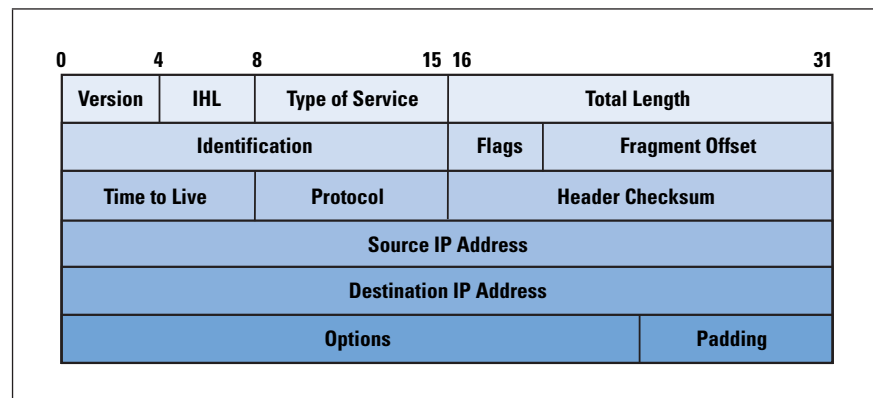
You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

IP Spoofing

by Farha Ali, Lander University

The *Internet Protocol*, or IP, is the main protocol used to route information across the Internet. The role of IP is to provide best-effort services for the delivery of information to its destination. IP depends on upper-level TCP/IP suite layers to provide accountability and reliability. The heart of IP is the IP *datagram*, a packet sent over the Internet in a connectionless manner. An IP datagram carries enough information about the network to get forwarded to its destination; it consists of a *header* followed by bytes of *data*. The header contains information about the type of IP datagram, how long the datagram should stay on the network (or how many hops it should be forwarded to), special flags indicating any special purpose the datagram is supposed to serve, the destination and source addresses, and several other fields, as shown in Figure 1.

Figure 1: The IP Header



Layers above IP use the source address in an incoming packet to identify the sender. To communicate with the sender, the receiving station sends a reply by using the source address in the datagram. Because IP makes no effort to validate whether the source address in the packet generated by a node is actually the source address of the node, you can spoof the source address and the receiver will think the packet is coming from that spoofed address. Many programs for preparing spoofed IP datagrams are available for free on the Internet; for example, *hping* lets you prepare spoofed IP datagrams with just a one-line command, and you can send them to almost anybody in the world. You can spoof at various network layers; for example, you can use *Address Resolution Protocol* (ARP) spoofing to divert the traffic intended for one station to someone else. The *Simple Mail Transfer Protocol* (SMTP) is also a target for spoofing; because SMTP does not verify the sender's address, you can send any e-mail to anybody pretending to be someone else. This article focuses on the various types of attacks that involve IP spoofing on networks, and the techniques and approaches that experts in the field suggest to contend with this problem.

Spoofing IP datagrams is a well-known problem that has been addressed in various research papers. Most spoofing is done for illegitimate purposes—attackers usually want to hide their own identity and somehow damage the IP packet destination. This article discusses ways of spoofing IP datagrams, various attacks that involve spoofed IP packets, and techniques to detect spoofed packets and trace them back to their original source; spoofing concerns for IPv6 are briefly addressed.

Spoofing an IP Datagram

IP packets are used in applications that use the Internet as their communications medium. Usually they are generated automatically for the user, behind the scenes; the user just sees the information exchange in the application. These IP packets have the proper source and destination addresses for reliable exchange of data between two applications. The IP stack in the operating system takes care of the header for the IP datagram. However, you can override this function by inserting a custom header and informing the operating system that the packet does not need any headers. You can use raw sockets in UNIX-like systems to send spoofed IP datagrams, and you can use packet drivers such as *WinPcap* on Windows. Some socket programming knowledge is enough to write a program for generating crafted IP packets. You can insert any kind of header, so, for example, you can also create *Transmission Control Protocol* (TCP) headers. If you do not want to program or have no knowledge of programming, you can use tools such as *hping*, *sendip*, and others that are available for free on the Internet, with very detailed documentation to craft any kind of packet. Most of the time, you can send a spoofed address IP packet with just a one-line command.

Why Spoof the IP Source Address?

What is the advantage of sending a spoofed packet? It is that the sender has some kind of malicious intention and does not want to be identified. You can use the source address in the header of an IP datagram to trace the sender's location. Most systems keep logs of Internet activity, so if attackers want to hide their identity, they need to change the source address. The host receiving the spoofed packet responds to the spoofed address, so the attacker receives no reply back from the victim host. But if the spoofed address belongs to a host on the same subnet as the attacker, then the attacker can “sniff” the reply. You can use IP spoofing for several purposes; for some scenarios an attacker might want to inspect the response from the target victim (called “nonblind spoofing”), whereas in other cases the attacker might not care (blind spoofing). Following is a discussion about reasons to spoof an IP packet.

Scanning

An attacker generally wants to connect to a host to gather information about open ports, operating systems, or applications on the host. The replies from the victim host can help the attacker in gathering information about the system.

These replies might indicate open ports, the operating system, or several applications running on open ports. For example, a response for connection at port 80 indicates the host might be running a Web server. The hacker can then try to *telnet* to this port to see the banner and determine the Web server version and type, and then try to exploit any vulnerability associated with that Web server. In the scanning case, attackers want to examine the replies coming back from the host, so they need to see the returned packet. If the spoofed address is actually an address of a host on the attacker's subnet, then the attacker can use a sniffer to see the packets.

Sequence-Number Prediction

If you establish the connection between two hosts by using TCP, the packets exchanged between the two parties carry sequence numbers for data and acknowledgments. The protocol uses these numbers to determine out-of-order and lost packets, thus ensuring the reliable delivery to the application layer as promised by TCP. These numbers are generated pseudo-randomly in a manner known to both the parties. An attacker might send several spoofed packets to a victim to determine the algorithm generating the sequence numbers and then use that knowledge to intercept an existing session. Again it is important for the attacker to be able to see the replies.

Hijacking an Authorized Session

An attacker who can generate correct sequence numbers can send a reset message to one party in a session informing that party that the session has ended. After taking one of the parties offline, the attacker can use the IP address of that party to connect to the party still online and perform a malicious act on it. The attacker can thus use a trusted communication link to exploit any system vulnerability. Keep in mind that the party that is still online will send the replies back to the legitimate host, which can send a reset to it indicating the invalid session, but by that time the attacker might have already performed the intended actions. Such actions can range from sniffing a packet to presenting a shell from the online host to the attacker's machine.

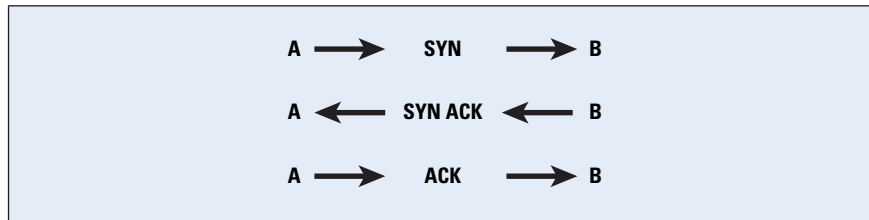
Determining the State of a Firewall

A firewall is used to protect a network from Internet intruders. Packets entering a firewall are checked against an *Access Control List* (ACL). TCP packets sent by a source are acknowledged by acknowledgment packets. If a packet seems like an acknowledgement to a request or data from the local network, then a stateful firewall also checks whether a request for which this packet is carrying the acknowledgment was sent from the network. If there is no such request, the packet is dropped, but a stateless firewall lets packets enter the network if they seem to carry an acknowledgment for a packet. Most probably the intended receiver sends some kind of response back to the spoofed address. Again, for this process to work, the attacker should be able to see the traffic returning to the host that has the spoofed address—and the attacker generally knows how to use the returned packet to advantage.

Denial of Service

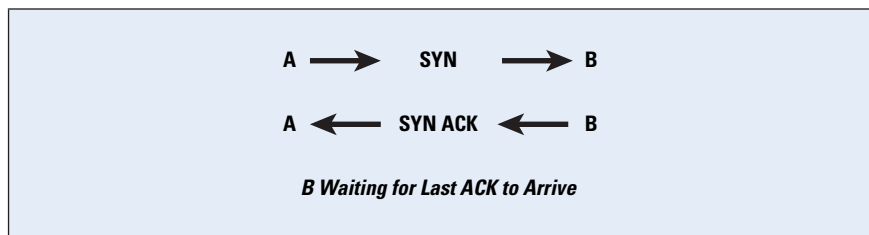
The connection setup phase in a TCP system consists of a *three-way handshake*. This handshake is done by using special bit combinations in the “flags” fields. If host A wants to establish a TCP connection with host B, it sends a packet with a SYN flag set. Host B replies with a packet that has SYN and ACK flags set in the TCP header. Host A sends back a packet with an ACK flag set, finishing the initial handshake. Then hosts A and B can communicate with each other, as shown in Figure 2.

Figure 2: A Normal TCP Connection
Request from A to B



The three-way handshake must be completed in order to establish a connection. Connections that have been initiated but not finished are called *half-open connections*. A finite-size data structure is used to store the state of the half-open connections. An attacking host can send an initial SYN packet with a spoofed IP address, and then the victim sends the SYN-ACK packet and waits for a final ACK to complete the handshake. If the spoofed address does not belong to a host, then this connection stays in the half-open state indefinitely, thus occupying the data structure. If there are enough half-open connections to fill the state data structure, then the host cannot accept further requests, thus denying service to the legitimate connections (Figure 3).

Figure 3: Half-Open TCP Connection



Setting a time limit for half-open connections and then erasing them after the timeout can help with this problem, but the attacker may keep continuously sending the packets. The attacked host will not have space to accept new incoming legitimate connections, but the connection that was established before the attack will have no effect. In this type of attack, the attacker has no interest in examining the responses from the victim. When the spoofed address does belong to a connected host, that host sends a reset to indicate the end of the handshake.

Flooding

In this type of attack an attacker sends a packet with the source address of the victim to multiple hosts. Responses from other machines flood the victim. For example, if an attacker uses the IP address of source A and sends a broadcast message to all the hosts in the network, then all of them will send a reply back to A, hence flooding it. The well-known *Smurf* and *fraggle* attacks used this technique.

Countermeasures for IP Spoofing

IP spoofing countermeasures include detecting spoofed IP packets and then tracing them back to the originating source. Detection of spoofed IP packets requires support of routers, host-based methods, and administrative controls, whereas tracing of IP packets involves special traceback equipment or traceback features in routers. The following section discusses both IP spoofing detection and IP spoofing traceback techniques.

Spoofed Packet Detection

Detection of a spoofed packet can start as early as at Layer 2. Switches with the *IP Source Guard* feature^[8] match the MAC address of the host with a *Dynamic Host Configuration Protocol* (DHCP)-assigned dynamic or administratively assigned static IP address. Packets that do not have the correct IP source address for that particular MAC address are dropped, thereby limiting the ability of hosts connected to such a switch to send a packet with their neighbor's address. The IP Source Guard feature works very well for interfaces with a single IP address, but one interface can be assigned multiple IP addresses, and that may cause problems. The same problems can occur with *Network Address Translation* (NAT), where hosts might get different IP addresses several times. Routers work at Layer 3 in networks, and they know which interface a network is connected to and what network addresses can be expected to come from that network. If the outgoing packet from an interface does not have the network address of that interface, then the packet is spoofed and the router can stop that packet at that point; however, if the attacker is spoofing an IP address of a host on the same network (most likely in the attacks where they will be sniffing the replies), then this technique is not really helpful. The same logic can be used for an incoming packet; if a packet destined for an interface has a source address of the same network as the interface, then it is a spoofed packet. Routers can detect spoofed packets only when the packets pass through them, and if the target and attacker are both on the same subnet then this technique does not work.

Hosts receiving a suspicious packet can also use certain techniques to determine whether or not the IP address is spoofed. The first (and easiest) one is to send a request to the address of the packet and wait for the response; most of the time the spoofed addressees do not belong to active hosts and hence no response is sent.

Another method is to check the *Time to Live* (TTL) value of the packet, and then send a request to the spoofed host. If the reply comes, you can compare the TTL of both packets. Most probably the TTL values will not match. But of course it is also possible that these TTL values are the same but the packet is coming from a different source, and conversely. Packets generated by different operating systems differ slightly in values of certain fields; for example, in *Internet Control Message Protocol* (ICMP) *ping* packets, you can examine the data payload to determine the operating system. Windows fills the packet with letters of the alphabet, whereas Linux puts numbers in the data portion. If the suspicious packet does not have the same characteristics as the legitimate packet, that is evidence it was not sent from the IP address that is in its source address field. You can also use IP identification numbers to determine whether a packet is actually coming from the said source. For legitimate packets the IP ID is close in value, but this method is not reliable because the attacker can ping the said source and determine the IP ID that it is using, and then craft packets that will seem legitimate. In all these techniques we are trying to determine only whether or not a packet is spoofed, and taking all these steps for all packets would be prohibitive from an overhead standpoint. Thus you should either randomly check packets or determine some suspicious activity that would trigger further investigation for spoofed-packet detection. The next section addresses measures you can take to trace a spoofed packet back to its real source.

Tracing Spoofed IP Packets

IP traceback technology plays an important role in discovering the source of spoofed packets. Hop-by-hop traceback and logging of suspicious packets in routers are the two main methods for tracing the spoofed IP packets back to their source.

When a node detects that it is a victim of flood attack, it can inform the *Internet Service Provider* (ISP). In flood attacks the ISP can determine the router that is sending this stream to the victim, and then it can determine the next router, and so on. It reaches either to the source of the flood attack or the end of its administrative domain; for this case it can ask the ISP for the next domain to do the same thing. This technique is useful only if the flood is ongoing.

As mentioned earlier, a router has an idea of the IP addresses that should be arriving at its interfaces. If it sees any packet that does not seem to belong to the address range for its interface, it can log the packet as suspicious. Appropriately timed broadcasts among different domains to detect spoofed packets can help administrators of different networks trace spoofed IP packets back to their source.

IP Spoofing and IPv6

IP spoofing detection, or in other words validating the source address of an IPv6 packet, is a little more complicated than the process for IPv4. A host using IPv6 may potentially have multiple addresses. Again the problem inside the Local Area Network is to associate the IPv6 address with the Layer 2 or MAC address. Among peers on the same network, you can use *Neighbor Discovery* or *Secure Neighbor Discovery* (SEND) advertisements to verify the source address in a packet. You can verify source addresses of packets arriving from nodes outside the network by using the *Authentication Header* (AH) in IPv6 datagrams. You can use agreed-upon parameters between source and destination to calculate authentication information on header fields that does not change during transit. Although this process will not prevent someone from signing a spoofed address, it does provide a means to authenticate the identity of the source.

IPv6 and IPv4 network interconnections will likely face spoofing problems. IPv6 packets are usually encapsulated in IPv4 packets to travel across the non-IPv6 supporting networks. The IPv6 interim mechanism “6to4”^[10, 11] uses automatic IPv6-to-IPv4 tunneling to interconnect networks using different IP versions. This mechanism uses 6to4 routers and 6to4 *Relay Routers* that accept and decapsulate IPv4 traffic from anywhere. There are no constraints on such embedded packets. Relay routers act as bridges between IPv6 and 6to4 networks and can be tricked into sending spoofed traffic anywhere. Also, anyone can send tunneled spoofed traffic to a 6to4 router, and the router will believe that it is coming from a legitimate relay. There is no simple way to prevent such attacks, and longer-term solutions are needed in both IPv6 and IPv4 networks.

Conclusion

IP spoofing is a difficult problem to tackle, because it is related to the IP packet structure. IP packets can be exploited in several ways. Because attackers can hide their identity with IP spoofing, they can make several network attacks. Although there is no easy solution for the IP spoofing problem, you *can* apply some simple proactive and reactive methods at the nodes, and use the routers in the network to help detect a spoofed packet and trace it back to its originating source.

References

- [1] Alaaeldin A. Aly, “Tracking and Tracing Spoofed IP Packets to Their Sources,” Proceedings of 6th annual conference, UAEU April 2005.
- [2] S.J. Templeton and K.E. Levitt, “Detecting Spoofed Packets,” DARPA Information Survivability Conference and Exposition, 2003.
- [3] “IP Spoofing an Introduction,”
<http://www.securityfocus.com/infocus/1674>
- [4] <http://www.phrack.org/issues.html?issue=48&id=14#article>
- [5] <http://www.hping.org>
- [6] <http://www.insecure.org/nmap>
- [7] <http://www.ietf.org/internet-drafts/draft-baker-sava-operational-00.txt>
- [8] <http://tools.ietf.org/html/draft-baker-sava-cisco-ip-source-guard-00>
- [9] <http://tools.ietf.org/id/draft-baker-sava-implementation-00.txt>
- [10] <http://tools.ietf.org/html/draft-ietf-v6ops-6to4-security-04>
- [11] Carpenter, B., Fink, B., and Moore, K., “Connecting IPv6 Routing Domains Over the IPv4 Internet,” *The Internet Protocol Journal*, Volume 3, No. 1, March 2000.
- [12] Wesley Eddy, “Defenses Against TCP SYN Flooding Attacks,” *The Internet Protocol Journal*, Volume 9, No. 4, December 2006.

FARHA ALI holds a BE in Computer Engineering from NED University, Pakistan, and an MS in Computer Engineering from Clemson University, South Carolina, with a focus area in Computer Communications. She is a member of American Mensa and ACM. Her research papers (co-authored with her advisor) as a PhD student at Clemson University’s Computer Science Department were published in IEEE’s Conferences on Web Intelligence and Web Services. She is a Sun Certified Java Programmer and a Certified Ethical Hacker. Her main interests are Distributed Computing, Network Security, and Semantic Web. Currently she is working as a faculty member at Lander University’s Department of Mathematics and Computing. She teaches mainly Networking and Programming courses.
E-mail: fali@lander.edu

Standards for Information Security Management

by William Stallings

To effectively assess the security needs of an organization and to evaluate and choose various security products and policies, the manager responsible for security needs some systematic way of defining the requirements for security and characterizing the approaches to satisfy those requirements. This process is difficult enough in a centralized data processing environment; with the use of local- and wide-area networks (LANs and WANs, respectively), the problems are compounded.

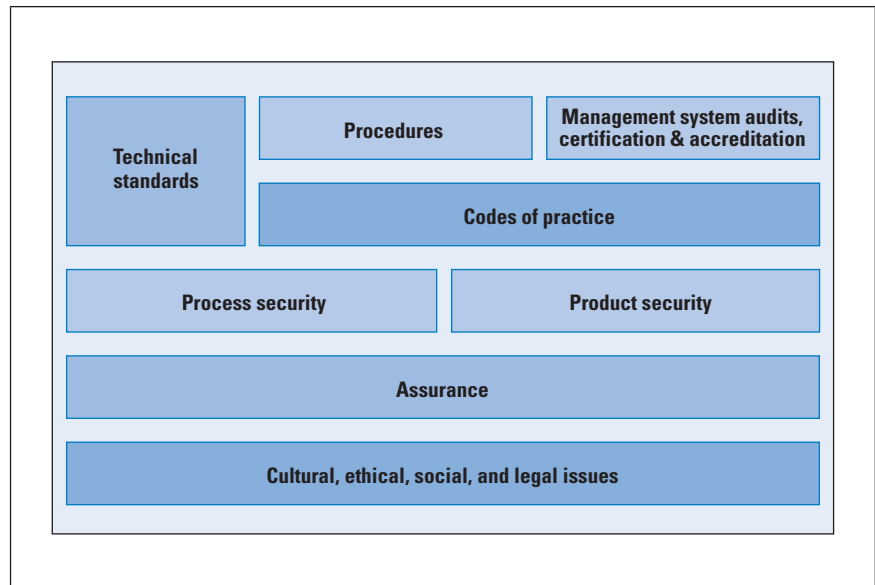
The challenges for management in providing information security are formidable. Even for relatively small organizations, information system assets are substantial, including databases and files related to personnel, company operation, financial matters, and so on. Typically, the information system environment is complex, including a variety of storage systems, servers, workstations, local networks, and Internet and other remote network connections. Managers face a range of threats always growing in sophistication and scope. And the range of consequences for security failures, both to the company and to individual managers, is substantial, including financial loss, civil liability, and even criminal liability.

Standards for providing information system security become essential in such circumstances. Standards can define the scope of security functions and features needed, policies for managing information and human assets, criteria for evaluating the effectiveness of security measures, techniques for ongoing assessment of security and for the ongoing monitoring of security breaches, and procedures for dealing with security failures.

Figure 1, based on [1], suggests the elements that, in an integrated fashion, constitute an effective approach to information security management. The focus of this approach is on two distinct aspects of providing information security: process and products. *Process security* looks at information security from the point of view of management policies, procedures, and controls. *Product security* focuses on technical aspects and is concerned with the use of certified products in the IT environment when possible. In Figure 1, the term *technical standards* refers to specifications that refer to aspects such as IT network security, digital signatures, access control, nonrepudiation, key management, and hash functions. Operational, management, and technical *procedures* encompass policies and practices that are defined and enforced by management. Examples include personnel screening policies, guidelines for classifying information, and procedures for assigning user IDs. *Management system audits, certification, and accreditation* deals with management policies and procedures for auditing and certifying information security products.

Codes of practice refer to specific policy standards that define the roles and responsibilities of various employees in maintaining information security. *Assurance* deals with product and system testing and evaluation. *Cultural, ethical, social, and legal issues* refer to human factors aspects related to information security.

Figure 1: Information Security Management Elements



Many standards and guideline documents have been developed in recent years to aid management in the area of information security. The two most important are *ISO 17799*, which deals primarily with process security, and the *Common Criteria*, which deals primarily with product security. This article surveys these two standards, and examines some other important standards and guidelines as well.

ISO 17799

An increasingly popular standard for writing and implementing security policies is *ISO 17799* “Code of Practice for Information Security Management.” (*ISO 17799* will eventually be reissued as *ISO 27002* in the new *ISO 27000* family of security standards). *ISO 17799* is a comprehensive set of controls comprising best practices in information security. It is essentially an internationally recognized generic information security standard. Table 1 summarizes the area covered by this standard and indicates the objectives for each area.

Table 1: ISO 17799 Areas and Objectives

<p>Security Policy Provide management direction and support for information security in accordance with business requirements and relevant laws and regulations.</p> <p>Organization of Information Security Manage information security within the organization. Maintain the security of the organization's information and information processing facilities that are accessed, processed, communicated to, or managed by external parties.</p> <p>Asset Management Achieve and maintain appropriate protection of organizational assets. Ensure that information receives an appropriate level of protection.</p> <p>Human Resources Security Ensure that employees, contractors, and third-party users (1) understand their responsibilities and are suitable for the roles they are considered for; (2) are aware of information security threats and concerns; (3) exit an organization or change employment in an orderly manner.</p> <p>Physical and Environmental Security Prevent unauthorized physical access, damage, and interference to the organization's premises and information. Prevent loss, damage, theft, or compromise of assets and interruption to the organization's activities.</p> <p>Communications and Operations Management Develop controls for operational procedures, third-party service delivery management, system planning, malware protection, backup, network security management, media handling, information exchange, e-commerce services, and monitoring.</p>	<p>Access Control Develop controls for business requirements for user access, user responsibilities, network access control, OS access control, application access control, and information access control.</p> <p>Information Systems Acquisition, Development, and Maintenance Develop controls for correct processing in applications, cryptographic functions, system file security, support process security, and vulnerability management.</p> <p>Information Security Incident Management Ensure information security events and weaknesses associated with information systems are communicated in a manner that allows timely corrective action to be taken. Ensure a consistent and effective approach is applied to the management of information security incidents.</p> <p>Business Continuity Management Counteract interruptions to business activities to protect critical business processes from the effects of major failures of information systems or disasters and to ensure their timely resumption.</p> <p>Compliance Avoid breaches of any law, statutory, regulatory, or contractual obligations, and of any security requirements. Ensure compliance of systems with organizational security policies and standards. Maximize the effectiveness of and minimize interference to and from the information systems audit process.</p>
--	---

With the increasing interest in security, ISO 17799 certification, provided by various accredited bodies, has been established as a goal for many corporations, government agencies, and other organizations around the world. ISO 17799 offers a convenient framework to help security policy writers structure their policies in accordance with an international standard.

Much of the content of ISO 17799 deals with security controls, which are defined as practices, procedures, or mechanisms that may protect against a threat, reduce a vulnerability, limit the effect of an unwanted incident, detect unwanted incidents, and facilitate recovery. Some controls deal with security management, focusing on management actions to institute and maintain security policies. Other controls are operational; they address the correct implementation and use of security policies and standards, ensuring consistency in security operations and correcting identified operational deficiencies. These controls relate to mechanisms and procedures that are primarily implemented by people rather than systems.

Finally, there are technical controls; they involve the correct use of hardware and software security capabilities in systems. These controls range from simple to complex measures that work together to secure critical and sensitive data, information, and IT systems functions. This concept of controls cuts across all the areas listed in Table 1.

To give some idea of the scope of ISO 17799, we examine several of the security areas discussed in that document. *Auditing* is a key security management function that is addressed in multiple areas within the document. First, ISO 17799 lists key data items that should, when relevant, be included in an audit log:

- User IDs
- Dates, times, and details of key events, for example, log-on and log-off
- Terminal identity or location if possible
- Records of successful and rejected system access attempts
- Records of successful and rejected data and other resource access attempts
- Changes to system configuration
- Use of privileges
- Use of system utilities and applications
- Files accessed and the kind of access
- Network addresses and protocols
- Alarms raised by the access control system
- Activation and deactivation of protection systems, such as antivirus systems and intrusion detection systems

It provides a useful set of guidelines for implementation of an auditing capability:

1. Audit requirements should be agreed upon by appropriate management.
2. The scope of the checks should be agreed upon and controlled.
3. The checks should be limited to read-only access to software and data.
4. Access other than read-only should be allowed only for isolated copies of system files, which should be erased when the audit is completed or given appropriate protection if there is an obligation to keep such files under audit documentation requirements.
5. Resources for performing the checks should be explicitly identified and made available.
6. Requirements for special or additional processing should be identified and agreed upon.

7. All access should be monitored and logged to produce a reference trail; the use of timestamped reference trails should be considered for critical data or systems.
8. All procedures, requirements, and responsibilities should be documented.
9. The person(s) carrying out the audit should be independent of the activities audited.

Under the area of communications and operations management, ISO 17799 includes *network security management*. One aspect of this management is concerned with network controls for networks owned and operated by the organization. The document provides implementation guidance for these in-house networks. An example of a control follows: Restoration procedures should be regularly checked and tested to ensure that they are effective and that they can be completed within the time allotted in the operational procedures for recovery. Similarly, the document provides guidance for security controls for network services provided by outside vendors. An example of guidance in this area follows: The ability of the network service provider to manage agreed-upon services in a secure way should be determined and regularly monitored, and the right to audit should be agreed upon.

As can be seen, some ISO 17700 specifications are detailed and specific, whereas others are quite general.

Common Criteria

The *Common Criteria for Information Technology Security Evaluation* (CC) is a joint international effort by numerous national standards organizations and government agencies^[3,4,5]. U.S. participation is by the *National Institute of Standards and Technology* (NIST) and the *National Security Agency* (NSA). CC defines a set of IT requirements of known validity that can be used in establishing security requirements for prospective products and systems. The CC also defines the *Protection Profile* (PP) construct that allows prospective consumers or developers to create standardized sets of security requirements that will meet their needs.

The aim of the CC specification is to provide greater confidence in the security of IT products as a result of formal actions taken during the process of developing, evaluating, and operating these products. In the development stage, the CC defines sets of IT requirements of known validity that can be used to establish the security requirements of prospective products and systems. Then the CC details how a specific product can be evaluated against these known requirements, to provide confirmation that it does indeed meet them, with an appropriate level of confidence. Lastly, when in operation the evolving IT environment may reveal new vulnerabilities or concerns. The CC details a process for responding to such changes, and possibly reevaluating the product.

Following successful evaluation, a particular product may be listed as CC certified or validated by the appropriate national agency, such as NIST or NSA in the United States. That agency publishes lists of evaluated products, which are used by government and industry purchasers who need to use such products.

The CC defines a common set of potential *security requirements* for use in evaluation. The term *Target of Evaluation* (TOE) refers to that part of the product or system that is subject to evaluation. The requirements fall into two categories:

- *Functional requirements*: Define desired security behavior. CC documents establish a set of security functional components that provide a standard way of expressing the security functional requirements for a TOE.
- *Assurance requirements*: The basis for gaining confidence that the claimed security measures are effective and implemented correctly. CC documents establish a set of assurance components that provide a standard way of expressing the assurance requirements for a TOE.

Both functional requirements and assurance requirements are organized into classes: A *class* is a collection of requirements that share a common focus or intent. Each of these classes contains numerous families. The requirements within each *family* share security objectives but differ in emphasis or rigor. For example, the audit class contains six families dealing with various aspects of auditing (for example, audit data generation, audit analysis, and audit event storage). Each family, in turn, contains one or more components. A *component* describes a specific set of security requirements and is the smallest selectable set of security requirements for inclusion in the structures defined in the CC.

For example, the cryptographic support class of functional requirements includes two families: cryptographic key management and cryptographic operation. The cryptographic key management family has four components, which are used to specify key generation algorithm and key size; key distribution method; key access method; and key destruction method. For each component, a standard may be referenced to define the requirement. Under the cryptographic operation family, there is a single component, which specifies an algorithm and key size based on an assigned standard.

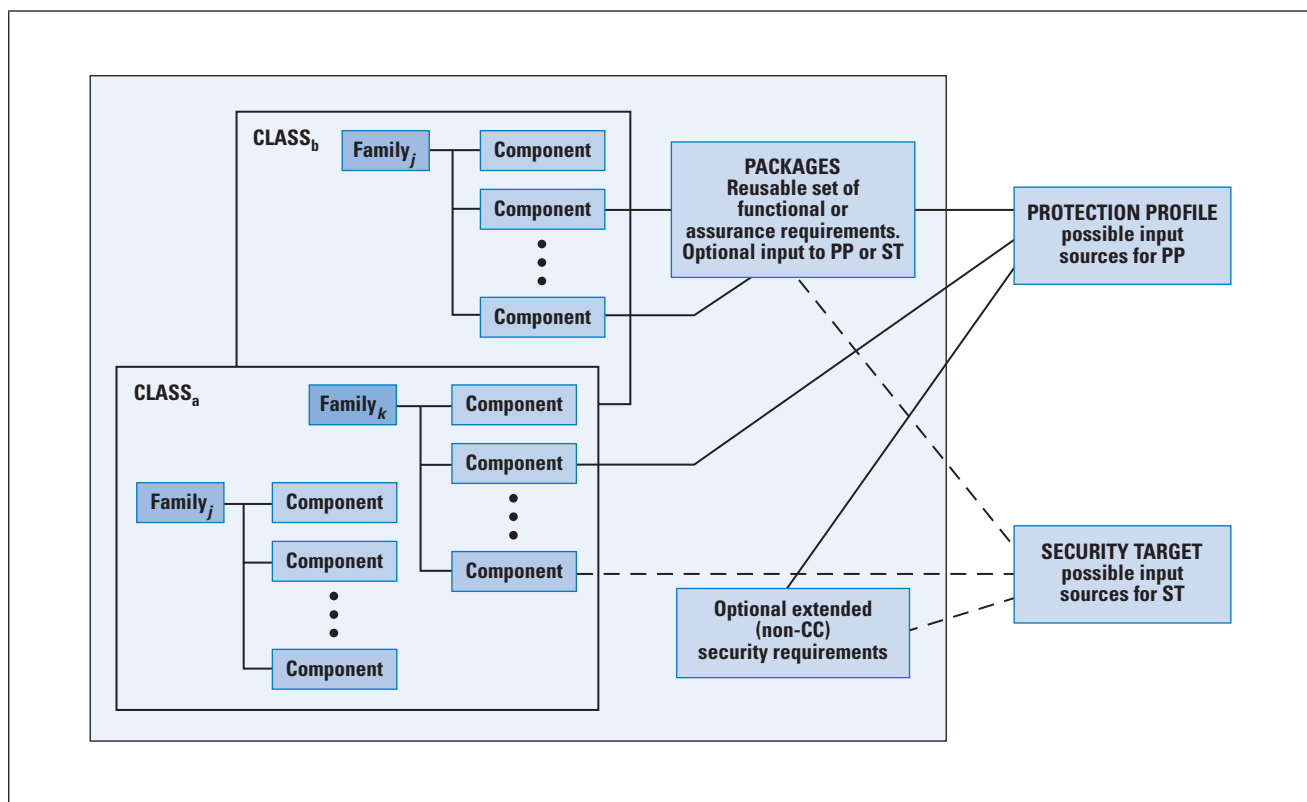
Sets of functional and assurance components may be grouped together into reusable packages, which are known to be useful in meeting identified objectives. An example of such a package would be functional components required for *Discretionary Access Controls*.

The CC also defines two kinds of documents that can be generated using the CC-defined requirements.

- *Protection profiles (PPs)*: Define an implementation-independent set of security requirements and objectives for a category of products or systems that meet similar consumer needs for IT security. A PP is intended to be reusable and to define requirements that are known to be useful and effective in meeting the identified objectives. The PP concept has been developed to support the definition of functional standards and as an aid to formulating procurement specifications. The PP reflects user security requirements.
- *Security targets (STs)*: Contain the IT security objectives and requirements of a specific identified TOE and define the functional and assurance measures offered by that TOE to meet stated requirements. The ST may claim conformance to one or more PPs and forms the basis for an evaluation. The ST is supplied by a vendor or developer.

Figure 2 illustrates the relationship between requirements on the one hand and profiles and targets on the other. For a PP, a user can select many components to define the requirements for the desired product. The user may also refer to predefined packages that assemble numerous requirements commonly grouped together within a product requirements document. Similarly, a vendor or designer can select numerous components and packages to define an ST.

Figure 2: Organization and Construction of Common Criteria Requirements



As an example for the use of the CC, consider the smart card. The protection profile for a smart card, developed by the *Smart Card Security User Group*, provides a simple example of a PP. This PP describes the IT security requirements for a smart card to be used in connection with sensitive applications, such as banking industry financial payment systems. The assurance level for this PP is *Evaluation Assurance Level* (EAL) 4, which is described subsequently. The PP lists *threats* that must be addressed by a product that claims to comply with this PP. The threats include the following:

- *Physical probing*: May entail reading data from the TOE through techniques commonly employed in *Integrated Circuit* (IC) failure analysis and IC reverse engineering efforts.
- *Invalid input*: Invalid input may take the form of operations that are not formatted correctly, requests for information beyond register limits, or attempts to find and execute undocumented commands. The result of such an attack may be a compromise in the security functions, generation of exploitable errors in operation, or release of protected data.
- *Linkage of multiple operations*: An attacker may observe multiple uses of resources or services and, by linking these observations, deduce information that may reveal security function data.

Following a list of threats, the PP turns to a description of *security objectives*, which reflect the stated intent to counter identified threats or comply with any organizational security policies identified. Nineteen objectives are listed, including the following:

- *Audit*: The system must provide the means of recording selected security-relevant events, so as to assist an administrator in the detection of potential attacks or misconfiguration of the system security features that would leave it susceptible to attack.
- *Fault insertion*: The system must be resistant to repeated probing through insertion of erroneous data.
- *Information leakage*: The system must provide the means of controlling and limiting the leakage of information in the system so that no useful information is revealed over the power, ground, clock, reset, or I/O lines.

Security requirements are provided to thwart specific threats and to support specific policies under specific assumptions. The PP lists specific requirements in three general areas: TOE security functional requirements, TOE security assurance requirements, and security requirements for the IT environment. In the area of *security functional requirements*, the PP defines 42 requirements from the available classes of security functional requirements.

For example, for security auditing, the PP stipulates what the system must audit; what information must be logged; what the rules are for monitoring, operating, and protecting the logs; and so on. Functional requirements are also listed from the other functional requirements classes, with specific details for the smart card operation.

The PP defines 24 *security assurance requirements* from the available classes of security assurance requirements. These requirements were chosen to demonstrate:

- The quality of the product design and configuration
- That adequate protection is provided during the design and implementation of the product
- That vendor testing of the product meets specific parameters
- That security functions are not compromised during product delivery
- That user guidance, including product manuals pertaining to installation, maintenance, and use, are of a specified quality and appropriateness

The PP also lists *Security Requirements of the IT Environment*. They cover the following topics:

- Cryptographic key distribution
- Cryptographic key destruction
- Security roles

The final section of the PP (excluding appendices) is a lengthy rationale for all the selections and definitions in the PP. The PP is an industrywide effort designed to be realistic in its ability to be met by a variety of products with a variety of internal mechanisms and implementation approaches.

The concept of *Evaluation Assurance* is a difficult one to define. Further, the degree of assurance required varies from one context and one function to another. To structure the need for assurance, the CC defines a scale for rating assurance consisting of seven EALs ranging from the least rigor and scope for assurance evidence (EAL 1) to the most (EAL 7). The levels are as follows:

- *EAL 1: Functionally tested:* For environments where security threats are not considered serious. It involves independent product testing with no input from the product developers. The intent is to provide a level of confidence in correct operation.
- *EAL 2: Structurally tested:* Includes a review of a high-level design provided by the product developer. Also, the developer must conduct a vulnerability analysis for well-known flaws. The intent is to provide a low to moderate level of independently assured security.

- *EAL 3: Methodically tested and checked:* Requires a focus on the security features, including requirements that the design separate security-related components from those that are not; that the design specifies how security is enforced; and that testing be based both on the interface and the high-level design, rather than a black box testing based only on the interface. It is applicable where the requirement is for a moderate level of independently assured security, with a thorough investigation of the TOE and its development without incurring substantial reengineering costs.
- *EAL 4: Methodically designed, tested, and reviewed:* Requires both a low-level and a high-level design specification; requires that the interface specification be complete; requires an abstract model that explicitly defines security for the product; and requires an independent vulnerability analysis. It is applicable in those circumstances where developers or users require a moderate to high level of independently assured security in conventional commodity TOEs, and there is willingness to incur some additional security-specific engineering costs.
- *EAL 5: Semiformally designed and tested:* Provides an analysis that includes all of the implementation. Assurance is supplemented by a formal model and a semiformal presentation of the functional specification and high-level design and a semiformal demonstration of correspondence. The search for vulnerabilities must ensure resistance to penetration attackers with a moderate attack potential. Covert channel analysis and modular design are also required.
- *EAL 6: Semiformally verified design and tested:* Permits a developer to gain high assurance from application of specialized security engineering techniques in a rigorous development environment, and to produce a premium TOE for protecting high-value assets against significant risks. The independent search for vulnerabilities must ensure resistance to penetration attackers with a high attack potential.
- *EAL 7: Formally verified design and tested:* The formal model is supplemented by a formal presentation of the functional specification and high-level design, showing correspondence. Evidence of developer “white box” testing of internals and complete independent confirmation of developer test results are required. Complexity of the design must be minimized.

The first four levels reflect various levels of commercial design practice. Only at the highest of these levels (EAL 4) is there a requirement for any source code analysis, and this analysis is required only for a portion of the code. The top three levels provide specific guidance for products developed using security specialists and security-specific design and engineering approaches.

National Institute of Standards and Technology

NIST has produced a large number of *Federal Information Processing Standards Publications* (FIPS PUBs) and special publications (SPs) that are enormously useful to security managers, designers, and implementers. Following are a few of the most significant and general. *FIPS PUB 200* “Minimum Security Requirements for Federal Information and Information Systems,” is a standard that specifies minimum security requirements in 17 security-related areas with regard to protecting the confidentiality, integrity, and availability of federal information systems and the information processed, stored, and transmitted by those systems^[6].

NIST *SP 800-100* “Information Security Handbook: A Guide for Managers,” provides a broad overview of information security program elements to assist managers in understanding how to establish and implement an information security program^[7]. Its topical coverage overlaps considerably with ISO 17799.

Several other NIST publications are of general interest. *SP 800-55* “Security Metrics Guide for Information Technology Systems,” provides guidance on how an organization, through the use of metrics, identifies the adequacy of in-place security controls, policies, and procedures^[8]. *SP 800-27* “Engineering Principles for Information Technology Security (A Baseline for Achieving Security),” presents a list of system-level security principles to be considered in the design, development, and operation of an information system^[9]. *SP 800-53* “Recommended Security Controls for Federal Information Systems,” lists management, operational, and technical safeguards or countermeasures prescribed for an information system to protect the confidentiality, integrity, and availability of the system and its information^[10].

Other Standards and Guidelines

Another important set of standards is the *Control Objectives for Information and Related Technology* (COBIT)^[11], a business-oriented set of standards for guiding management in the sound use of information technology. It has been developed as a general standard for information technology security and control practices and includes a general framework for management, users, IS audit, and security practitioners. COBIT also has a process focus and a governance flavor; that is, management’s need to control and measure IT is a focus point. COBIT was developed under the auspices of a professional organization, the *Information Systems Audit and Control Association* (ISACA). The documents are quite detailed and provide a practical basis for not only defining security requirements but also implementing them and verifying compliance.

Another excellent source of information is “The Standard of Good Practice for Information Security” from the *Information Security Forum*. The standard is designed as an aid to organizations in understanding and applying best practices for information security. Because it addresses security from a business perspective, The Standard appropriately recognizes the intersection between organizational factors and security factors.

In addition to these standards, numerous informal guidelines are widely consulted by organizations in developing their own security policy. The *CERT Coordination Center* (www.cert.org) has an Evaluations and Practices section of its Website with a variety of documents and training aids related to information security for organizations. The *Chief Information Officers Council* (cio.gov) has published a collection of Best Practices and other documents related to organizational security.

References

- [1] Eloff, J., and Eloff, M., “Information Security Management,” *Proceedings of SAICSIT 2003*, South African Institute of Computer Scientists and Information Technologists, 2003.
- [2] International Organization for Standardization, “ISO/IEC 27001 – Information technology – Security Techniques – Information security management systems – Requirements,” June 2005.
- [3] Common Criteria Project Sponsoring Organisations, “Common Criteria for Information Technology Security Evaluation, Part 1: Introduction and General Model,” CCIMB-2004-01-001, January 2004.
- [4] Common Criteria Project Sponsoring Organisations, “Common Criteria for Information Technology Security Evaluation, Part 2: Security Functional Requirements,” CCIMB-2004-01-002, January 2004.
- [5] Common Criteria Project Sponsoring Organisations, “Common Criteria for Information Technology Security Evaluation, Part 3: Security Assurance Components,” CCIMB-2006-09-003, September 2006.
- [6] National Institute of Standards and Technology, “Minimum Security Requirements for Federal Information and Information Systems,” FIPS PUB 200, March 2006.
- [7] National Institute of Standards and Technology, “Information Security Handbook: A Guide for Managers,” NIST Special Publication 800-100, October 2006.

- [8] “Security Metrics Guide for Information Technology Systems,” NIST Special Publication 800-55, July 2003.
- [9] National Institute of Standards and Technology, “Engineering Principles for Information Technology Security (A Baseline for Achieving Security),” NIST Special Publication 800-27, June 2004.
- [10] National Institute of Standards and Technology, “Recommended Security Controls for Federal Information Systems,” NIST Special Publication 800-53, February 2005.
- [11] IT Governance Institute, “COBIT 4.0.,” USA, 2005.
- [12] Information Security Forum, “The Standard of Good Practice for Information Security,” 2005.

WILLIAM STALLINGS is a consultant, lecturer, and author of more than a dozen books on data communications and computer networking. His latest book, with Lawrie Brown, is *Computer Security: Principles and Practice* (Prentice Hall, 2007). He maintains a computer science resource site for computer science students and professionals at WilliamStallings.com/StudentSupport.html and is on the editorial board of *Cryptologia*. He has a Ph.D. in computer science from M.I.T. He can be reached at ws@shore.net.

Looking Toward the Future

by Vint Cerf, Google

The *Internet Corporation for Assigned Names and Numbers* (ICANN) was formed 9 years ago, following a period of considerable debate about the institutionalization of the basic functions performed by the *Internet Assigned Numbers Authority* (IANA)^[1]. Nearly simultaneous with the inauguration of ICANN in September 1998 came the unexpected and untimely death of the man, Jonathan B. Postel^[2], who had responsibility for these functions for more than a quarter century. The organization began with very limited sources of funds, a small and overworked staff, and contentious debate about its organizational structure, policy apparatus, and operational procedures. The organization underwent substantial change through its *Evolution and Reform Process* (ERP)^[3]. Among the more difficult constituencies to accommodate in the organization's policy-making process was the general public. An *At-Large Advisory Committee* (ALAC)^[4] emerged from the ERP and has recently formed *Regional At-Large Organizations* (RALOs) in all of ICANN's five regions.

Today, ICANN is larger, more capable, more international, and better positioned to fulfill its mandate. It stands for one global, interoperable Internet, and the model of stakeholder representation has worked. But the Internet and its vast user population have grown during the same time by a factor of more than 20 in all dimensions. The 50 million users of 1997 have become nearly 1.2 billion users today. The 22 million hosts on the network have increased to nearly 500 million today. The bandwidth of the core data circuits in the Internet have grown from 622 million bits per second to between 10 and 40 billion bits per second. This dramatic growth in physical size has been accompanied by an equally dramatic growth in the number and diversity of applications running on the Internet. All forms of media now appear on and are carried by Internet packets. Consumers of information are producing more and more of it themselves with e-mail, blogs, instant messaging, social and game-playing Websites, video uploads, and podcasts. The Internet continues to evolve and although ICANN has achieved more than most people realize, it must continue to evolve along with it.

Operational Priorities

ICANN's primary responsibility is to contribute to the security and stability of the Internet system of unique identifiers. In the most direct way, it carries out this mandate through its operation of the IANA. There is no doubt that the conduct of this function in an exemplary fashion is essential not only to ICANN's mission but also to inspiring confidence in ICANN as an organization.

But ICANN's role in the Internet goes beyond these specific IANA functions. ICANN is an experiment in the balancing of multiple stakeholder interests in policy about the implementation, operation, and use of the *Domain Name System* (DNS) and the address spaces of the Internet. Its policy choices can directly affect the business models of operating entities involved in the management of domain names and Internet addresses. The privacy and Internet-related rights of registrants and, more generally, Internet users may also be directly affected. Some policy choices raise public policy concerns in the view of governments and methods and will be needed to factor such concerns into the making of ICANN policy.

Effective, fair, and timely policy development should be a priority for ICANN. That this policy development needs to be achieved in a global setting is simply another challenge to be met. ICANN leadership and staff must seek to maintain and improve the ability of all of ICANN's many constituencies to achieve consensus or at least to prepare the ICANN Board to make choices when consensus may not be forthcoming. Because policies often have technical, economic, social, and governance implications, it is vital that ICANN's practices draw on expertise in all these domains.

Clarity in the roles and responsibilities of the many participants in the Internet arena, especially those with specific interest in ICANN policies and practices, will be helpful and should be documented. In some cases, the documentation might take the form of relatively formal relationships such as the contracts between ICANN and domain-name registries and registrars. In other cases, they may need only to characterize in plain terms the roles that each party plays.

In some areas, such as root-zone operation, excellence can be measured in such terms as responsiveness, scalability, resilience to disruption, and ability to adapt to changing needs such as *Domain Name System Security* (DNSSEC)^[5], *Internationalized Domain Names* (IDNs)^[6], and the addition of IPv6 records to the root zone. Many parties currently play a role in the maintenance of the root-zone file, and clear documentation of responsibility and lines of authority will be beneficial. As the technology of the Internet continues to evolve, the roles of various parties may need to change to meet the objectives of stability and security of the Internet system of unique identifiers. Managing the evolution of these roles represents another priority for policy development and implementation.

Because of the potential effect of decisions made through the ICANN policy process, it is important to implement checks and balances that make all aspects of ICANN's operation accountable and transparent. Work is still necessary in this area so that independent review of legitimate concerns arising out of policy making is possible when deemed necessary.

At the same time, it is vital that the mechanisms chosen do not have the effect of locking up the policy-making process and preventing any decisions from being made. We need to seek a balance between a potentially unfair tyranny of the majority and an equally unacceptable tyranny of the minority.

The general success of the *Uniform Dispute Resolution Process* (UDRP)^[7] suggests that ICANN should seek mechanisms for resolving disputes arising in connection with implementing ICANN policies that scale, permit choice without abusive “forum shopping,” and make efficient use of ICANN resources.

Outreach, transparency, and broadly participatory processes on an international basis are not inexpensive. It is vital for ICANN to continue to refine its models for sustainable operation, accounting for the economics of the various actors in the Internet arena that rely on ICANN’s operation, and fairly apportioning costs of ICANN operation to appropriate sources of support. Not all of the beneficiaries of ICANN’s work derive the same level of revenue from the Internet (and some, none at all). ICANN must account for this discrepancy when devising mechanisms for supporting its operation, and it should work to make transparent the need to provide services to parties who may not be able to contribute commensurate with cost. Adequate and stable funding for ICANN is necessary if ICANN is to fulfill its charter. Over the past several years, ICANN has significantly increased its ability to staff vital functions, contributing to the effectiveness of the organization. It should be a priority to assure adequate reserves to weather unanticipated expenses or periods of decreased income.

Organizational Perspectives

ICANN is a multistakeholder institution operating in the private sector but including the involvement of governments. Throughout its history, ICANN has sought to draw on international resources and to collaborate, coordinate, and cooperate with institutions whose expertise and responsibilities can assist ICANN in the achievement of its goals. ICANN should seek to establish productive relationships with these institutions, cementing its own place in the Internet universe while confining its role to its principal responsibilities.

As part of its normal operation, ICANN engages in self-examination and external review of the effectiveness of its organizational structure and processes. Improvements in all aspects of ICANN operation and structure will increase confidence in the organization and its ability to sustain long-term operation.

Finding and engaging competent participants and leaders in each of ICANN’s constituent parts must be a priority. ICANN should seek to improve its ability to identify from around the world and attract highly qualified staff, executive leadership, board, and supporting organization participants. It is possible and even likely that improvements in the processes by which this process is done today will have significant payoff in the future.

Although ICANN does not bear a specific responsibility for achieving the *Millennium Development Goals* (MDGs) developed during the conduct of the *World Summit on the Information Society* (WSIS)^[8], it has an opportunity to contribute to them both directly and indirectly. Its operation of its IANA functions and support for actors in the domain-name, Internet-address, and standards-development areas provides ICANN with a specific opportunity. Participation in forums dedicated to developing policies for Internet expansion and use offer indirect ways for ICANN to draw upon and provide expertise in these areas.

It has been demonstrated that the presence of ICANN staff in various regions and time zones around the world and familiarity with local languages and customs has been beneficial to parties reliant on ICANN for its services. ICANN should continue to seek ways to improve its effectiveness in this area. The introduction of the Fellowship program that supports the participation of qualified candidates in ICANN-related activities is a vital step in facilitating ICANN's outreach to the developing world. We should pursue expansion of this program through partnerships with other like-minded organizations in the interest of the globalization of ICANN.

It is possible that the present formulation of ICANN as a not-for-profit, charitable research and education entity under California law could be beneficially adapted to a more international framework. As part of its long-term strategic development, ICANN should evaluate a variety of alternatives on the possibility that a change could increase the effectiveness of its operation.

The successful creation of five Regional At-Large Organizations, one in each of ICANN's five regions, needs to be followed by a serious effort to engage these entities in the formulation of ICANN policies and in dialog with the general user community. The various constituency reviews that form part of ICANN's normal processes should address the role of these entities in the conduct of ICANN business. To the extent that civil society is not fully represented through the *Governmental Advisory Committee* (GAC)^[9] and the ALAC/RALO system, an organizational home may be needed to accommodate the interests of that constituency.

The five *Regional Internet Registries* (RIRs)^[10] represent a key element in the Internet and ICANN pantheon. The RIRs have responsibility for allocating IP address space to Internet service providers and sometimes individual end-user organizations. They are the means by which bottom-up global policy is developed and recommended, through the *Number Resource Organization* (NRO)^[11], to ICANN. It will require substantial coordination and cooperation between the RIRs and ICANN to work through the coming years of depletion of available new IPv4 address space and the rising implementation of the new IPv6 address space.

There is little doubt that economic incentives will emerge that will distort fair and neutral IPv4 address-space allocations as the available space is depleted. Minimizing the effect of this transition will be the joint responsibility of ICANN and the RIRs.

Similarly, ICANN's cooperative relationship with the *Root Server Operators*^[12] will also demand coordination and capacity building as IPv4 and IPv6 addresses are associated with old and new domain names and as the IPv6 infrastructure grows. A vital objective is to assure that the IPv6 Internet and the IPv4 Internet are, to the extent possible, completely and totally coterminous. Every termination needs to be reachable through both address spaces. In the absence of this uniformity, some IPv6 addresses may be unreachable from others, defeating the goal of a single, interoperable, and fully reachable network.

Meeting the Challenges

As ICANN approaches the close of its first decade, the operational Internet will be turning 25. In the course of its evolution, it has become a global digital canvas on which a seemingly endless array of applications has been painted. Despite the broad swath of its current applications, it is almost certain that many, many more will be invented. All of them will rely, for the foreseeable future, on the basic architecture of the system, including the global Internet address space and the DNS. But the structure will become more complex. Two parallel address spaces, IPv4 and IPv6, will be in use. ICANN needs to promote the adoption of IPv6 so as to limit the side effects of the exhaustion of the unique address space provided by IPv4.

A vast and new range of non-Latin, internationalized domain names may be registered, certainly at the second or lower levels in the domain-name hierarchy, and many will be proposed for the top level. Their diversity will create new challenges for the protection of users from confusing and potentially abusive registrations. New dispute resolution principles may be needed to deal with domain-name registrations and delegations of new top-level domains. The exposure of ASCII *punycode* strings in browsers or other applications may produce additional stresses in the intellectual property arena (for example **xn--cocacola**).

Digital signatures will play an increasingly important role in validating the assignment of domain names and Internet addresses, and new protocols are certain to be invented and their parameters recorded by the IANA. Infrastructure for the management of digital certificates or other authentication mechanisms will be needed to realize the value of the DNSSEC concept.

More generally, the multilayer architecture of the Internet shows vulnerabilities of various kinds that demand redress. Attacks against the DNS root servers, name resolvers, and general name servers at all levels must be mitigated.

Some of the components of the DNS are actually used to exacerbate the effects of *Denial-of-Service* (DoS) attacks. Although ICANN does not have responsibility for developing the Domain Name technology, it can use its visibility and area of responsibility to highlight the need for increased security measures for the protection of the technical infrastructure of the Internet and to facilitate its implementation where ICANN has a direct involvement in its operation.

An increasing number of mobile devices will become Internet-enabled, as will appliances of all kinds. Access speeds will increase, enabling many new applications and enhancing older ones. All of this activity will contribute to increasing reliance on the Internet for a wide range of functions by an increasingly larger user population. Electronic commerce will continue to expand, placing high priority on the stable, secure, and reliable operation of all aspects of the Internet, including those within ICANN's purview.

Although some of these aspects of the evolution of the Internet will be of direct concern to ICANN, the ICANN organization and processes will need to pay attention to additional matters as well. The business processes that sustain the management of the Internet address space and domain names will almost certainly need to adapt to account for new applications. Some of these applications will monetize various aspects of the Internet in unexpected and innovative ways that will challenge existing policy and procedures. It will be extremely important for ICANN to evolve and strengthen its implementation of multistakeholder policy development. The interests of a wide range of entities must be balanced in the process.

Although adherence to a set of technical standards has allowed millions of component networks and systems to interwork on the Internet, it is also the case that many varying business models have sustained their operation. The richness and diversity of these models is one of the reasons that the Internet has proved to be so resilient in many dimensions. ICANN's policy-development processes need to account for an informed understanding of the economics of these varying business models and the ways in which ICANN policy may affect them.

On the Domain Name side, the development of market-savvy rules of operation for operators will be essential. ICANN needs to assure compliance with policies developed through the ICANN consensus process to establish confidence in the policy processes and their execution. Clear rules for the creation of new *Top Level Domains* (TLDs) of all kinds must be adopted and enforced.

The roles of registrars, registries, wholesale registry operators, root-server operators, regional Internet address registries, governments, and standards and technical research and development bodies, among others, need to be characterized so as to set expectations and permit the establishment of practical working relationships. The documentation of best practices will be beneficial, especially where the introduction of the Internet is new.

In matters of public policy—including but not limited to public safety, security, privacy, law enforcement, conduct of electronic commerce, protection of digital property, and freedom of speech—broad and international agreements may be needed if the Internet is to serve as a useful, global infrastructure. Many of these matters lie outside the formal purview of ICANN, but some ICANN policies and resulting operational practices will contribute to the global framework for life online. ICANN must seek to contribute to public confidence in the Internet and the processes that govern its operation. It cannot accomplish this objective alone. The coordinated and cooperative efforts of many distinct entities will be essential to achieving this goal. At the same time, ICANN must protect its processes from capture or abuse by interests that are inimical to the openness and accessibility of the Internet for everyone.

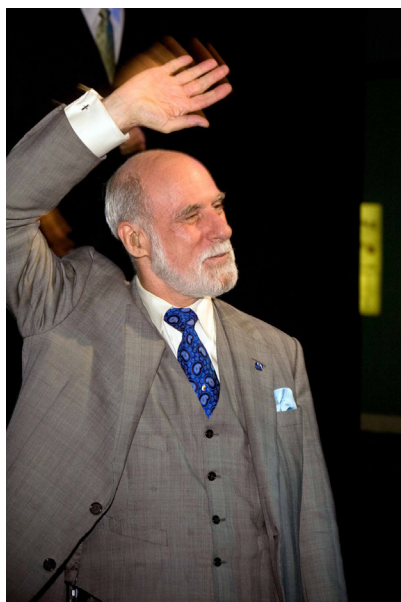
A Collective Goal

As of this writing, only about 1.2 billion people around the world use the Internet. Over the course of the next decade, that number could conceivably quintuple to 6 billion, and users will be depending on ICANN, among many others, to do its part to make the Internet a productive infrastructure that invites and facilitates innovation and serves as a platform for egalitarian access to information. It should be a platform that amplifies voices that might otherwise never be heard and creates equal opportunities for increasing the wealth of nations and their citizens.

ICANN's foundation has been well and truly fashioned. It is the work of many heads and hands. It represents a long and sometimes difficult journey that has called for personal sacrifices from many colleagues and bravery from others. It has demanded long-term commitments, long hours, days, months, and years. It has called upon many to transform passion and zeal into constructive and lasting compromises. ICANN has earned its place in the Internet universe. To those who now guide its path into the future comes the challenge to fashion an enduring institution on this solid foundation. I am confident that this goal is not only attainable but now also necessary. The opportunity is there: make it so.

For Further Reading

- [1] <http://www.iana.org>
- [2] Vint Cerf, "I Remember IANA," *The Internet Protocol Journal*, Volume 1, No. 3, December 1998. Also published as RFC 2468, October 1998.
- [3] <http://www.icann.org/committees/evol-reform/>
- [4] <http://alac.icann.org/>
- [5] Miek Gieben, "DNSSEC: The Protocol, Deployment, and a Bit of Development," *The Internet Protocol Journal*, Volume 7, No. 2, June 2004.
- [6] <http://icann.org/topics/idn/>
- [7] <http://icann.org/udrp/>
- [8] <http://www.itu.int/wsis/index.html>
- [9] <http://gac.icann.org/web/index.shtml>
- [10] Daniel Karrenberg, Gerard Ross, Paul Wilson, and Leslie Nobile, "Development of the Regional Internet Registry System," *The Internet Protocol Journal*, Volume 4, No. 4, December 2001.
- [11] <http://nro.org/>
- [12] <http://www.root-servers.org/>



Photographer: Vanessa Stump

VINTON G. CERF is vice president and chief Internet evangelist for Google. In this role, he is responsible for identifying new enabling technologies to support the development of advanced Internet-based products and services from Google. He is also an active public face for Google in the Internet world. Cerf is the former senior vice president of Technology Strategy for MCI. In this role, he helped guide corporate strategy development from a technical perspective. Previously, he served as MCI's senior vice president of Architecture and Technology, leading a team of architects and engineers to design advanced networking frameworks, including Internet-based solutions for delivering a combination of data, information, voice, and video services for business and consumer use.

Widely known as one of the "Fathers of the Internet," Cerf is the co-designer of the TCP/IP protocols and the architecture of the Internet. In December 1997, President Clinton presented the U.S. National Medal of Technology to Cerf and his colleague, Robert E. Kahn, for founding and developing the Internet. Kahn and Cerf were named the recipients of the ACM Alan M. Turing Award, sometimes called the "Nobel Prize of Computer Science," in 2004 for their work on the Internet protocols. In November 2005, President George Bush awarded Cerf and Kahn the Presidential Medal of Freedom for their work. The medal is the highest civilian award given by the United States to its citizens.

Prior to rejoining MCI in 1994, Cerf was vice president of the Corporation for National Research Initiatives (CNRI). As vice president of MCI Digital Information Services from 1982 to 1986, he led the engineering of MCI Mail, the first commercial e-mail service to be connected to the Internet.

During his tenure from 1976 to 1982 with the U.S. Department of Defense Advanced Research Projects Agency (DARPA), Cerf played a key role leading the development of Internet and Internet-related packet-data and security technologies.

Vint was seated on the ICANN Board of Directors at the 1999 annual meeting, having been selected by the Protocol Supporting Organization. He was then selected by the nominating committee for a term on the board of directors that ran from June 2003 through the 2004 annual meeting. At the end of that term, he was selected by the 2004 nominating committee to an additional term, which ran from the end of the 2004 annual meeting through the conclusion of the ICANN annual meeting in 2007. He served as founding president of the Internet Society from 1992 to 1995, and in 1999 served a term as chairman of the board. In addition, Cerf is honorary chairman of the IPv6 Forum, dedicated to raising awareness and speeding introduction of the new Internet Protocol. Cerf served as a member of the U.S. Presidential Information Technology Advisory Committee (PITAC) from 1997 to 2001 and serves on several national, state, and industry committees focused on cyber security. Cerf sits on the board of directors for the Endowment for Excellence in Education, Avanex Corporation, and the ClearSight Systems Corporation. Cerf is a Fellow of the IEEE, ACM, and American Association for the Advancement of Science, the American Academy of Arts and Sciences, the International Engineering Consortium, the Computer History Museum, and the National Academy of Engineering.

Cerf is a recipient of numerous awards and commendations in connection with his work on the Internet, including the Marconi Fellowship, Charles Stark Draper Award of the National Academy of Engineering, the Prince of Asturias Award for science and technology, the National Medal of Science from Tunisia, the Alexander Graham Bell Award presented by the Alexander Graham Bell Association for the Deaf, the NEC Computer and Communications Prize, the Silver Medal of the International Telecommunications Union, the IEEE Alexander Graham Bell Medal, the IEEE Koji Kobayashi Award, the ACM Software and Systems Award, the ACM SIGCOMM Award, the Computer and Communications Industries Association Industry Legend Award, installation in the Inventors Hall of Fame, the Yuri Rubinsky Web Award, the Kilby Award, the Yankee Group/Interop/Network World Lifetime Achievement Award, the George R. Stibitz Award, the Werner Wolter Award, the Andrew Saks Engineering Award, the IEEE Third Millennium Medal, the Computerworld/Smithsonian Leadership Award, the J.D. Edwards Leadership Award for Collaboration, the World Institute on Disability Annual Award, and the Library of Congress Bicentennial Living Legend medal. In December 1994, People magazine identified Cerf as one of that year's "25 Most Intriguing People."

In addition to his work on behalf of MCI and the Internet, Cerf has served as a technical advisor to production for the "Gene Roddenberry's Earth: Final Conflict" television series and made a special guest appearance on the program in May 1998. Cerf has appeared on television programs NextWave with Leonard Nimoy and on World Business Review with Alexander Haig and Caspar Weinberger. He is also a distinguished visiting scientist at the Jet Propulsion Laboratory, where he is working on the design of an interplanetary Internet.

Cerf holds a Bachelor of Science degree in Mathematics from Stanford University and Master of Science and Ph.D. degrees in Computer Science from UCLA. He also holds honorary doctorate degrees from the Swiss Federal Institute of Technology (ETH), Zurich; Luleå University of Technology, Sweden; University of the Balearic Islands, Palma; Capitol College, Maryland; Gettysburg College, Pennsylvania; George Mason University, Virginia; Rovira i Virgili University, Tarragona, Spain; Rensselaer Polytechnic Institute, Troy, New York; the University of Twente, Enschede, The Netherlands; Brooklyn Polytechnic; and the Beijing University of Posts and Telecommunications.

Cerf's personal interests include fine wine, gourmet cooking, and science fiction. Cerf and his wife Sigrid were married in 1966 and have two sons, David and Bennett.

E-mail: vint@google.com

Remembering Itojun: The IPv6 Samurai

by Bob Hinden, Nokia

“Itojun” (Dr. Junichiro Hagino) passed away on October 29, 2007. He was 37 years old. Memorial events were held in Tokyo in November and in Vancouver at the IETF meeting in December.

Itojun was an active participant in the IETF and a member of the IAB from 2003 to 2005. He worked as a Senior Researcher at the *Internet Initiative Japan* (IIJ) and was a member of the board of the *Widely Integrated Distributed Environment* (WIDE) project. He was a strong supporter of open standards development and open software, working as a core researcher at the KAME project, a joint effort of six companies in Japan to provide a free stack of IPv6, IPsec, and Mobile IPv6 for BSD variants, from 1998 to 2006.

Itojun was totally dedicated to the development and deployment of IPv6. Most of his work was centered around building a much larger worldwide Internet based on IPv6. He was simply the “IPv6 Samurai.”



Photographer: Diane Bruce

Quotes from Internet Colleagues

Steve Deering: “Those of us who got to know Itojun through his work in the Internet Engineering Task Force have lost a dear friend and much-admired colleague. From the day he arrived at his first IETF meeting, he won the respect of all in the way most honored by Internet engineers: by helping to build consensus based on running code. Moreover, he provided the best possible example of collaboration, generosity, and leadership, making not only extraordinary technological contributions but also many friends and a better world. His untimely passing is a huge loss to all who knew him, and to all those who will never have that chance.”

Randy Bush: “An open heart, a big soul, and very kind and patient. A very special person. He wrote a lot of great code and got great joy from doing so.”

Marc Blanchet: “Itojun adopted the Samurai’s philosophy in his life: *Bushido*, which consists of values such as Honesty, Justice, Courtesy, Heroic Courage, Honor, Compassion, Sincerity, Duty, and Loyalty. Very difficult to achieve, he encompassed all these. Moreover, he was always available to help, anyone, without judging. His intelligence, his competency, and his dedication has inspired a generation of network engineers for the project he took as a mission: IPv6. Many computers in the world now run his code. My family always enjoyed meeting Itojun. He was always interested in sharing his knowledge with my children, even with the French-to-Japanese-through-English language barrier. Itojun, it was an honor to know you and to meet you. You will always be a source of inspiration to me, to my family, and to many network engineers in the world. We miss you.”

Rod Van Meter: “I didn’t know Itojun very well; I met him for the first time about five years ago at an IPv6 meeting in the Silicon Valley, once or twice in between, and then spent three days at the WIDE Camp this past September co-supervising (with Bill Manning, Brad Huffaker, and Kenji Saito) a group of students trying to establish long-term goals for WIDE in the area of naming. Itojun was gentle but insistent with students, a good mentor. That was the last time I saw him. Go in peace, Itojun.”

Joel Jaeggli: “He cared more for the people who were going to use the code and the product of his and our labor than anyone would have had a right to expect. The Itojun that I know, our friend, has been taken from us, but we’ll be the beneficiary of the fact that he cared, for decades.”

Itojun IPv6 Fund

Itojun's family has expressed sincere appreciation to all who attended the memorial and funeral services. His family has set up a memorial fund in Itojun's name under the directorship of the IETF/Internet Society. The fund will be used to award an R&D grant to a person who has contributed to the deployment and further advancement of IPv6. ISOC has set up an e-mail address to accept commitments for the *Itojun IPv6 Fund*. The address is: **itojun-fund@isoc.org**

The procedure for making contributions is being developed; if you wish to contribute now, please send a note to the e-mail address describing the amount you want to contribute (and in what currency), and ISOC will collect the funds.

ROBERT HINDEN is a Nokia Fellow at Nokia and is located in Mountain View, California, USA. He has been involved in the Internet since it was a research project at ARPA. He developed one of the first TCP/IP implementations and his team at Bolt, Beranek, and Newman, Inc. built and operated the routers that formed the early Internet backbone. He was co-recipient of the 2008 IEEE Internet Award "For pioneering work in the development of the first Internet routers." He has been active in the IETF since 1985 and is the author of 35 RFCs. He was recently appointed to a position on the IETF Administrative Oversight Committee (IAOC) and co-chairs the IPv6 Maintenance (6man) working group. Prior to this he served on the Internet Architecture Board (IAB), was Area Director for Routing in the Internet Engineering Steering group from 1987 to 1994, and chaired the IPv6, Virtual Router Redundancy Protocol, Simple Internet Protocol Plus, IP over ATM, and Open Routing working groups. Hinden is also a member of the RFC Editorial Board. He holds a B.S.E.E. and a M.S. in Computer Science from Union College, Schenectady, New York. E-mail: **bob.hinden@nokia.com**

Book Review

Network Routing

Network Routing: Algorithms, Protocols, and Architectures, by Deepankar Medhi and Karthikeyan Ramasamy, Morgan Kaufmann Publishers, ISBN-13:978-0120885886, 2007,
<http://www.NetworkRouting.net>

Routing is a fundamental architectural component of any network, and in this book the authors examine in detail the routing technologies of the Internet and the *Public Switched Telephone Network* (PSTN).

Organization

The book is divided into five parts, with an additional advanced section provided on CDROM. The first part examines the fundamentals of routing technology, looking in detail at the basic approaches of distance-vector and link-state routing. The second part looks at the routing protocols used in the Internet today, as well as Traffic Engineering. The third part addresses routing in the PSTN, examining the SS7 signaling protocol and the overall architecture of the PSTN. The next part explores the internal architecture of routers, address-lookup algorithms, and packet-classification techniques. Finally, the authors consider topics encompassed in the so-called “Next-Generation Network,” including *Quality-of-Service Routing*, *Multiprotocol Label Switching* (MPLS), and *Voice over IP* (VoIP). The advanced-topic section includes a more detailed examination of packet-switching approaches, scheduling, and conditioning. This book is positioned as a graduate-level text, and each chapter is accompanied by exercises that review the material.

The book covers a broad range of material: each topic has been the subject of entire books. The level of detail in the book varies considerably. In some instances, such as in the area of IP Traffic Engineering, it presents a highly detailed mathematical analysis of aspects of the topic, whereas in other instances, such as in the treatment of the *Border Gateway Protocol* (BGP), the material appears to be obviously condensed. I was expecting a little more use of algorithms to illustrate routing concepts, and found at times the mathematical analysis to be unhelpful in terms of understanding the underlying problem space being described.

Comparison

In this area of Internet routing, any publication is inevitably compared to Radia Perlmann’s book *Interconnections: Bridges, Routers, Switches and Internetworking Protocols*, and this book is no exception. To my mind it falls a little short of this rather demanding standard. Radia spends some time discussing the underlying rationale as to why a particular technology was devised for a given problem space, and also discusses the strengths and limitations of the technology in various areas of application.

In *Network Routing* the authors limit their approach to a description of the technology by looking at packet payloads and protocol interactions and numerous deployment scenarios that illustrate the features of this particular routing technology. The consideration of choices made in the development of the protocol, and the consequent implications of such design choices, are missing in such a treatment, and the reader is often left wondering why a routing protocol has chosen to support certain functions but not others.

I found this to be a very ambitious book, because it appears to position itself both as a reference publication on routing technologies and architecture and also on the description of routing protocols, while at the same time wanting to encompass the role of a course text. This goal could have been attainable if the book had chosen a tighter focus, but the all-encompassing approach that led to the inclusion of considerations of the PSTN topics makes the outcome less than fully satisfying.

Recommended

However, the book manages to bring together the basic topics in routing in both the Internet and the PSTN, and it not only includes a good description of the routing technologies in use today, but also looks at some of the advanced topics in routing today. I found the major strength of the book in its role as a graduate-course text, where there is sufficient description of the topic to lead into further reading of current research papers and more-detailed technical material. Although the book has some shortcomings, I'd certainly recommend it as a suitable addition to the shelves of any professional in the area of Internet routing technologies and architecture.

—Geoff Huston
gih@apnic.net

The Author Responds:

I thank Geoff Huston for writing a well-thought-out review; in general, this review is fair. This book was certainly an ambitious project. I wanted to do it as I've investigated various routing protocols for almost two decades—and many people I talked to thought that it would be useful to have such a book. In fact, Dave Clark, when he read the original book proposal, wrote "It is ambitious—there may be issues of how much depth they can get on all these topics in one book," but felt that "...the approach is distinctive and very valuable. So I support the idea." As can be observed from the book, the depth on different topics remained a major trade-off we pondered without making the book go over 1000 pages (with 140 pages on CD-ROM it came pretty close).

There were a few “design” decisions I deliberately made in organizing and writing this book. One of them was based on years of teaching and interacting with industry folks: I decided to divide materials broadly on “how and why” away from “what;” this approach is somewhat surprising, but people’s learning style seems to fall into these two categories (certainly there are overlaps). Therefore, details on “how and why” of different protocols went into Chapter 3 (and for algorithms into Chapter 2), while details on “what” went with chapters on specific protocols such as OSPF or BGP. Similarly, I also separated out the topic of “how” routing in the global Internet works and is organized (such as public exchange points) from the chapter on BGP. Secondly, we separated math parts from non-math parts—this way, those who are interested, for example, in detailed Traffic Engineering modeling can read the relevant chapters. Others may skip them and read just the first couple of overview sections; it should be noted that math-oriented chapters are generally organized from simple concepts to difficult concepts. Thirdly, we covered address lookup, packet filtering and classification, and router architectures separately because they can be read independently; Karthik brought his wealth of experience in writing these chapters.

I want to take this opportunity to respond to a few of Geoff’s comments:

1. “...expecting a little more use of algorithms to illustrate routing concepts.” I suspect that Geoff didn’t think that Chapters 2 and 3 covered enough, although these chapters included details illustrative of distance-vector protocols, link-state protocols, path-vector protocols (and their pitfalls), and so on. As stated previously, by design of the book, illustrative examples of routing concepts were separated from specific protocols so that readers can read different portions of the material according to their interests. As an indicator to the reader, each chapter starts with a brief “reading guideline” (which is a unique feature of this book) that states how the material is organized and its relation to other chapters or sections in the rest of the book.
2. “The consideration of choices made in the development of the protocol, and the consequent implications of such design choices are missing in such a treatment.” We did indeed cover these aspects in many instances. For example, the book covers why, for I-BGP scalability, the route reflector or the confederation approach are needed; why route flap damping was developed; why ROUTE-REFRESH was added; what MPLS was trying to solve that IP-only couldn’t do at that time; the need for age with Sequence Number field in link-state protocols; what led to the development of dynamic call routing from hierarchical routing in the PSTN, and so on.

That said, I did not include certain discussions because some choices on protocols have been based on personality clashes and “camps;” I felt that this is not easily explainable in many instances—trying to do so would require quite a bit of discussion, and could potentially divert from the main focus of the book. For example, I explained why the route reflector or confederation approach was needed for I-BGP scalability, but I didn’t discuss why both route reflector and confederation approaches were developed simultaneously when both convey the same idea conceptually.

3. “... the all-encompassing approach that led to the inclusion of considerations of the PSTN topics into this book make the outcome less than fully satisfying.” I included routing in the PSTN because of its historic context, and particularly to make readers aware of the evolution from hierarchical routing to dynamic routing and recent changes in routing due to Local Number Portability—these lessons are important ones to learn for anyone interested in routing or designing future routing protocols. Secondly, many concepts in MPLS/GMPLS have parallels in the PSTN, thus certain aspects in MPLS/GMPLS are easily explainable if a reader is familiar with PSTN details. We therefore felt it was appropriate to include all this material in one place. Furthermore, control- and data-path separation in GMPLS is strikingly similar to separation of signaling in PSTN through SS7 from actual voice communication. Thus, lessons learned from failure propagation from SS7 to voice paths are relevant lessons to be aware of for anyone involved in deploying GMPLS-based networks. Lastly, to discuss VoIP routing, it is critical to tie into PSTN because in the real operational environment PSTN-Internet interworking for VoIP routing is expected to remain prevalent for years to come.

Finally, the “barrier to entry” in learning about routing is very high, especially for entry-level professionals—I’ve attempted to position the book as both a text and a reference for professionals. Thus, I very much appreciated Geoff’s concluding comment “... as a suitable addition to the shelves of any professional in the area of Internet routing technologies and architecture.”

—Deep Medhi
dmedhi@umkc.edu

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. Contact us at ipj@cisco.com for more information.

Nii Quaynor Receives 2007 Postel Service Award

The *Internet Society* (ISOC) has awarded pioneering Internet engineer Nii Quaynor the prestigious *Jonathan B. Postel Service Award* for 2007 for his leadership in advancing Internet technology in Africa and galvanizing technologists to improve Internet access and capabilities throughout the continent. ISOC presented the award, including a \$20,000 [USD] honorarium, during the 70th meeting of the *Internet Engineering Task Force* (IETF) in Vancouver, BC, Canada.

“Dr. Quaynor has selflessly pioneered Internet development and expansion throughout Africa for nearly two decades, enabling profound advances in information access, education, healthcare and commerce for African countries and their citizens,” said ISOC president Lynn St. Amour. “Today, Dr. Quaynor continues to champion not just technological advances but also African involvement in Internet standards, processes and deployments, discussion on Internet policies and regulations, and ensuring African interests are well-represented globally. He has shaped a community of Africans who share his vision and reflect the dedication shown by Jon Postel.”

“I am humbled by the award and what Jon Postel represents to our community in Africa. Jon Postel’s efforts and the global view he maintained on the operation of the *Domain Name System* and the numbering services assured that Africa would share in the Internet growth and early. I thank the Internet Society for the recognition and am very pleased to be associated with Jon’s memorial,” said Dr. Nii Quaynor. “We will work to develop more African engineers to meet the fast network growth needs of the region, being a late starter, and to join the technical policy processes. Our overall objective is to strengthen education and research in network technologies in Africa.”

The annual ISOC award is named after Dr. Jonathan B. Postel to commemorate his extraordinary stewardship exercised throughout his thirty-year career in networking. Between 1971 and 1998, Postel managed, nurtured and transformed the RFC series of notes, which encompasses the technical specifications and recommendations for the Internet and was created by Steve Crocker in 1969 as a part of his work on the ARPANET, the forerunner of today’s Internet. Postel was a founding member of the Internet Architecture Board and the first individual member of the Internet Society, where he also served as a trustee until his untimely death.

Dr. Quaynor is chairman of *Network Computer Systems* (NCS) Ghana.COM and a professor of computer science at University of Cape-Coast, Ghana. He is also the convener of the *African Network Operators Group* (AfNOG), a network technology transfer institution since 2000 and the founding chairman of AfriNIC, the African numbers registry.

Dr. Quaynor began his pioneering Internet work in Africa in 1993 when he returned to his home country of Ghana to establish the first Internet Service operated by NCS in West Africa. At NCS, he and his team worked on the early development of the Internet in Africa. Today, there are more than 43 million Internet users in Africa.

Prior to NCS, Dr. Quaynor worked with Digital Equipment Corporation in the United States from 1977 till 1992. In 1979, he established the Computer Science department at the University of Cape Coast, Ghana. Dr. Quaynor graduated from Dartmouth College in 1972 with B.A (Engineering Science) and received a Ph.D. (Computer Science) in distributed systems in 1977 from State University of New York at Stony Brook.

The Jonathan B. Postel Service Award was established by the Internet Society to honor those who, like Postel, have made outstanding contributions in service to the data communications community. The award is focused on sustained and substantial technical contributions, service to the community, and leadership. With respect to leadership, the nominating committee places particular emphasis on candidates who have supported and enabled others in addition to their own specific actions.

Previous recipients of the Postel Award include Jon himself (posthumously and accepted by his mother), Scott Bradner, Daniel Karrenberg, Stephen Wolff, Peter Kirstein, Phill Gross, Jun Murai, Bob Braden, and Joyce K. Reynolds. The award consists of an engraved crystal globe and \$20,000 [USD]. This year's award is sponsored in part by Afilias Global Registry Services. For more information about ISOC, please visit: www.isoc.org

Steps Taken for Multilingual Internet

The *Internet Corporation for Assigned Names and Numbers* (ICANN), the *International Telecommunication Union* (ITU), and the *United Nations Educational, Scientific and Cultural Organization* (UNESCO) will collaborate on global efforts to forge universal standards towards building a multilingual cyberspace. The three agencies organized a workshop on this subject during the second *Internet Governance Forum* (IGF) which took place in Rio de Janeiro, Brazil from 12 to 15 November 2007.

The Internet is a key factor in developing a more inclusive and development-oriented information society, which stresses plurality and diversity instead of global uniformity. Multilingualism is a key concept to ensure cultural diversity and participation for all linguistic groups in cyberspace. There is growing concern that hundreds of local languages may be sidestepped, albeit unintentionally in the radical expansion of Internet communication and information. The *World Summit on the Information Society* (WSIS) recognized the importance attached to linguistic diversity and local content, with UNESCO given the responsibility to coordinate implementation of the *Summit Action Line*.

“The discussions at this multilingualism workshop—combined with our current evaluation of *Internationalized Domain Names* (IDNs)—are going to help ICANN keep moving toward full implementation of Internationalized Domain Names,” said Dr Paul Twomey, ICANN’s President and CEO. “ICANN is in the midst of the largest ever evaluation of IDNs at the top level.”

Thanks to ICANN’s evaluation of Internationalized Domain Names, Internet users around the globe can now access wiki pages (see <http://idn.icann.org/>) with the domain name **example.test** in the 11 test languages—Arabic, Persian, Chinese (simplified and traditional), Russian, Hindi, Greek, Korean, Yiddish, Japanese and Tamil. The wikis will allow Internet users to establish their own sub pages with their own names in their own language; one suggestion is: **example.test/yourname**

Domain Names, which are currently mainly limited to characters from the Latin or Roman scripts, are seen as an important element in enabling the multilingualization of the Internet, reflecting the diverse and growing language needs of all users. “ITU is fully committed to assist its membership in promoting the diversity of language scripts for domain names,” said Dr Hamadoun Touré, Secretary-General of ITU. “This workshop represents an important opportunity to strengthen the need for cooperation with relevant organizations, such as UNESCO, the *World Intellectual Property Organization* (WIPO) and ICANN among others to ensure Internet use and advancement across language barriers.”

The Plenipotentiary Conference of ITU, which took place in Antalya, Turkey in November 2006, recognized the need to make Internet content available in non-Latin based scripts. Internet users are more comfortable reading or browsing through texts in their own language and a multilingual Internet is essential to make it more widely accessible. The WSIS outcomes also focused on the commitment to work towards multilingualization of the Internet as part of a multilateral, transparent and democratic process involving governments and all stakeholders.

UNESCO, joined by both ITU and ICANN, seeks to convene all major stakeholders around the world towards an agreement on universal standards regarding language issues in cyberspace. Such issues are far broader than the single issue of IDNs as they extend to standards for fonts and character sets, text encoding, language implementations within major computer operating systems, content development tools, automatic translation software, and search engines across languages. Ultimately, equitable access to information can be only achieved if we resolve language barriers at the same time we build communications infrastructures and capacity building programs.

RIPE Community Resolution on IPv4 Depletion and Deployment of IPv6

During the RIPE 55 meeting in Amsterdam in October 2007, the RIPE community agreed to issue the following statement on IPv4 depletion and the deployment of IPv6:

“Growth and innovation on the Internet depends on the continued availability of IP address space. The remaining pool of unallocated IPv4 address space is likely to be fully allocated within two to four years. IPv6 provides the necessary address space for future growth. We therefore need to facilitate the wider deployment of IPv6 addresses.

While the existing IPv4 Internet will continue to function as it currently does, the deployment of IPv6 is necessary for the development of future IP networks.

The RIPE community has well-established, open and widely supported mechanisms for Internet resource management. The RIPE community is confident that its *Policy Development Process* meets and will continue to meet the needs of all Internet stakeholders through the period of IPv4 exhaustion and IPv6 deployment.

We recommend that service providers make their services available over IPv6. We urge those who will need significant new address resources to deploy IPv6. We encourage governments to play their part in the deployment of IPv6 and in particular to ensure that all citizens will be able to participate in the future information society. We urge that the widespread deployment of IPv6 be made a high priority by all stakeholders.”

For more information, see: <http://ripe.net/ripe/>

Upcoming Events

The next *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) will be held in Taipei, Taiwan from February 20th to 29th, 2008. As usual, this conference is co-located with an APNIC Open Policy Meeting. For more information about these events see: <http://www.apricot2008.net/> and <http://www.apnic.net/meetings/25/index.html>

The *Internet Engineering Task Force* (IETF) will meet in Philadelphia, Pennsylvania, March 9–14 and “somewhere in Europe” July 27–August 1. (The announcement of the exact location is expected soon). The final IETF meeting in 2008 will take place in Minneapolis, Minnesota, November 16–21. For more information see: <http://www.ietf.org/meetings/0mtg-sites.txt>

The *Internet Corporation for Assigned Names and Numbers* (ICANN) will meet in New Delhi, India, February 10–15, and in Paris, France, June 22–27. See: <http://icann.org/meetings/>

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2007 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive, M/S SJ-7/3
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA
--

The Internet Protocol Journal

March 2008

Volume 11, Number 1

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
IDNs	2
LISP	23
Book Review.....	37
Fragments	39
Call for Papers.....	43

The *Domain Name System* (DNS) was not designed to support anything beyond 7-bit ASCII characters. Thus my middle name, Jørgen, or my colleague's surname, Fältström, cannot be used in a domain name. In fact, even using such strings on the left side of the @-sign—or in the body of an e-mail message—is problematic. We often find ourselves ignoring this limitation, using either “Jorgen” and “Faltstrom” or in some cases the two-letter convention “Joergen” and “Faeltstroem.” As Scandinavians, Mr. Fältström and I are relatively lucky in that our languages contain only three characters in addition to those that can be represented by 7-bit ASCII. This, of course, isn't true for such languages as Arabic, Chinese, Japanese, or Korean, to name just a few. The IETF, ICANN, and others have been working hard to design and deploy a system that will allow native characters to appear in the DNS. Our first article discusses these efforts, known collectively as *Internationalized Domain Names* (IDNs). Geoff Huston gives an overview of IDNs and describes the many technical and political challenges that must be overcome in order to deploy such a system.

Recent activities have focused much attention on IPv6 deployment. Experiments have been conducted at several major Internet events (NANOG, APRICOT, and IETF) to “turn off” IPv4 for a period of time to test connectivity and interoperability to the outside world. You can read more about these experiments in our “Fragments” section on page 41. Such experiments provide valuable information about what works and what doesn't, and several more IPv4 “outages” are planned for 2008 and beyond. At the same time, researchers have been looking at ways to scale the routing system of the Internet, regardless of IP protocol version. One such approach is the *Locator/Identifier Separation Protocol* (LISP), which Dave Meyer describes in our second article.

The next issue of *The Internet Protocol Journal*, to be published sometime in June 2008, will be our Tenth Anniversary issue. We would love to hear your reflections on the last ten years of this journal and about the Internet as a whole over the same time period. Send your Letters to the Editor to ipj@cisco.com

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

Internationalizing the Domain Name System

by Geoff Huston, APNIC

Considering the global reach of the Internet, internationalizing the network sounds like a tautology. Surely the Internet is already truly “international,” isn’t it? The Internet reaches around the globe to every country, doesn’t it? And no matter where you may travel these days, an Internet café is just around the corner. How much more “international” can you get?

But maybe I’m just being too parochial here when I call it a tautology. I use a dialect of the English language, and all the characters I need are contained in the *Western Latin* character set. Therefore, I avoid using a non-English language on the Internet; the only language I use on the Internet is English, and all the characters I need are encompassed in the ASCII character set. If I tried to use the Internet with a language that has a non-Latin character set and a different script, my experience would probably be different—and acutely frustrating. If my native language used a different script and a different text flow than English, I would probably give the Internet an extremely low score for ease of use. It is not as simple as managing glyph sets to represent the characters of the language; although it is relatively easy to present pictures of characters in a variety of fonts and scripts, using them in an intuitive and natural way in the context of the Internet becomes more challenging.

Mostly what is needed is good *localization*, or adapting the local computing environment to suit local linguistic needs. This environment may include support for additional character sets and additional language scripts, and perhaps altering the direction of text flow, or even the entire layout of the information.

For example, Japanese is traditionally written in a format called *Tategaki*. In this format, the text flows in columns going from top to bottom, with columns ordered from right to left. Modern Japanese also uses another writing format, called *Yokogaki*. This writing format is identical to that of European languages such as English, where the text flows from left to right in successive rows from top to bottom.

Today, the left-to-right direction is dominant in Japanese *Kana*, Chinese characters, and Korean *Hangul* for horizontal writing. This change is due partly to the influence of English, and partly to the increased use of computerized typesetting and word-processing software, most of which does not directly support right-to-left layout of East Asian languages. It would appear that even *Yokogaki* is an outcome of the lack of capability of IT systems to correctly cope with localization.^[1]

One topic, however, does not appear to have a compellingly obvious localization solution in this multilingual environment: the *Domain Name System* (DNS). The subtle difference here is that the DNS is the “glue” that binds all users’ language symbols together, and performing localized adaptations to suit local language use needs is not enough. The DNS spans the entire network, so what works for me in the DNS must also work for you. What we need is a means to allow the use of all of these language symbols within the same system, or *internationalization*.

The DNS is the most prevalent means of initiating a network transaction, whether it is a *BitTorrent* session, the Web, e-mail, or any other form of network activity. But the DNS name string is not just an arbitrary string of characters. What you find in the DNS is most often a sequence of words or their abbreviations, and the words are generally English words, using characters drawn from a subset of the Latin character set. Perhaps unsurprisingly, some implementations of the DNS also assume that all DNS names must be constructed only from this ASCII character set, and these implementations are incapable of supporting a larger character repertoire. If you want to use a larger character set in order to represent various diacritics, such as acute and grave symbols, umlauts and similar marks, then the deployed DNS can be resistant to this use, and may provide incorrect responses to queries that include such characters. And if you want to use words drawn from languages that do not use the western script for their characters, such as Japanese or Thai, for example, then the DNS is highly resistant to this form of multilingual use.

Latin and Roman Alphabets

The default Latin alphabet is the Roman^[2] alphabet, supplemented with G, J, U, W, Y, Z, and lowercase variants. Additional letters may be formed:

- As *ligatures*, as W was from VV, for example *Æ* (*ash*) from AE, *oethel* *ƿ* from OE, *eszett* *ß* from fz (long s + z), *engma* *ŋ* from NG, *ou* *Ů* from OU, *Ñ* from NN, or *ä* from ae
- By *diacritics*, such as Å, Č, and Ů
- As *digraphs*, such as fi and fl
- By modification, as J was from I, G from C, Ø from O, *eth* *Ð* from D, *yogh* *ȝ* from G, or *schwa* *ə* from E
- By borrowing from another alphabet entirely, as *thorn* *Þ* and *wynn* *ƿ* were from Futhark (Runic)

Over the years we have done a reasonable job of at least displaying non-Latin-based scripts within many applications, and although at times it appears to represent a less-than-reasonable compromise, it is possible to enter non-Latin characters on computer keyboards. So it appears to be possible to customise a local computing environment to use a language other than English in a relatively natural way.

But what happens when we extend the scope to consider multilingual support in the wider world of the Internet?

Again the overall story is not all that bad. We can use non-Latin character scripts in e-mail, in all kinds of Web documents, and in a wide variety of network applications. We can tag content with a language context to allow display of the content in the correct language using the appropriate character sets and presentation glyphs. However, until recently, one area continued to stick steadfastly to its ASCII roots: the DNS. This article addresses DNS internationalization, or *Internationalized Domain Names* (IDNs).

What do we mean when we talk of “internationalizing the DNS”? It refers to an environment where English, and the Latin character set, is just one of many languages and scripts in use, and where a communication is initiated in one locale and then the language and presentation are preserved wherever the communication is received.

Terminology

The following terms are used in this article:

Language: A language uses characters drawn from a collection of scripts.

Script: A script is a collection of characters that are related in their use by a language.

Character: A character is a unit of a script.

Glyph: The presentation of a character within the style of a font is called a glyph.

Font: A font is a collection of glyphs encompassing a script character set that share a consistent presentation style.

Multiple languages can use a common script, and any locale or country may use many languages, reflecting the diversity of its population and the evolution of local dialects within communities.

It is also useful to remember the distinction between internationalization and localization. *Internationalization* is concerned with providing a common substrate that many—preferably all—languages and all users can use, whereas *localization* is concerned with the use of a particular language within a particular locale and within a defined user population. Unsurprisingly, the two concepts are often confused, particularly when true internationalization is often far more difficult to achieve than localization.

Internationalizing the DNS

The objective is the internationalization of the DNS, such that the DNS can support the union of all character sets while preserving the absence of ambiguity and uncertainty in terms of resolution of any individual DNS name. We need to describe all possible characters in all languages and allow their use in the DNS. So the starting point is the “universal character set,” and that appears to be Unicode.

One of the basic building blocks for internationalization is a character set that is the effective union of all character sets. *Unicode*^[3] is intended to be such a universal encoding of characters (and symbols) in the contexts of all scripts and all languages. The current version of the *Unicode Standard*, Version 5.0, contains 98,884 distinct coded graphic characters.

A sequence of Unicode code points can be represented in multiple ways by using different character encoding schemes in a *Unicode Transformation Format* (UTF). The most commonly used schemes are UTF-8 and UTF-16.

UTF-8 is a variable-length encoding using 8-bit words, meaning that different code points require different numbers of bytes. The larger the index number of a code point, the more bytes are required to represent it using UTF-8. For example, the first 127 Unicode code points, which correspond exactly to the values used by the ASCII character set (which maps only 127 characters), can be represented using only 8 bits in UTF-8, using the same 8-bit values as in ASCII. UTF-8 can require up to 32 bits to encode certain code points. A criticism of UTF-8 is that it “penalizes” certain scripts by requiring more bytes to represent their code points. The IETF has made UTF-8 its preferred default character encoding for internationalization of Internet application protocols.

UTF-16 is a variable-length character encoding using 16-bit words. Characters in the *Basic Multilingual Plane* are mapped into a single 16-bit word, with other characters mapped into a pair of 16-bit words.

UTF-32 is a fixed-length encoding that uses 32 bits for every code point. This encoding tends to make for a highly inefficient coding that is, generally, unnecessarily large, because most language uses of Unicode draw characters from the Basic Multilingual Plane, making the average code size 16 bits in UTF-16 as compared to the fixed-length 32 bits in UTF-32. For this reason UTF-32 is far less commonly used than UTF-8 and UTF-16.

But languages, which we humans change in various ways every day, are not always definitive in their use of characters, and Unicode has some weaknesses in terms of identifying a context of a script and a language for a given character sequence. The common approach to using Unicode encodings in application software is to use an associated “tag,” allowing content to be tagged with a script and an encoding scheme. For example, a content tag might read: “This text has been encoded using the KOI-8 encoding of the CYRILLIC script.”

Tagging allows for decoding of the encoded characters in the context of a given script and a given language. This decoding has been useful for e-mail or Web page content, but tagging breaks down in the context of the DNS. There is no natural space in DNS names to contain language and script tags, implying that attempting to support internationalization in the DNS has to head toward a “universal” character set and a “universal” language context. Another way of looking at this situation is that the DNS must use an implicit tag of “all characters and all languages.”

The contexts of the use of DNS names have numerous additional artefacts. What about domain-name label separators? This “dot” between DNS “words,” or a DNS label separator, is an ASCII period character. In some languages, such as Thai, for example, there is no natural use of such a label separator. In a similar vein, are URLs intended to be visible to end users? If so, then we may have to transform the punctuation components of the URL into the script of the language. Therefore, we may need to understand how to manage protocol strings, such as “http:” and separators such as the “/” character. To complete the integrity of the linguistic environment, these elements may also require local presentation transformations.

For example, the Thai alphabet uses 44 consonants and 15 basic vowel characters, which are horizontally placed, from left to right, with no intervening space, to form syllables, words, and sentences. Vowels associated with consonants are nonsequential: they can be located before, after, above, or below their associated consonant, or in a combination of these positions. The latter in particular causes problems for computer encoding and text rendering^[4].

The DNS name string reads left to right, and not right to left or top to bottom as in other script and language cultures. How much of this string you can encode in the DNS and how much must be managed by the application is part of the problem here. Is the effort to internationalize the DNS with multiple languages restricted to the “words” of the DNS, leaving the implicit left-to-right ordering and the punctuation of the DNS unaltered? If so, how much of this ordering and punctuation is a poor compromise, in that these DNS conventions in such languages are not natural translations?

The Unicode UTF-8, UTF-16, and UTF-32 encodings all require an “8-bit clean” storage and transmission medium. Because “traditional” DNS domain names are representable with 7-bit ASCII characters, not all applications that process domain names preserve the status of the eighth bit; in other words, they are not 8-bit clean. This situation stimulated significant debate in the IETF’s *IDN Working Group* and influenced the direction of the standards development into the area of application assistance: the group took a very conservative view of the capabilities of the DNS as a restricted ASCII code application.

Accordingly, we now see the DNS itself as a heavily restricted “language.” The prudent use of the DNS specifies, in RFC 1035^[5], a sequence of “words” (or “labels”), where each label conforms to the “Letter, Digit, Hyphen” (LDH) restriction. Each DNS label must begin with a letter, restricted to the Latin character subset of “A” through “Z” and “a” through “z”, followed by a sequence of letters, digits, or hyphens, with a trailing letter or digit, and no trailing hyphen. Furthermore, the case of the letter is not important to the DNS, so, within the DNS “a” is equivalent to “A”, and so on, and all characters are encoded in monospace ASCII. The DNS uses a left-to-right ordering of these labels, with the ASCII period as the label delimiter. This restriction is often referred to as the *LDH Convention*.

The challenge posed with the effort of *internationalizing* the DNS is one of attempting to create a framework that allows Internet applications—and the DNS in particular—to be set in the user’s own language in an entirely natural fashion, and yet allow the DNS to operate in a consistent and deterministic manner within its restricted “language.” In other words, we all should be able to use browsers and e-mail systems using our own language and scripts, yet still be able to communicate naturally with others who may be using a different language interface.

The most direct way of stating the choice set of IDN design is that IDNs either change the “prudent use” of the deployed DNS into something quite different by permitting a richer character repertoire in all parts of the DNS, or IDNs change the applications that want to support a multilingual environment such that they have to perform some form of encoding transfer to map between a language string using Unicode characters and an “equivalent” string using the restricted DNS LDH character-set repertoire. It appears that options other than these two lead us into fragmented DNS roots, and having already explored that particular concept in the past, not many of us want to return to that subject. So if we want to maintain a cohesive and unified symbol space for the DNS, then either the deployed DNS has to become 8-bit clean, or applications have to do the work and present to the DNS an encoded form of the Unicode sequences that conform to the restricted DNS character repertoire.

The IDN Framework

If you are an English language user with the ASCII character set, the DNS name you enter into the browser—or the domain part of an e-mail address—is almost the same string as the string that is passed to the DNS resolver to resolve into an address (the difference is the conversion of the characters into monospace). If you want to send a mail message, you might send it to `user@example.com`, for example, and the domain name part of this address, `example.com`, is the string used to query the DNS for an *MX Resource Record* in order to establish how to actually deliver the message.

But what if you want to use a domain name that is expressed in another language? What if the e-mail address is `user@記念.com`? The problem here is that this domain name cannot be “naturally” expressed in the restricted syntax of the DNS, and although this domain name may have a perfectly reasonable Unicode code sequence, this encoded sequence is not a strict LDH sequence, nor is it case-insensitive (whatever “case” may mean in an arbitrary non-Latin script). It is here that IDNs depart from the traditional view of the DNS and use a hybrid approach to the task of mapping these language strings into network addresses.

The IDN Working Group of the IETF was formed in 2000 with the goal of developing standards to internationalize domain names. The working group’s charter was to specify a set of requirements and develop IETF standards-track protocols to allow use of a broader range of characters in domain names. The outcome of this effort was the *IDN in Applications* (IDNA) framework, published as RFCs 3454, 3490, 3491, and 3492.^[6,7,8,9]

Rather than attempting to expand the character repertoire of the DNS itself, the IDN working group used an *ASCII Compatible Encoding* (ACE) to encode the binary data of Unicode strings that would make up IDNs into an ASCII character encoding. The concept is similar to the Base64 encoding used by the *Multipurpose Internet Mail Extension* (MIME) e-mail standards, but whereas Base64 uses 64 characters from ASCII, including uppercase and lowercase, the ACE approach requires the smaller DNS-constrained LDH subset of ASCII.

The working group examined various ACE algorithms in its efforts to converge to a single standard (because different encoding algorithms have different compression goals and yields) and encode the data using slightly different subsets of ASCII. Most proposals specified a prefix to the ACE coding to tag the fact that this string was, in fact, an encoded Unicode string. The IETF adopted *punycode* as its standard IDN ACE^[9]. Punycode was chosen for its efficient encoding compression properties that produce short ACE strings. For example, the domain name of `記念.com` encodes with punycode to `xn--h7tw15g.com`.

IDN in Applications

Although an ASCII-compatible encoding of Unicode characters allows representation of an IDN in a form that will probably not be corrupted by the deployed DNS infrastructure on the Internet, an ACE alone is not a full solution. The IDN approach also needs to specify how and where the ACE should be applied.

The overall approach to IDNs is relatively straightforward. In IDN the application has a critical role to play. The application takes a domain name that is expressed in a particular language using a particular script—and potentially in a particular character and word order that is related to that language—and produces an ASCII-compatible LDH-encoded version of this DNS name. Equally, when presenting a DNS string to the user, the application should take the LDH-encoded DNS name and transform it to a presentation sequence of glyphs that correspond to the original string in the original script.

It is critical that all applications perform this encoding and decoding function correctly, deterministically, and uniformly. In fact, this capability is critical to the entire IDN framework.

The basic shift in the DNS semantics that IDNs bring to the DNS is that the actual name itself is no longer in the DNS. An encoded version of the canonical name form sits in the DNS, and applications need to perform the canonical name transformation, as well as the mapping between the Unicode character string and the encoded DNS character string. So we need to agree on what are the “canonical” forms of name strings in every language. We also need to agree on the encoding method, and our various applications must have precise equivalents of these canonical name and encoding algorithms, or the symbolic consistency of the DNS will fail. The problem here is that the DNS does not perform approximate matches or return a set of possible answers to a query. The DNS is a deterministic system that performs a precise match on the query in order to generate a response. The implication here is that if we want the same IDN character sequence to map to the same network response in all cases and all contexts, then all applications must perform precisely the same operations on the character sequence in order to generate the ACE-equivalent label sequence.

RFC 3454^[6] defines a presentation layer in IDN-aware applications that is responsible for the punycode ACE encoding and decoding. This new layer in the application architecture is responsible for encoding any internationalized input in domain names into punycode format before the corresponding LDH encoded domain name is passed to the DNS for resolution. This presentation layer is also responsible for decoding the punycode format in IDNs and rendering the appropriate glyphs for the user.

It is a matter of personal perspective whether this solution is an elegant one or it simply shifts an unresolved problem from one area of the IETF to another. The IDNA approach assumes that it is easier to upgrade applications to all behave consistently in interpreting IDNs than it is to change the underlying DNS infrastructure to be 8-bit clean in a manner that would support direct use of Unicode code points in the DNS.

The Presentation Layer Transform for IDNs

The objective here is to define a reliable and deterministic algorithm that takes a Unicode string in a given language and produces a DNS string as expressed in the LDH character repertoire. This algorithm should not provide a unique 1:1 mapping, but should group “equivalent” Unicode strings, where “equivalence” is defined in the context of the language of use, into the same DNS LDH string. Any reverse mapping from the DNS LDH string into the Unicode string should deterministically select the single “canonical” string from the group of possible IDN strings.

Stringprep

The first part of the presentation layer transform is to take the original Unicode string and apply numerous transformations to it to produce a “regular” or “canonical” form of the IDN string. This form of the string is then transformed using the punycode ACE into an encoded DNS string form. The generic name of this process is, in IDN language, “stringprep,”^[6] and the particular profile of transformations used in IDNAs is termed “nameprep.”^[8]

This transform of a Unicode string into a canonical format is based on the observation that many languages have a variety of ways to display the same text and a variety of ways to enter the same text. Although we humans are unconcerned about this concept of expressing an idea in multiple ways, the DNS is an exact equivalence match operation and it cannot tolerate imprecision. So how can the DNS tell that two text strings are intended to be identical, even though their Unicode strings are different? The IDN approach is to transform the string so that all equivalent strings are mapped to the same canonical form, or “stringprep” the string. The stringprep specification is not a complete algorithm, and it requires a “profile” that describes the applicability of the profile, the character repertoire (at the time of writing RFC 3454, it was Unicode 3.2, although the Unicode Consortium has subsequently released Unicode Version 4.0, 4.1, and 5.0), mapping tables normalization, and prohibited output characters.

Mapping

In converting from a string to a *normal*, or canonical, form, the first step is to map each character into its *normalized* equivalent, using a mapping table. This table is conventionally used to map characters to their lowercase equivalent value to ensure that the DNS string comparison is case-insensitive.

Other characters are removed from the string by using this mapping operation because their presence or absence in the string does not affect the outcome of a string-equivalence operation, such as characters that affect glyph choice and placement, but without semantic meaning.

The mapping function will create monospace (specifically lowercase) outcomes and also will eliminate non-significant code points (such as, for example, the Unicode code point 1806; MONGOLIAN TODO SOFT HYPHEN or the Unicode code point 200B; ZERO WIDTH SPACE, if you really wanted to know what a non-significant code point was).

Normalization

Numerous languages use different character sequences for the same meaning. Characters may appear the *same* in presentation format as a glyph sequence, yet have *different* underlying code points. This may be associated with variable ways of combining diacritics, or using canonical code points, or using compatibility characters, and, in some language contexts, performing character reordering. For example, the character Ä can be represented by a single Unicode code point 00C4; LATIN CAPITAL A WITH DIAERESIS. Another valid representation of this character is the code point 0041; LATIN CAPITAL LETTER A followed by the separate code point 0398; COMBINING DIAERESIS.

The intent of normalization is to ensure that every class of character sequences that are equivalent in the context of a language is translated into a single canonical, consistent format. This consistency of format allows the equivalence operator to perform at the character level using direct comparison without additional language-dependent equivalence operations.

Languages in daily use are not rigid structures, and human use patterns of languages change. Normalization is no more than a best-effort process to detect equivalences in a rigid, rule-managed manner, and it may not always produce predictable outcomes. This unpredictability can be a problem with regard to namespace collisions in the DNS, because it does not increase the confidence level of the DNS as a deterministic exact-match information-retrieval system. IDNs introduce some forms of name approximation into the DNS environment, and the DNS is extremely ill-suited to the related “fuzzy-search” techniques that accompany such approximations.

Filtering Prohibited Characters

The last phase in string preparation is removal of prohibited characters, including the various Unicode white-space code points, control code points and joiners, private-use code points, and other code points used as surrogates or tags.

Right-to-Left Characters

As an option for a particular stringprep profile, you can perform a check for right-to-left displayed characters, and if any are found, make sure that the whole string satisfies the requirements for bidirectional strings. The Unicode standard has an extensive discussion of how to reorder glyphs for display when dealing with bidirectional text such as Arabic or Hebrew. All Unicode text is stored in logical order as distinct from the display order.

Nameprep: A Stringprep Profile for the DNS

The nameprep profile^[8] specifies stringprep for internationalized domain names, specifying a character repertoire (in this case the specification references Unicode 3.2) and a profile of mappings, normalization (form “KC”), prohibited characters, and bidirectional character handling. The outcome is that two-character sequences can be considered equivalent in the context of IDNs if, by following the sequences of operations defined by the nameprep profile, the resultant sequences of Unicode code points are identical. These code point sequences are the “canonical” forms of names that the DNS uses.

The Punycode ASCII-Compatible Encoding

The next step in the processing of IDN names by the application is to transform this canonical form of the Unicode name string into a LDH-equivalent string using an ACE. The algorithm used, *punycode*, uses a highly efficient encoding, attempting to limit the extent to which Unicode sequences become extended-length ACE strings.

The algorithm first divides the input code points into a set of “basic” code points that require no further encoding, and the set of “extended” code points. The algorithm takes the basic code points and reproduces this sequence in the encoded string: the “literal portion” of the string. A delimiter is then added to the string. This delimiter is a basic code point that does not occur in the remainder of the string. The extended code points are then added to the string as a series of integers expressed through an encoding into the basic (LDH) code set.

These additions of the extended code points are done primarily in the order of their Unicode values, and secondarily in the order in which they occur in the string. The encoding of the code point and its insertion position is done by using a difference, or offset, encoding, so that sequences of clustered code points, such as would be found in a single language, encode efficiently.

For example, the German language string *bücher* uses basic codes for all characters except the *ü* character. The punycode algorithm copies all the basic codes, followed by a “-”. The value and position of the *ü* insertion now has to follow.

The encoded form for *ü* (code 252) is at the position between the first and second basic characters. Using the punycode^[10] algorithm gives a delta code of 745, a value that can be expressed in base 35 as $(21 \times 35) + 10$. This code point and the position information are expressed in base 35 notation as (10,22,1), or in reverse notation, with the encoding **kva**. So the punycode encoding of *bücher* is **bcher-kva**. The internationalized domain-name format prepends the string **xn--** to the punycode string, resulting in the encoded IDN domain-name form of **xn--bcher-kva**.

IDNS and Our Assumptions About the DNS

At this stage it should be evident that we have the code points for characters drawn from all languages, and the means to create canonical forms of various words and express them in an encoded form that the DNS can resolve.

However, there is more to IDNs than the encoding algorithm. Although a massive number of discrete code points exist in the realm of Unicode, all these distinct characters are not necessarily displayed in unique ways. Indeed, given a relatively finite range of glyphs, the same glyph can display numerous discrete code points.

The often-quoted example with IDNs and name confusion is the name **paypal**. What is the difference between **www.paypal.com** and **www.paypal.com**? There is a subtle difference in the first “a” character, where the second domain name has replaced the Latin *a* with the Cyrillic *a*. Did you spot the difference? Of course not. These *homoglyphs* are cases where the underlying domain names are distinct, yet their appearance is indistinguishable. In the first case the domain name **www.paypal.com** is resolved in the DNS with the query string **www.paypal.com**, yet in the second case the query string **www.paypal.com** is translated by the application to the DNS query string **www.xn--pypal-4ve.com**. How can you tell one case from the other?

This example is by no means a unique case in the IDN realm. The reports “Unicode Security Considerations” (Unicode Technical Report 36) and “Unicode Security Mechanisms” (Unicode Technical Report 39) provide many more examples of postnormalization homographs.

There is no clear and unique relationship between characters and glyphs. Cyrillic, Latin, and Greek share numerous common glyphs. Glyphs may change their shape depending on the character sequence, multiple characters may produce a single glyph, such as the character pair *fl* being displayed as the single glyph *fl*, and a single character may generate multiple glyphs.

Homoglyphs extend beyond a conventional set of characters and include syntax elements as well. For example, the Unicode point 0244 FRACTION SLASH is often displayed using the slash glyph, allowing URLs of the form `http://a.com/e.com`. Despite its appearance, this is not a reference to `a.com` with a locator suffix of `e.com`, but is a reference to the domain `a.com/e.com`.

The basic response is that if you maintain IDN integrity at the application level, then the user just cannot tell. The punycode transform of `www.paypal.com` into `www.xn--paypal-4ve.com` is intended to be a secret between the application and the DNS, because this ASCII-encoded form is simply meaningless to the user. But if this encoded form remains invisible to the user, how can the user detect that the two identically presented name strings are indeed different? Sadly, the only true “security” we have in the DNS is the “look” of the DNS name that is presented to the user, and the user typically works on the principle that if the presented DNS string looks like the real thing, then it must be the real thing.

When this homoglyph problem was first exposed, the response from many browser implementations was to turn off all IDN support in their browser. The next response was to deliberately expose the punycode version of the URL in the browser address bar, so that directing the browser to `http://www.paypal.com` would display in the address bar the URL value of `http://www.xn--paypal-4ve.com`.

The distinction between the two equivalently displayed names was then visible to the user, but the downside was that we were back to displaying ASCII names again, and in this case ASCII versions of punycode-encoded names. If trying to “read” Base64 was difficult, then the displaying—and understanding—of displayed punycode names is surely equally as difficult, if not more so. The encoded names can be completely devoid of any form of useful association or meaning. Although the distinction between ASCII and Cyrillic may be evident by overt differences in their ASCII-encoded names, what happens when the homoglyph occurs across two non-Latin languages? The punycode strings are different, but which string is the “intended” one? Did you mean `http://xn--21bm41.com` or `http://xn--q2buub.com` when you enter a Hindi script URL?

Using ASCII as the fall-back to resolve name confusion in response to the problem of ambiguities in non-ASCII script names appears to be a nonsensical solution. We appear to be back to guessing games in the DNS again, unfortunately, and particularly impossible guessing games at that.

These days most popular browsers display the glyphs, rather than the ASCII punycode, but once more we are back to the homoglyph problem.

If the intention in the IDN effort was to preserve the deterministic property of DNS resolution, such that a DNS query can be phrased deterministically and not have the query degenerate into a search term or require the application of fuzzy logic to complete the query, then we are not quite there yet.

The underlying observation is that languages are indeed human-use systems. They can be tricky, and they invariably use what appear to be rules in strange and inconsistent ways. They are also resistant to automated processing and the application of rigid rule sets. The canonical name forms that are produced by nameprep-like procedures are not comprehensive, nor does it appear that such a rigidly defined rule-driven system can produce the desired outcomes in all possible linguistic situations. And if the intention of the IDN effort was to create a completely “natural” environment using a language environment other than English and a display environment that is not reliant on ASCII and ASCII glyphs, while preserving all the other properties of the DNS, then the outcome does not appear to match our original IDN expectations.

The underlying weakness here is the implicit assumption that in the DNS “what you see is what you get,” and that two DNS names that look identical are indeed references to the same name, and when resolved in the DNS produce precisely the same resolution outcome. When you broaden the repertoire of appearances of the DNS, such that the entire set of glyphs can be used in the DNS, then the mapping from glyph to underlying code point is not unique. Any effort to undertake such a mapping needs additional context in the form of a language and script context. But the DNS does not carry such a context, making the task of maintaining uniqueness and determinism of DNS name translation essentially impossible if we also want to maintain the property that it is the appearance, or presentation format, of DNS names to the user that is the foundation stone of the integrity of our trust in the DNS.

Some concerns still remain in this space, including the inclusion of various forms of character codes that are in effect invisible. In addition, homoglyphs could be better managed by using a refined definition of IDN labels that lists which Unicode code points can be used in the context of IDNs, excluding all others. It would be helpful if confusing and non-reversible character mappings were removed from the IDN space, including the consistent treatment of ligatures and diacritics, refining the treatment of right-to-left and left-to-right scripts, and removing the dependency on a particular version of the Unicode standard. This effort is under way in the IETF in the context of revisions to the IDNA specification documents.

IDNs, TLDs, and the Politics of the DNS

So why is there a very active debate, particularly within ICANN-related forums, about putting IDN codes into the root of the DNS as alternative *top-level domains* (TLDs)?

I have seen two major lines of argument here; namely the argument that favors the existence of IDNs in all parts of the DNS, including the TLDs, and the argument that favors a more restricted view of IDNs in the root of the DNS that links their use to that of an existing (ASCII-based) DNS label in the TLD zone.

Apparently, those who favor the approach of using IDNs in the top-level zone as just another DNS label see this as a natural extension of adding punycode-encoded name entries into lower levels of the DNS. Why should the root of the DNS be any different, in terms of allowing IDNs? Why should a non-Latin script user of the Internet have to enter the TLD code in its ASCII text form, while entering the remainder of the string in a local language? And in right-to-left scripts, where does this awkward ASCII appendage sit when a user attempts to enter it into an application?

Surely, goes the argument, the more natural approach is to allow any DNS name to be wholly expressible in the user's language, implying that all parts of the DNS should be able to carry native language-encoded DNS names. After all, コンピュータは予約する.jp looks wrong as a monolingual domain name. What is that .jp appendage doing there in that DNS name? Surely a Japanese user should not have to resort to an ASCII English abbreviation to enter in the country code for Japan, when 日本 is obviously more “natural” in the context of a Japanese user using Japanese script. If we had punycode TLDs then, goes the line of argument, users could enter the entire domain name in their language and have the punycode encoding happen across the entire name string, and then successfully perform a DNS lookup on the punycode equivalent. This way the user would enter the Japanese character sequence: コンピュータは予約する.日本 and have the application translate this entry to the DNS string **xn--88j0bve5g9-bxg1ewerdw490b930f.xn--wgv71a**. For this process to work in its entirety uniformly and consistently, the name **xn--wgv71a** needs to be a TLD name.

We can always take this thought process one step further and question the ASCII string **http** and the punctuation symbols **://** for precisely the same reason, but I have not heard (yet) calls for multilingual equivalents of protocol identifier codes. The multilingual presentation of these elements remains firmly in the provenance of the application, rather than attempting to alter the protocol identifiers in the relevant standards.

The line of argument also encompasses the implicit threat that if the root of the DNS does not embrace TLDs as expressed in the language of the Internet's users, then language communities will break away from a single DNS root and meet their linguistic community's requirements in their own DNS hierarchy. Admitting such encoded tags into the DNS root is the least problematic, including the consequence of inactivity, which is cited as being tantamount to condoning the complete fragmentation of the Internet's symbol set.

Of course having an entirely new TLD name in an IDN name format does not solve all of the potential problems with IDNs. How can a user tell what domain names are in the ASCII top level, and what are in the "equivalent" IDN-encoded TLDs? Are any two name spaces that refer to the same underlying name concept equivalent? Is **xn--88j0bve5g9bxg1ewerdw490b930f** appropriately a subdomain of **.jp**, or a subdomain of **xn--wgv71a**? Should the two domains be tightly synchronized with respect to their zone content and represent the same underlying token set, or should they be independent offerings to the marketplace, and allow registrants and the end-user base make implicit choices here? In other words, should the pair of domain names, namely **xn--88j0bve5g9bxg1ewerdw490b930f.xn--wgv71a** and **xn--88j0bve5g9bxg1ewerdw490b930f.jp**, reference precisely the same DNS zone, or should they be allowed to compete, and each find their own "natural" level of market support based on decoupled TLD names of **.jp** and **.xn--wgv71a**?

What does the term *equivalence* really imply here? Is equivalence something as loose as the relationship between **.com** and **.biz**, namely being different abbreviations of words that reflect similar concepts with different name-space populations that reflect market diversity and a competitive supply industry? Or is equivalence a much tighter binding in that equivalent names share precisely the same subdomain name set, and a registration in one of these equivalence names is in effect a name registration across the entire equivalence set?

Even this subject is not readily resolvable given our various interpretations of *equivalence*. In theory, the DNS root zone is populated by ISO two-letter country codes and numerous "generic" TLDs. Under what basis, and under what authority, is **xn--wgv71a** considered an "equivalent" of the ISO 3166 two-letter country code JP? Are we falling into the trap once again of making up the rules as we go along? Is the distinction between **.com** and **.biz** apparent only in English? And why should this distinction apply only to non-Latin character sets? Surely it makes more sense for a native German language speaker to refer to commercial entities as *kommerze*, and the abbreviated TLD name as **.kom**? When we say "multilingual" are we in fact ignoring "multilingual" and looking exclusively at "multiscript"?

Let's put aside the somewhat difficult concept of name equivalence for a second, and assume that this equivalence problem is solved. Also suppose that we want tight coupling across equivalence sets of names.

In other words, what we want is that a name registered in any of the elements of the equivalent domain-name set in all scripts is, in effect, registered in all the equivalent DNS zones. The question is: how should it be implemented in the DNS? One approach that could support tight synchronization of equivalence is to use the DNAME record^[11] to create these TLD name aliases for their ASCII equivalents, thereby allowing a single name registration to be resolvable using a root name expressed in any of the linguistic equivalents of the original TLD name. The DNAME entry for all but the “canonical” element of the equivalence set effectively translates all queries to a query on the canonical name. The positive aspects of such an approach is uniformity across linguistic equivalents of the TLD name form—a single name delegation in a TLD domain becomes a name within all the linguistic equivalents of the TLD name without any further delegation or registration required.

Using DNAME as a tool to support sets of equivalent names in the DNS is still in the early stages. The limited experience so far with DNAME indicates that CNAME synthesis places load back on the name servers that would otherwise not be there, and the combination of this synthetic record and DNSSEC starts to get very unwieldy. Also, the IETF is reviewing the DNAME specification with the intention to remove the requirement to perform CNAME synthesis. All of these factors may explain why there is no immediate desire to place DNAMEs in the DNS root zone.

Different interpretations of equivalence in IDN names are possible. The use of DNAMEs as aliases for existing TLDs in effect “locks up” IDNs into the hands of the incumbent TLD name-registry operators. Part of the IDN debate, is, as usual, a debate over the generic TLD registry operators and the associated perception of incumbent monopolies. An alternative approach is to associate a single registrar with each IDN variant of the same generic TLD, allowing a form of “competition” between the various registrars. From the perspective of a coherent symbol space where the same symbol, expressed in any language script, resolves in the same fashion, such independent registries are not overly consistent with such a model of registry diversity in a multilingual environment. In this case such an artifice of IDN “competition” may well do more harm than good for Internet users.

It appears that another line of argument is that the DNS top-level name space is very conservatively managed, and new entries into this space are not made lightly. There are concerns of stability of operation, of attempting to conserve a coherent namespace, and the ever-present consideration that if we manage to “break” the DNS root zone it would be an irrevocable act.

This line of argument recognizes the very hazy nature of name equivalence in a multilingual environment and is based on the proposition that the DNS is incapable of representing such imprecision with any utility. The DNS is not a search engine, and the DNS does not handle imprecision at all well. Again, goes the argument, if this is the case then can we push this problem back to the application rather than trying to bend the DNS? If an application is capable of translating, say, 日本 into `xn--wgv71a`, and considering that the TLD name space is relatively small, it appears that having the application performing a further translation of this intermediate form punycode string into the ASCII string `jp` is not a particularly challenging form of table lookup. In such a model no new TLD aliases or equivalences are required in the root zone of the DNS. If we are prepared to pass the execution of the presentation layer of the DNS to the application layer to perform, then why not also ask this same presentation layer to perform the step of further mapping the punycode ACE equivalents of the TLDs to the actual ASCII TLDs, using some richer language context that the application may be aware of that is not viable strictly within the confines of the DNS?

So, with respect to the question of whether IDN TLDs should be loaded into the DNS at all, and, if so, whether they should represent an opportunity for further diversity in name supply or be constrained to be aligned to existing names, and precisely how name equivalence is to be interpreted in this context, then it appears that ICANN has managed to place itself in a challenging situation. In not making a decision, those with an interest in having diverse IDN TLDs appear to derive some pleasure in pointing out that the political origins of ICANN and its strong linguistic bias to English are influencing it to ignore non-English language use and non-English language users of the Internet. Where dramatic statements are called for, such statements often use terms such as “cultural imperialism” to illustrate the nature of the linguistic insult. The case has been made repeatedly, in support of IDN TLDs, that an overwhelming majority of Internet users and commercial activity of the Internet is in languages other than native English, and the imposition of ASCII labels on the DNS is an unnatural imposition on the overwhelming majority of Internet users.

On the other hand, most decisions to permit some form of entry in the DNS are generally seen as irrevocable, and building a DNS that is littered with the legacy of various non-enduring name technologies and poor ad hoc decisions to address a particular concern or problem without any context of a longer-term framework seems also to represent a step along a direction leading to a heavily littered and fragmented Internet where, ultimately, users cannot communicate with each other.

What about global interoperability and the Internet? Should we just take the easy answer and simply give up on the entire concept? Well of course not! But, taking a narrower perspective, are IDNs simply not viable in the DNS? I would suggest that not only is this question one that was overtaken by events years ago, but even if we want to reconsider it now, then the answer remains that any users using their local language and local script should have an equally “natural” experience. IDNs are a necessary and valuable component of the symbol space of any global communications system, and the Internet is no exception. However, we also should recognize that we do need combinations of both localization and globalization, and that we are voicing some pretty tough objectives. Is the IDNA approach enough? Is our assumption that an unaltered DNS with application-encoded name strings represents a rich enough platform to preserve the essential properties of the DNS while allowing true multilingual use of the DNS? On the other hand, taking a pragmatic view of the topic, is what we have with IDNA enough for us to work on, and is the alternative of reengineering the entire fabric of the DNS into an 8-bit clean system just not a viable option?

I suspect that the framework of IDNA is now the technology for IDNs for the Internet, and we simply have to move on from here and deliberately take the stance of understanding the space from users’ perspectives when we look at the policy concerns of IDNs. The salient questions from such perspectives include: “What is the “natural” thing to do?” and “What causes a user the least amount of surprise?” Because in this world, what works for the user is what works for the Internet as a whole.

Further IDN News

IDNs are by no means completed work. Development continues in the Unicode forum on elaboration of character sets, and there are further proposals in the IETF to continue a complementary standards activity of refining the IDN documents.

In February 2008 the *Applications Area* of the IETF announced a proposal for further work on IDNs. The proposal has noted that the existing RFC documents are tied to version 3.2 of Unicode, while the Unicode Consortium has released version 5.0.0.

The proposed work is to consider revision of the IDN documents to untie the Internet specifications that define validity based on Unicode properties from specific versions of Unicode using algorithms. It is also proposed that these updates study revision of bi-directional algorithms, and to permit the use of some scripts that were inadvertently excluded by the original Internet specification.

This is not intended to be a major rewrite of the IDN approach, and, in particular, IDNs will continue to use the **xn-** prefix, the same Punycode ASCII-compatible encoding, and the bidirectional algorithm is intended to follow the same design as presently specified.

Further Reading

It is possible to reference an overwhelming amount of commentary on this topic, so I have deliberately kept this list of further reading on the topic of IDNs relatively brief:

- [A] John Klensin, “Internationalizing Top-Level Domain Names: Another Look,” ISOC Member Briefing, September 2004, <http://www.isoc.org/briefings/018/>
- [B] John Klensin, “National and Local Characters for DNS Top Level Domain (TLD) Names,” RFC 4185, October 2005.
- [C] Papers submitted to the ICANN IDN TLD workshop, held in November 2005: <http://www.icann.org/announcements/announcement-17nov05.htm>
- [D] Internet Architecture Board, “Review and Recommendations for Internationalized Domain Names (IDNs),” RFC 4690, September 2006.
- [E] “ICANN’s IDN Roadmap Announcement—Progress and Future,” <http://www.icann.org/announcements/announcement-1-01nov06.htm>
- [F] “An Important Step Toward the Implementation of IDN Top-Level Domains: New Versions of IDNA Protocol Revision Proposals Posted,” <http://www.icann.org/announcements/announcement-26nov07.htm>
- [G] ICANN’s IDN Evaluation Gateway. Eleven new internationalized domains representing the name **example.test** entirely in scripts other than the Latin characters: <http://idn.icann.org/>

References

- [1] http://en.wikipedia.org/wiki/Horizontal_and_vertical_writing_in_East_Asian_scripts
- [2] http://en.wikipedia.org/wiki/Roman_script
- [3] <http://unicode.org>
- [4] <http://www.omniglot.com/writing/thai.htm>
- [5] Mockapetris, P., “Domain Names—Implementation and Specification,” RFC 1035, November 1987.
- [6] Hoffman, P., and Blanchet, M., “Preparation of Internationalized Strings (“stringprep”),” RFC 3454, December 2002.
- [7] Hoffman, P., Fältström, P., and Costello, A., “Internationalizing Domain Names in Applications (IDNA),” RFC 3490, March 2003.
- [8] Hoffman, P., and Blanchet, M., “Nameprep: A Stringprep Profile for Internationalized Domain Names (IDN),” RFC 3491, March 2003.
- [9] Costello, A., “Punycode: A Bootstring encoding of Unicode for Internationalized Domain Names in Applications (IDNA),” RFC 3492, March 2003.
- [10] <http://en.wikipedia.org/wiki/Punycode>
- [11] Crawford, M., “Non-Terminal DNS Name Redirection,” RFC 2672, August 1999.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. The author of numerous Internet-related books, he is currently the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He was a member of the Internet Architecture Board from 1999 until 2005, and served on the Board of the Internet Society from 1992 until 2001. E-mail: gih@apnic.net

The Locator Identifier Separation Protocol (LISP)

by David Meyer, Cisco Systems

The Internet Architecture Board's (IAB)'s October 2006 *Routing and Addressing Workshop*^[8] renewed interest in the design of a scalable routing and addressing architecture for the Internet. Many concerns prompted this renewed interest, including the scalability of the routing system and the impending exhaustion of the IPv4 address space. Since the IAB workshop, several proposals have emerged that attempt to address the concerns expressed both at the workshop and in other forums^[7,9,12,13,14]. All of these proposals are based on a common concept: the separation of locator and identifier in the numbering of Internet devices, often termed the "Loc/ID split." This article focuses on one proposal for implementing this concept: the *Locator/Identifier Separation Protocol* (LISP)^[3].

The basic idea behind the Loc/ID split is that the current Internet routing and addressing architecture combines two functions: *Routing Locators* (RLOCs), which describe how a device is attached to the network, and *Endpoint Identifiers* (EIDs), which define "who" the device is, in a single numbering space, the IP address. Proponents of the Loc/ID split argue that this "overloading" of functions makes it virtually impossible to build an efficient routing system without forcing unacceptable constraints on end-system use of addresses. Splitting these functions apart by using different numbering spaces for EIDs and RLOCs yields several advantages, including improved scalability of the routing system through greater aggregation of RLOCs. To achieve this aggregation, we must allocate RLOCs in a way that is congruent with the topology of the network ("Rekhter's Law"). Today's "provider-allocated" IP address space is an example of such an allocation scheme. EIDs, on the other hand, are typically allocated along organizational boundaries. Because the network topology and organizational hierarchies are rarely congruent, it is difficult (if not impossible) to make a single numbering space efficiently serve both purposes without imposing unacceptable constraints (such as requiring renumbering upon provider changes) on the use of that space.

LISP, as a specific instance of the Loc/ID split, aims to decouple location and identity. This decoupling will facilitate improved aggregation of the RLOC space, implement persistent identity in the EID space, and, in some cases, increase the security and efficiency of network mobility.

Implementing the Locator/ID Separation

There are two basic approaches to implementing the Loc/ID split: *map-and-encap* and *address rewriting*. Each is briefly discussed in the following sections.

Map-and-encap

In the map-and-encap scheme (generally considered to have evolved from Bob Hinden's ENCAPS protocol^[24]), when a source sends a packet to the EID of a destination outside of the source domain, the packet traverses the domain infrastructure to a border router (or other border element). The border router maps the destination EID to a RLOC that corresponds to an entry point in the destination domain (hence an EID-to-RLOC mapping system is needed; proposals are discussed later in the article). This phase is the “map” phase of map-and-encap. The border router then encapsulates the packet and sets the destination address to the RLOC returned by the mapping infrastructure (if any; it may be statically configured as well). This phase is the “encap” phase of the map-and-encap model.

Thus map-and-encap works by appending a new header to the existing packet; the “inner-header” source and destination addresses are EIDs, and the “outer-header” source and destination addresses are in most cases RLOCs. When an encapsulated packet arrives at the destination border router, the router decapsulates the packet and sends it on to its destination. Note that this process suggests that EIDs may need to be routable in some scope (likely scoped to the domain).

Map-and-encap schemes have the desirable property that they do not in general require host changes or changes to the core routing infrastructure. In addition, map-and-encap schemes work with both IPv4 and IPv6, and retain the original source address (a feature that is useful in various filtering scenarios). Controversy remains, however, as to whether or not the encapsulation overhead of map-and-encap schemes is problematic; opinions exist on both sides of this topic (see, for example, [18]).

Address Rewriting

The basic idea behind the address-rewriting schemes, originally proposed by Dave Clark and later by Mike O'Dell in his 8+8/GSE specification^[11], is to take advantage of the 128-bit IPv6 address and use the top 64 bits as the routing locator (“Routing Goop,” or RG), and the lower 64 bits as the endpoint identifier (hence rewriting works only for IPv6). In this scheme, when a host emits a packet destined for another domain, the source address contains its identifier (frequently a IEEE MAC address) in the lower 64 bits, and a special value (meaning unspecified) in the RG. The destination address contains the fully specified destination address (RG and EID).

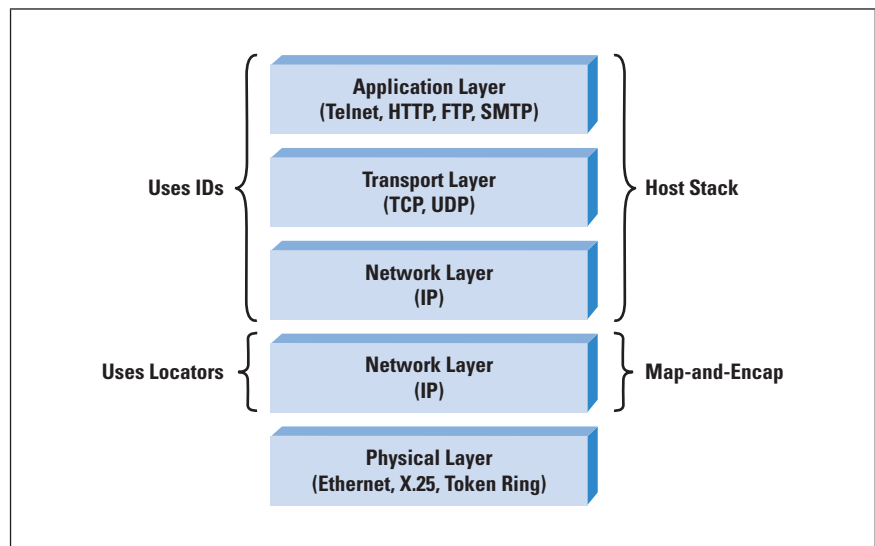
When a packet destined for a remote domain arrives at the local domain egress router, the source RG is filled in (forming a full 128-bit address), and the packet is routed to the remote domain. On ingress to the remote domain, the destination RG is rewritten with the unspecified value, ensuring that the host does not know what its RG is.

This process, in theory, would enable the ease of renumbering that would be required to maintain congruence between prefix assignment and physical network topology that is required for the kind of “aggressive” renumbering envisioned in the 8+8/GSE specification.

The Locator/Identifier Separation Protocol (LISP)

LISP is designed to be a simple, incremental, network-based map-and-encap protocol that implements separation of Internet addresses into EIDs and RLOCs. Because LISP is a map-and-encap protocol, it requires no changes to host stacks and no major changes to existing database infrastructures. It is designed to be implemented in a relatively small number of routers. LISP is also an instance of what is architecturally called a “jack-up,” because the existing network layer is “jacked up” and a new network layer is inserted below it (the term “jacked up” is attributed to Noel Chiappa). The LISP jack-up is depicted in Figure 1.

Figure 1: LISP is a Jack-Up



The LISP design aims to improve site multihoming (for example, by controlling site ingress without complex protocols), improve *Internet Service Provider* (ISP) multihoming, decouple site addressing from provider addressing, and reduce the size and dynamic properties of the core routing tables.

The LISP data plane (the map-and-encap operation) and the LISP control plane (the EID-to-RLOC mapping system) are very modular. In particular, although the base LISP specification defines the format of messages to query the mapping system and to receive responses from that system, it makes no assumptions on the architecture of potential mapping systems. As a result, several mapping systems have been proposed^[0,1,4,5,6,10].

LISP Network Elements

The LISP specification defines two network elements: The Egress Tunnel Router (ETR) and the Ingress Tunnel Router (ITR).

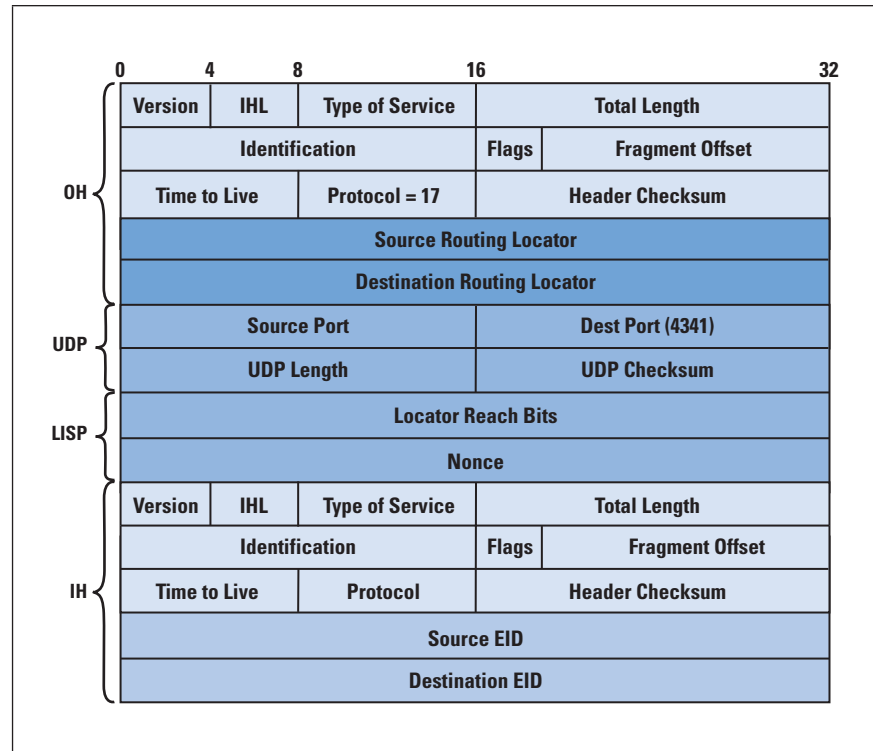
A LISP *Egress Tunnel Router* (ETR) receives LISP-encapsulated IP packets from the Internet on one side and sends decapsulated IP packets to site end systems on the other side. In particular, an ETR accepts an IP packet where the destination address in the “outer” IP header is one of its own RLOCs. The router strips the “outer” header and forwards the packet based on the next IP header found.

A LISP *Ingress Tunnel Router* (ITR) accepts IP packets from site end systems on one side and sends LISP-encapsulated IP packets toward the Internet on the other side. In particular, an ITR accepts an IP packet with a single IP header (more precisely, an IP packet that does not contain a LISP header). The router treats this “inner” IP destination address as an EID and performs an EID-to-RLOC mapping lookup if necessary (that is, it does not already have an EID-to-RLOC mapping for the EID). The router then prepends an “outer” IP header with one of its globally routable RLOCs in the Source Address field and the result of the mapping lookup in the Destination Address field. Note that this destination RLOC may be an intermediate, proxy device that has better knowledge of the EID-to-RLOC mapping closest to the destination EID.

LISP Data-Plane Operation

When a host in a LISP-capable domain emits a packet, it puts its EID in the packet source address, and EID of the correspondent host in its destination address (note that hosts will typically look up EIDs in the *Domain Name System* [DNS]). If the destination of the packet is in another domain, the packet traverses the source domain infrastructure to one of its ITRs. The ITR maps destination EID to a RLOC that corresponds to an ETR that is either in the destination domain or a proxy for the destination domain (how this mapping is accomplished in LISP is discussed later in the article). The ITR then encapsulates the packet, setting the destination address to the RLOC of the ETR returned by the mapping infrastructure or by static configuration. Note that LISP is address family-agnostic and as such can be used with both IPv4 and IPv6 (or any other address family). Figure 2 depicts the LISP IPv4 in IPv4 encapsulation.

Figure 2: LISP Header Format



When the packet arrives at the destination ETR, it decapsulates the packet and sends it on to its destination. Again, note that this scenario implies that EIDs need to be routable in some scope (likely scoped to the domain).

As mentioned previously, the LISP specification defines three packet types designed to support an EID-to-RLOC mapping system. The first type of packet, the *Data Probe*, is a data packet that an ITR may send into the mapping system to probe for the mapping; the authoritative ETR responds to the ITR with a Map-Reply message when it receives such a data packet. Note that in this case the ETR detects that the packet is a Data Probe by noticing that the inner *Destination Address* (DA) was copied to the outer DA by the ITR, that is, the inner DA equals the outer DA and is an EID. The second type of LISP packet used to support the mapping system is the *Map Request*. An ITR may query the mapping system by sending a Map-Request message into the mapping system to request a particular EID-to-RLOC mapping. As in the Data Probe case, the authoritative ETR responds with a Map-Reply message.

The third type of LISP packet used to support the mapping system is the *Map Reply*. An ETR emits a Map Reply under two conditions. First, if the ETR receives a LISP-encapsulated packet in which the outer-header destination address is the same as that of the inner header, it knows that the packet is a Data Probe and can respond with a Map Reply to the source ITR. The ETR may also receive a Map Request, in which case it replies to the requesting ITR with the mapping.

LISP Control Plane

Both map-and-encap and address-rewriting models rely on an additional level of indirection in the addressing architecture to make the routing system scale reasonably. Because packets are sourced with an EID in the Destination Address field and EIDs are not in general routable on the global Internet, the destination EID must be mapped to an RLOC in order to deliver the packet to another domain (that is, across the Internet). In the case of the map-and-encap schemes, it is a direct translation: an EID is mapped to a RLOC. The situation is subtly different for the rewriting schemes; in general such schemes must look up the entire destination address (usually proposed to reside in the DNS)^[11,13], but must somehow determine the source RG when rewriting the source address at the domain border.

In either Loc/ID split model, an EID-to-RLOC mapping service is needed to make the system scale reasonably and to make it operationally viable. There are three important scale parameters to consider when architecting a mapping service: the rate of updates to the mapping database, the state of the mapping service required, and the latency incurred during database lookup. The scaling properties of the database are frequently characterized as a $(Rate \times State)$ problem (ignoring for the moment the subject of lookup latency); because most estimates put the size of the mapping database at $O(10^{10})$, the database update rate must be small (note that this situation is a primary reason that current mapping proposals do not incorporate reachability information into the mapping database). In addition, the choice of push vs. pull also affects latency: if you push the entire database close to the edge, you improve lookup latency at the cost of increased state; if you architect a service that requires a mapping request and you find an authoritative server for that mapping (that is, pull), you reduce state at the cost of increased lookup latency.

LISP-Alternative-Topology: A LISP Control Plane

The basic idea behind *LISP-Alternative-Topology* (LISP-ALT)^[4] is to build an alternative logical topology for managing EID-to-RLOC mappings for LISP. This logical topology uses existing technology and tools, specifically the *Border Gateway Protocol* (BGP)^[17] and its multiprotocol extension^[15], along with the *Generic Routing Encapsulation* (GRE)^[16] protocol to construct an overlay network of devices that advertise EID prefixes only.

As was the case for the LISP data plane, an important design goal of LISP-ALT is to minimize the number of changes to existing hardware and software that are required to deploy the mapping system. Therefore, LISP-ALT requires modifications to neither BGP nor GRE.

Note that LISP-ALT is a hybrid push/pull architecture. Aggregated EID prefixes are “pushed” among the LISP-ALT routers and, optionally, to ITRs (which may elect to receive the aggregated information, as opposed to simply using a default mapping). Specific EID-to-RLOC mappings are “pulled” by ITRs either by Map Requests or Data Probes, both of which are routed over the alternate topology and result in Map Replies being generated by ETRs.

The basic idea behind in LISP-ALT, then, is to use BGP running over a GRE overlay to build the reachability required to route Data Probes, Map Requests, and Map Replies over the alternate topology. The *ALT Routing Information Base* (RIB) comprises EID prefixes and associated next hops. The LISP-ALT routers talk *External BGP* (eBGP) to each other in order to propagate EID prefix update information, which is learned either over eBGP connections from the authoritative ETR or by configuration. ITRs may also eBGP peer with one or more LISP-ALT routers in order to route Data Probe packets or Map Requests.

In summary, the LISP-ALT uses BGP to propagate EID-prefix reachability information used by ITRs and ETRs to forward Map Requests, Map Replies, and Data Probes. This reachability is carried as IPv4 or IPv6 *Network Layer Reachability Information* (NLRI) without modification (because the EID space has the same syntax as IPv4 or IPv6). LISP-ALT routers eBGP peer with one another, forming the overlay network. A LISP-ALT router near the edge learns EID prefixes that originate with authoritative ETRs. In general then, LISP-ALT routers aggregate EID prefixes, and forward Data Probes, Map-Requests, and Map-Replies.

Threat Models and Mitigation

As in any Loc/ID split approach, a critical operation is the creation of locator-to-ID binding state that devices will use over time. In the case of LISP, the critical operation is the creation of EID-to-RLOC mappings in the ITR and the ETR. We can obtain these mappings in three ways:

- By using the information obtained from a LISP data packet
- By using the information contained in the Map-Reply message
- By using an EID-to-RLOC mapping database

LISP mitigates attacks on the first two techniques by including a *nonce* in the LISP header; the nonce is a 32-bit randomly generated number (generated by the source ITR) that is used to test route returnability.

More specifically, an ETR echoes the nonce back to the ITR in a Map-Reply message. That is, the nonce, combined with the ITR accepting only solicited Map Replies, provides a base level of authentication for Map Replies. Note however, that these techniques do not protect against man-in-the-middle attacks.

The LISP design assumes that many (if not most) security mechanisms are part of the mapping database service when using control-plane procedures for obtaining EID-to-RLOC mappings. *Denial-of-Service* (DoS) attack prevention, on the other hand, depends on the ability of an implementation to rate-limit Map Requests and Map Replies (in the control plane), as well as its ability to rate limit the number of data-triggered Map Replies (for example, in response to Data Probe packets).

Refer to [19] for a more detailed preliminary threat analysis for LISP.

LISP and Fast Endpoint Mobility

Fast endpoint mobility occurs when an endpoint moves relatively rapidly, changing its IP layer network attachment point, and maintenance of session continuity is a goal. Mobile IPv4^[20] and Mobile IPv6^[21,22,27] mechanisms can be used in this case; note however, that the interaction of Mobile IP with LISP needs further exploration. Refer to the LISP specification^[3] for additional details.

In summary, the major problem introduced by a Loc/ID split scheme is that as an endpoint moves, changes to the mapping between its EID and a set of RLOCs for its new network location may be required. When this change is added to the overhead of mobile IP binding updates, some packets might be delayed or dropped. In general, the problem is controlling the update rate (that is, the $[Rate \times State]$ product described previously), and is an area of ongoing research.

Multicast

A multicast group address, as defined in the original Internet architecture, is an identifier of a grouping of topologically independent receiver host locations. The address encoding itself does not determine the location of the receiver(s). The multicast routing protocol and the network-based state the protocol creates determine the location of the receivers.

In the LISP context, a multicast group address is both an EID and a RLOC. As such, no specific action is necessary for destination addresses; a group address that appears in an inner IP header (built by a source host) is used as the destination EID by an ITR as a destination address when it LISP-encapsulates the packet (that is, the ITR uses the same group address as the destination RLOC).

The source RLOC, as is usually the case, is the ITR IP address (that is, one of its RLOCs).

At the receiving side, *Protocol Independent Multicast* (PIM)^[23] has to translate the source-address Join/Prune messages from RLOCs to EIDs when multicast packets are forwarded by the ETR. However, in contrast to the unicast case (where a Map Request is sent by the ITR at forwarding time), a Map Request can be sent when the multicast tree is being built.

Putting It All Together: A Day in the Life of a LISP Packet

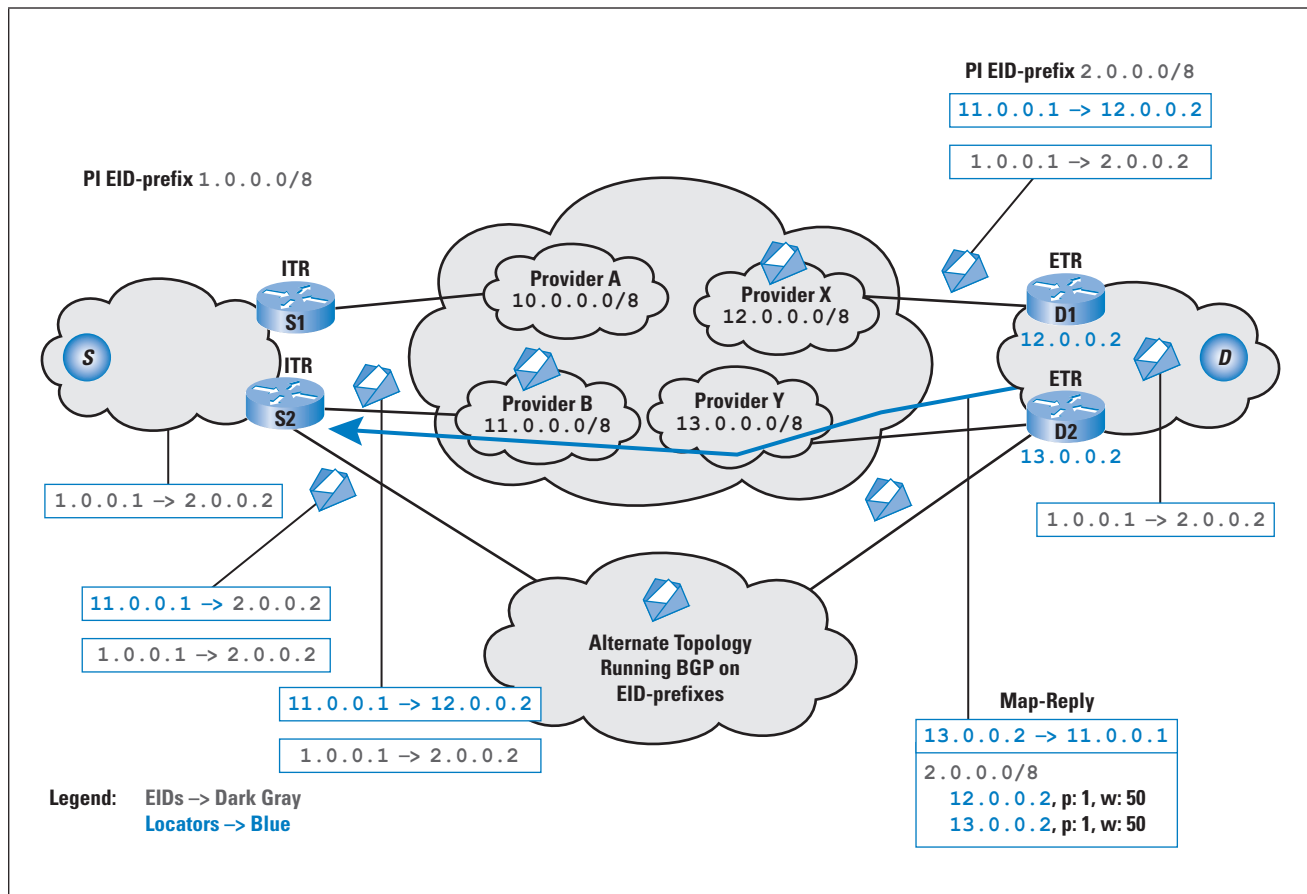
When a host in a LISP-capable domain wants to send a packet, it first looks up the correspondent host's EID in the DNS. It then puts its EID in the packet source address, and EID of the correspondent host in its destination address; if the destination of the packet is in another domain, the packet traverses the source domain infrastructure to one of the domain ITRs.

If the ITR has cached the EID-to-RLOC mapping for the destination EID, it sets the destination RLOC in the outer (encapsulated) header to the cached RLOC, and the source RLOC to its RLOC (note that the inner header has the source host's EID as the source and the destination's EID in the Destination field). The packet is then sent over the Internet to the ETR indicated in the destination RLOC, which decapsulates the packet and sends it on to the destination EID.

If, on the other hand, the ITR does not have a EID-to-RLOC mapping for the destination EID, it encapsulates the packet in a LISP header in which the destination address is the same as the inner header destination address, namely, the EID of the destination host. This packet is a Data Probe packet, and is routed over the LISP-ALT topology to the LISP-ALT router (typically an ETR, but this type of router is not required) that is authoritative for the EID-to-RLOC mapping. When the ETR receives the Data Probe packet, it decapsulates the packet and sends it on to the destination EID and sends a Map Reply to the source ITR so subsequent packets are sent natively over the Internet (as opposed to over the LISP-ALT overlay network). This query/response transaction is required only for the first packet sent between sites; all subsequent packets are sent LISP-encapsulated directly between the ITR and the ETR (and in particular, not over the LISP-ALT topology). Finally, note that the ITR could also preload its cache with mappings for popular destinations using the Map-Request message, avoiding the Data Probe packet (and associated latency, if any) altogether.

For example, consider the scenario depicted in Figure 3. In this case, a source S with EID $1.0.0.1$ wants to send a packet to destination D whose EID is $2.0.0.2$. The packet arrives at ITR S2, which does not have an EID-to-RLOC mapping for $2.0.0.2$. S2 LISP-encapsulates the packet with the outer header having its RLOC ($11.0.0.1$) as the source address, copies the destination EID ($2.0.0.2$) from the inner header to the outer-header destination, and sends the data packet (a Data Probe) into the LISP-ALT topology. The packet follows the paths computed by BGP in the LISP-ALT topology to ETR D2. When D2 receives the packet, it decapsulates it and forwards the packet to the destination $2.0.0.2$; D2 also responds with a Map-Reply message that tells S2 ($11.0.0.1$) that the EID-to-RLOC mapping for $2.0.0.0/8$ has two elements, ETR D1 (whose RLOC is $12.0.0.2$) and ETR D2 (whose RLOC is $13.0.0.2$). After receiving the Map Reply, ITR S2 can send LISP-encapsulated packets natively over the Internet (that is, not over the ALT topology).

Figure 3: A Day in the Life of a LISP Packet



Note that the mapping has priority (p) and weight (w) attributes. Priorities tell the ITR which ETRs to use in which order, and weights tell the ITR how to split load across ETRs of a given priority (w is a percentage of traffic that should go to each ETR). In this case, both ETRs have the same priority (1), and have weight 50 (that is, each ETR should receive 50 percent of the traffic).

New Functions Enabled by the Mapping System

Weights and priorities provide new capabilities for multihomed sites, which can use these features to control how traffic ingressing to the site is spread across its links without the complexity and overhead of running BGP. In particular, a multihomed site can configure its mapping database so that its links are used in an “active-active” configuration (that is, both links are in use). This situation is depicted in Figure 3, where the mapping databases entry `2.0.0.0/8` has two ETRs at the same priority that are equally weighted, meaning that the ITR will spread flows equally among the two ETRs.

This function is particularly attractive for *Small Office or Home Office* (SOHO) sites that desire both redundancy in their Internet connections and the ability to easily load share across those links in an active-active configuration, without the complexity and operational expense of running BGP.

Another interesting functionality enabled by the LISP control plane is the ability to mitigate some types of DoS attacks. In particular, if an ETR notices that it is the subject of a DoS attack from behind an ITR (that is, DoS packets are destined to an EID-prefix for which it is authoritative), it can use the LISP locator reachability bits (see Figure 2) to tell the source ITR that the RLOC for that EID-prefix is not available. The ETR accomplishes this by sending a locator-reachability bit of zero for the RLOC to the offending ITR. Note that this functionality is similar to Ioannidis and Bellocin’s “ICMP Pushback” proposal^[25].

Performance Considerations

LISP and its associated mapping protocol(s) have two primary performance concerns:

- Encapsulation overhead
- EID-to-RLOC lookup latency and packet loss

In the case of encapsulation overhead, the concern is that the addition of the LISP header will cause the encapsulated packet to exceed the path *Maximum Transmission Unit* (MTU). As mentioned previously, this area of research is still active (see, for example, [18]).

In the case of lookup latency and packet loss, because LISP-ALT uses BGP to find a particular EID-to-RLOC mapping, there could be latency associated with the first few packets in the first flow between sites (note that it is only the first flow; subsequent flows can use the mapping installed in the ITR). However, this latency is mitigated, and the initial packets are not lost because LISP can send the first few data packets over the control plane; these packets are the Data Probe packets. There is additional latency associated with the time required for the destination ETR to return the Map Reply. However, after this initial transaction is completed, no additional latency is injected by the mapping system.

As mentioned previously, there is a trade-off in the mapping system among the state required to be held by network elements, the rate of updates to the mapping system, and the latency incurred when looking up an EID-to-RLOC mapping. LISP-ALT is a hybrid (push/pull) architecture that attempts to minimize the state requirements on ITRs, while at the same time minimizing lookup latency.

Conclusions

LISP is a new protocol that implements the Loc/ID split using a map-and-encap protocol. It obtains the advantages of the level of indirection afforded by the Loc/ID split while minimizing changes to hosts and to the core routing system. In addition, LISP enables new functions such as BGP-free multihoming in an active-active configuration.

Acknowledgments

The LISP specification and supporting documents are the work of many people, including Scott Brim, Noel Chiappa, Dino Farinacci, Vince Fuller, Eliot Lear, Darrel Lewis, and Dave Oran.

References

- [0] Brim, S., et al., "EID Mappings Multicast Across Cooperating Systems for LISP," Internet Draft, Work in Progress, **draft-curran-lisp-emacs-00.txt**
- [1] Brim, S., et al., "LISP-CONS: A Content distribution Overlay Network Service for LISP," Internet Draft, Work in Progress, **draft-meyer-lisp-cons-03.txt**
- [2] Chiappa, N., "Endpoints and Endpoint Names: A Proposed Enhancement to the Internet Architecture,"
<http://ana.lcs.mit.edu/~jnc//tech/endpoints.txt>
- [3] Farinacci, D., et al., "Locator/ID Separation Protocol (LISP)," Internet Draft, Work in Progress, **draft-farinacci-lisp-06.txt**
- [4] Fuller, V., et al., "LISP Alternative Topology (LISP-ALT)," Internet Draft, Work in Progress, **draft-fuller-lisp-alt-01.txt**
- [5] Jen, D., et al., "APT: A Practical Transit Mapping Service," Internet Draft, Work in Progress, **draft-jen-apt-01.txt**
- [6] Lear, E., "NERD: A Not-so-Novel EID to RLOC Database," Internet Draft, Work in Progress, **draft-lear-lisp-nerd-03.txt**

- [7] Massey, D., Wang, L., Zhang, B., and L. Zhang, "A Proposal for Scalable Internet Routing and Addressing," Internet Draft, Work in Progress, **draft-wang-ietf-efit-01.txt**
- [8] Meyer, D., et al., "Report from the IAB Workshop on Routing and Addressing," RFC 4984, September 2007.
- [9] Narten, T., et al., "Routing and Addressing Problem Statement," Internet Draft, Work in Progress, **draft-narten-radir-problem-statement-01.txt**
- [10] Nordmark, E., "Shim6: Level 3 Multihoming Shim Protocol for IPv6," Internet Draft, Work in Progress, **draft-ietf-shim6-proto-09.txt**
- [11] O'Dell, M., "GSE - An Alternate Addressing Architecture for IPv6," <http://www.watersprings.org/pub/id/draft-ietf-ipngwg-gseaddr-00.txt>
- [12] Templin, F., "The IPvLX Architecture," Internet Draft, Work in Progress, **draft-templin-ipv1x-08.txt**
- [13] Vogt, C., "Six/One: A Solution for Routing and Addressing in IPv6," Internet Draft, Work in Progress, **draft-vogt-rrg-six-one-01.txt**
- [14] Whittle, R., "Ivip (Internet Vastly Improved Plumbing) Architecture," Internet Draft, Work in Progress, **draft-whittle-ivip-arch-01.txt**
- [15] Bates, T., et al., "Multiprotocol Extensions for BGP-4," RFC 2858, June 2000.
- [16] Farinacci, D., et al., "Generic Routing Encapsulation (GRE)," RFC 2784, March 2000.
- [17] Rekhter, Y., (Ed.), et al., "A Border Gateway Protocol 4 (BGP-4)," RFC 4271, January 2006.
- [18] Templin, F., "Subnetwork Encapsulation and Adaptation Layer," Internet Draft, Work in Progress, **draft-templin-seal-02.txt**
- [19] Bagnulo, M., "Preliminary LISP Threat Analysis," Internet Draft, Work in Progress, **draft-bagnulo-lisp-threat-01.txt**
- [20] Perkins, C., "IP Mobility Support for IPv4, revised," Internet Draft, Work in Progress, **draft-ietf-mip4-rfc3344bis-05.txt**

- [21] Johnson, D., Perkins, C., and J. Arkko, "Mobility Support in IPv6," RFC 3775, June 2004.
- [22] Arkko, J., Vogt, C., and W. Haddad, "Enhanced Route Optimization for Mobile IPv6," RFC 4866, May 2007.
- [23] Fenner, B., et al., "Protocol Independent Multicast - Sparse Mode (PIM-SM): Protocol Specification (Revised)," RFC 4601, August 2006.
- [24] Hinden, R., "New Scheme for Internet Routing and Addressing (ENCAPS) for IPNG," RFC 1955, June 1996.
- [25] Ioannidis John, and Bellovin, S., "Pushback: Router-Based Defense Against DDoS Attacks,"
<http://citeseer.ist.psu.edu/420554.html>
- [26] Huston, G., "More ROAP: Routing and Addressing at IETF68," *The Internet Protocol Journal*, Volume 10, No. 2, June 2007.
- [27] Carlos J. Bernardos, Ignacio Soto, and María Calderón, "IPv6 Network Mobility," *The Internet Protocol Journal*, Volume 10, No. 2, June 2007.

DAVID MEYER is currently a Director in the Advanced Research and Technologies Group at Cisco Systems, where he works on future directions for Internet technologies. E-mail: dmm@cisco.com

Book Review

Patterns in Network Architecture

Patterns in Network Architecture: A Return to Fundamentals, by John Day, ISBN-10: 0132252422, ISBN-13: 9780132252423, Prentice Hall, 2007. <http://www.informit.com/store/product.aspx?isbn=0132252422>

It isn't every day (pun intended) that one of the true Old Guard writes and publishes a book, and it behooves us to take notice. In this case, the author's expertise and his subject matter are of particular timeliness, because of the worldwide resurgence of activities with regard to next-generation network architectures, that is, a replacement, or upgrade to the Internet (dare one say "Internet 2.0"?).

John Day is a well-known scholar of historical cartography, and this book, in a way, is a roadmap of network architecture. The roadmap starts back in 1970, tracing from the roots of connectionless packet-switched dynamically routed systems such as Cyclades, and the ARPANET, through to recent discussions on multihoming, multicast, and mobility, with a view along the way of naming, addressing, protocol stack design, protocol design, and concepts of layering.

That description makes the book sound fairly standard in terms of structure and content, but it isn't. The book includes many discursive elements whose intent is to provide a collection of *patterns*. Design patterns originated in the building trade as a way for crafts people to pass on successful methods of construction (in the sense of affordable and noncollapsing) to less-inventive people (or people who want to spend their inventive efforts in different areas). Software engineers picked up on this idea, applying the techniques in both the microscopic world: patterns allow you to decide what algorithm is applicable in solving a problem in the small; and the macroscopic world: architectural patterns allow you to decide on an approach to breaking down a large system into the right kind of components.

Essentially, this book does the same thing, at the protocol stack level, and at the system level, with a collection of historical and contemporary examples to support the arguments.

The book makes interesting reading, especially as it represents a fair balance in reporting the early ideas that came not just from the United States, and restates the importance of the *Opens Systems Interconnection* (OSI) model (not the ISO protocols) in understanding layering and beads-on-a-string, as well as reasserting the use of the model in clarifying the perennially confusing concepts of names, addresses, and routes.

The book begins with a discussion of seven principles that emerged through the early history of networking (I won't spoil the book for readers by listing them here), and ends in the tenth and final chapter, entitled "Backing Out of a Blind Alley," with an appeal to fundamentals. Essentially, the author points out that researchers (especially academics) are strongly motivated to keep moving on with claims of ever-newer tricks, but rarely to consolidate these tricks into a set of principles that stand for a long time (because then they would have to completely change the topic of their research). Thus uncovering a foundational theory of networking would put a whole generation of networkers out of work (or funding at least).

The book is peppered (saltily) with fine quotes and fascinating asides from philosophy (for this reader, especially, the Chinese diversions were most novel and illuminating). Illustrative of the range is that one finds Wittgenstein and Dave Clark, Confucius, and Dr. Seuss—Frege's useful reminder that "The sign '=' should be read as 'is easily confused with'" would make an excellent IETF T-shirt.

I found the book extremely readable and enjoyable, and although I might argue with some of the opinions in the book, I think that this is just more evidence that I should recommend the book to anyone interested in knowing why we are where we are in networking, and being better informed about where we should go next.

—Jon Crowcroft, *University of Cambridge*
Jon.Crowcroft@cl.cam.ac.uk

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the "networking classics." In some cases, we may be able to get a publisher to send you a book for review if you don't have access to it. Contact us at ipj@cisco.com for more information.

ICANN Recovers Large Block of Internet Address Space

The *Internet Corporation for Assigned Names and Numbers* (ICANN) has found a little breathing room in the IPv4 address space with its recovery of a block of 16 million IPv4 addresses.

The IP addresses recovered were once used to connect older protocol packet-data networks with the fledgling Internet. The block of addresses, technically referred to as **14.0.0.0/8**, is also known as “Net-14.”

“Net-14 was the easiest network to reclaim, the so-called low hanging fruit,” said Barbara Roseman, General Manager with the *Internet Assigned Numbers Authority* (IANA), which is operated by ICANN. “None of the other legacy assignments in the IPv4 space are likely to be completely reclaimed as they are all in active use.”

A small percentage of the addresses in Net-14 had been assigned, most more than 15 years ago. The assignments were so old that finding people who knew about them was a lengthy process. Nearly 50 organizations worked cooperatively with ICANN staff throughout 2007 to confirm that the 984 registrations were no longer in use. IANA undertook the reclamation effort to ensure that the greatest number of IPv4 addresses can be made available to Internet users as the overall free pool of IPv4 addresses is depleted. IANA allocates IPv4 and IPv6 addresses to *Regional Internet Registries* (RIRs). The five RIRs allocate addresses to network operators in their local regions. IANA allocated more than one /8 (16m IPv4 addresses) per month in 2007 and the rate of allocation is not expected to slow in 2008. The reclamation of Net-14 means there are now 43 unallocated /8s left.

“The recovery of these addresses offers some breathing room as the four billion addresses in IPv4 space are depleted, but it is only a temporary solution,” added Roseman. “The real and lasting solution is the technical move to IPv6—the protocol that will make 340 trillion trillion unique IP addresses available.”

IPv6 Address Added for Root Servers in the Root Zone

ICANN recently took another step along the path of deployment for the next-generation IPv6 Internet addressing system. IPv6 addresses were added for six of the world’s 13 *root server* networks (A, F, H, J, K, M) to the appropriate files and databases. This move allows for the possibility of fuller IPv6 usage of the *Domain Name System* (DNS). Prior to today, those using IPv6 had needed to retain the older IPv4 addressing system in order to be able to use domain names.

“The ISP community welcomes this development as part of the continuing evolution of the public Internet,” said Tony Holmes, chair of ICANN’s Internet Service and Connectivity Provider Constituency. “IPv6 will be an essential part our future and support in the root servers is essential to the growth, stability, and reliability of the public Internet.”

Name server software relies on the root servers as a key part in translating domains like `icann.org` into the routing identifiers used by computers to connect to one another. In 2007 the ICANN *Security and Stability Advisory Committee* concluded that ICANN should move forward with the enhancement of the DNS root service by adding IPv6 addresses for the root servers. “The addition of IPv6 addresses for the root servers enhances the end-to-end connectivity for IPv6 networks, and furthers the growth of the global interoperable Internet,” added David Conrad, ICANN’s Vice President of Research and IANA Strategy. “This is a major step forward for IPv6-only connectivity and the global migration to IPv6.”

Further technical information on the move is available at:

<http://www.iana.org/reports/root-aaaa-announcement.html>

RIPE NCC Publishes Case Study of YouTube Hijack

As you may be aware from recent news reports, traffic to the `youtube.com` Website was “hijacked” on a global scale on Sunday February 24, 2008. The incident was a result of the unauthorized announcement of the prefix `208.65.153.0/24` and caused the popular video sharing Website to become unreachable from most, if not all, of the Internet. The RIPE NCC conducted an analysis into how this incident was seen and tracked by the RIPE NCC’s *Routing Information Service* (RIS) and has published a case study at:

<http://www.ripe.net/news/study-youtube-hijacking.html>

The RIPE NCC RIS is a service that collects *Border Gateway Protocol* (BGP) routing information from roughly 600 peers at 16 *Internet Exchange Points* (IXPs) across the world. Data is stored in near real-time and can be instantly queried by anyone to provide multiple views of routing activity for any point in time. The RIS forms part of the RIPE NCC’s suite of Information Services, which together provide a deeper insight into the workings of the Internet. The RIPE NCC is a neutral and impartial organization, and commercial interests therefore do not influence the data collected. The RIPE NCC Information Services suite also includes the *Test Traffic Measurement* (TTM) service, the *DNS Monitoring* (DNSMON) service and Hostcount. All of these services are available to anyone, and most of them are offered free of charge.

More information about RIPE NCC Information Services can be found at: <http://is-portal.ripe.net>

IETF Examines Future of the Internet by Going IPv6 Native

The *Internet Engineering Task Force* (IETF) put a spotlight on the next generation of Internet addressing when it switched off attendees' access to IPv4 during its March 2008 meeting. For an hour, Internet engineers at the meeting could only access the Internet using an IPv6 network.

During this event, IETF participants were encouraged to explore the Internet as it appears today in the IPv6 environment. The purpose of this exploration was to determine the next steps necessary toward deployment of IPv6 as the next generation of Internet addressing. The IETF undertook this activity at a time when IPv6-implementation is becoming a matter of global importance for the Internet. The event provided all IETF meeting attendees a first-hand opportunity to work with the Internet over an exclusive IPv6 network. "We get a lot of reports from members of our community who use IPv6, but this was an opportunity for everyone to observe and discuss the technical issues as a group," said Russ Housley, Chair of the IETF. "This first-hand data helps to inform our engineering decisions."

Some members of the Internet technical community assert that the ongoing deployment of IPv6 has been held back by a lack of IPv6-accessible Websites, creating the classic first-step dilemma for network operators. "It has been incredible to observe as members of the community organized themselves and updated their home networks to be ready for this event," said Leslie Daigle, Chief Internet Technology Officer at the Internet Society. "As we continue to solve the engineering and implementation obstacles to IPv6 deployment, creative engineers around the world will develop new uses for the Internet, through IPv6, in ways we can't yet imagine."

The IETF has provided dual stack IPv4/IPv6 network connectivity at its meetings for years, which has been useful for its regular IPv6-using attendees. The difference during this meeting was that a strictly IPv6 network was made available as well, and all attendees were encouraged to explore and experiment with the Internet as seen from IPv6. This focus was heightened when IPv4 access was deliberately shut off for an hour, leaving only IPv6 for connectivity. Following this—and other similar experiments—the engineering community expects to have a better understanding of the next steps necessary in the development of protocols and standards to support the continued deployment of IPv6 in support of the global Internet. The Comcast Corporation provided the facilities to conduct the live test of IPv6 and was the host sponsor of IETF-71 in Philadelphia.

For more information about this event, and similar events please see:

http://www.isoc.org/educpillar/resources/ipv6_faq.shtml

http://wiki.tools.isoc.org/IETF71_IPv4_Outage

<http://www.civil-tongue.net/clusterf/>

Postel Network Operator's Scholarship 2008

The *North American Network Operators' Group* (NANOG) and the *American Registry for Internet Numbers* (ARIN) have been unique and successful cooperative fora for Internet builders in North America and other parts of the world. Senior practitioners from around the world contribute their time to NANOG and ARIN as presenters, teachers and trainers, to produce consistent non-commercial conferences of high-quality.

Since 2007, the generosity of an anonymous donor and the administration of the Internet Society, have allowed NANOG and ARIN to provide financial support to a person from a developing country to participate in the October joint NANOG/ARIN meeting through the *Postel Network Operator's Scholarship*.

The Scholarship Committee cordially invites suitable applicants to apply for fellowship funding to participate in the October 2008 joint NANOG/ARIN meeting. The Scholarship targets personnel from developing countries who are actively involved in Internet development, in any of the following roles: Engineers (Network Builders), Operational and Infrastructure Support Personnel, and Educators, Teachers, and Trainers

Successful applicants will be provided with transportation to and from the meetings and a reasonable allowance for food and accommodation. In addition all fees for participation in the conferences, tutorials, and social events will be waived. Applicants from any part of the world will be considered. The deadline for application is June 1, 2008, and the awardee will be informed by July 1, 2008.

To apply for the fellowship please read <http://www.nanog.org/postel-scholarship.html> and submit your application by e-mail to PostelNOS@nanog.org

For more information about NANOG and ARIN meetings, see: <http://www.nanog.org/> and <http://www.arin.net/>

JPNIC Releases IPv4 Exhaustion Report

The *Japan Network Information Center* (JPNIC) has released a report entitled "Study Report on the IPv4 Address Space Exhaustion Issue (Phase I)." The report can be downloaded from the following link:

<http://www.nic.ad.jp/en/ip/ipv4pool/ipv4exh-report-071207-en.pdf>

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, fire-walls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2008 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA
--

The Internet Protocol Journal

June 2008

Volume 11, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
A Decade of Internet Evolution.....	2
A Decade in the Life of the Internet	7
Mobile WiMAX	19
Letters to the Editor.....	36
Fragments	39

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

FROM THE EDITOR

Ten years ago we published the first issue of *The Internet Protocol Journal* (IPJ). Since then, 41 issues and a total of 1,612 pages have been produced. Today, IPJ has about 37,000 subscribers all around the world. Although most of our readers prefer the paper edition, a growing number of subscribers are reading IPJ online or downloading the PDF version. This shift in reading habits may be related to the changes in technology over the last 10 years. Lower costs and higher-resolution displays and printers, as well as improvements in Internet access technologies, have made the online “experience” a lot better than in 1998.

Publishing is by no means the only area that has seen dramatic changes in the last decade. We asked Vint Cerf and Geoff Huston to reflect on Internet developments in this period, and the resulting articles, “A Decade of Internet Evolution” and “A Decade in the Life of the Internet,” are included in this issue.

Let me take this opportunity to thank all those people who have made IPJ possible. Our authors deserve a round of applause for carefully explaining both established and emerging technologies. They are assisted by an equally insightful set of reviewers and advisors who provide feedback and suggestions on every aspect of our publications process. The process itself relies heavily on two individuals: Bonnie Hupton, our copy editor, and Diane Andrada, our designer. Thanks go also to our printers and mailing and shipping providers. Last, but not least, our readers provide encouragement, suggestions, and feedback. This journal would not be what it is without them.

Because we are considering some Internet history in this issue, I would like to announce a project that takes us even further back. Before joining Cisco in 1998 I worked at the Interop Company, where I was responsible for the monthly publication of *ConneXions—The Interoperability Report*, published from 1987 through 1996. Unlike IPJ, *ConneXions* was produced in the “old-fashioned way” using various pieces of text and artwork assembled onto paste-up boards, and then photographed for subsequent plate making and offset printing. Thus no PDF files were produced at the time, but I am pleased to announce that *The Charles Babbage Institute* at the University of Minnesota has scanned the complete collection (117 issues) and it is now available at: <http://www.cbi.umn.edu/hostedpublications/Connexions/index.html>

Our final article is a look at Mobile WiMAX. WiMAX is an emerging technology that was originally designed as a fixed wireless broadband technology, a “DSL replacement,” but has evolved to support mobility.

— Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

A Decade of Internet Evolution

by Vinton G. Cerf, Google

In 1998 the Internet had about 50 million users, supported by approximately 25 million servers (Web and e-mail hosting sites, for example, but not desktops or laptops). In that same year, the *Internet Corporation for Assigned Names and Numbers* (ICANN)^[1] was created. Internet companies such as Netscape Communications, Yahoo!, eBay, and Amazon were already 3 to 4 years old and the Internet was in the middle of its so-called “dot-boom” period. Google emerged that year as a highly speculative effort to “organize the world’s information and make it accessible and useful.” Investment in anything related to the Internet was called “irrational exuberance” by the then head of the U.S. Federal Reserve Bank, Alan Greenspan.

By April 2000, the Internet boom ended—at least in the United States—and a notable decline in investment in Internet application providers and infrastructure ensued. Domino effects resulted for router vendors, Internet service providers, and application providers. An underlying demand for Internet services remained, however, and it continued to grow, in part because of the growth in the number of Internet users worldwide.

During this same period, access to the Internet began to shift from dial-up speeds (on the order of kilobits to tens of kilobits per second) to broadband speeds (often measured in megabits per second). New access technologies such as digital subscriber loops and dedicated fiber raised consumer expectations of Internet capacity, in turn triggering much interest in streaming applications such as voice and video. In some locales, consumers could obtain gigabit access to the Internet (for example, in Japan and Stockholm). In addition, mobile access increased rapidly as mobile technology spread throughout the world, especially in regions where wireline telephony had been slow to develop.

Today the Internet has an estimated 542 million servers and about 1.3 billion users. Of the estimated 3 billion mobile phones in use, about 15 percent are Internet-enabled, adding 450 million devices to the Internet. In addition, at least 1 billion personal computers are in use, a significant fraction of which also have access to the Internet. The diversity of devices and access speeds on the Internet combine to produce challenges and opportunities for Internet application providers around the world. Highly variable speeds, display areas, and physical modes of interaction create a rich but complex canvas on which to develop new Internet applications and adapt older ones.

Another well-documented but unexpected development during this same decade is the dramatic increase in user-produced content on the Internet. There is no question that users contributed strongly to the utility of the Internet as the World Wide Web made its debut in the early 1990s with a rapidly growing menu of Web pages.

But higher speeds have encouraged user-produced audio and video archives (*Napster* and *YouTube*), as well as sharing of all forms of digital content through peer-to-peer protocols. Voice over IP, once a novelty, is very common, together with video conferencing (*iChat* from Apple, for example).

Geographically indexed information has also emerged as a major resource for Internet users. In the scientific realm, *Google Earth* and *Google Maps* are frequently used to display scientific data, sensor measurements, and so on. Local consumer information is another common theme. When I found myself in the small town of Page, Arizona, looking for saffron to make paella while in a houseboat on Lake Powell, a Google search on my Blackberry quickly identified markets in the area. I called one of them and verified that it had saffron in stock. I followed the map on the Website and bought 0.06 ounces of Spanish saffron for about \$12.99. This experience reinforced my belief that having locally useful information at your fingertips no matter where you are is a powerful ally in daily living.

New business models based on the economics of digital information are also emerging. I can recall spending \$1,000 for about 10 MB of disk storage in 1979. Recently I purchased 2 TB of disk storage for about \$600. If I had tried to buy 2 TB of disk storage in 1979, it would have cost \$200 million, and probably would have outstripped the production capacity of the supplier. The cost of processing, storing, and transporting digital information has changed the cost basis for businesses that once required the physical delivery of objects containing information (books, newspapers, magazines, CDs, and DVDs). The Internet can deliver this kind of information in digital form economically—and often more quickly than physical delivery. Older businesses whose business models are based on the costs of physical delivery of information must adapt to these new economics or they may find themselves losing business to online competitors. (It is interesting to note, however, that the Netflix business, which delivers DVDs by postal mail, has a respectable data rate of about 145 kbps per DVD, assuming a 3-day delivery time and about 4.7 GB per DVD. The CEO of Netflix, Reed Hastings, told me nearly 2 years ago that he was then shipping about 1.9 million DVDs per day, for an aggregate data rate of about 275 Gbps!)

Even the media that have traditionally been delivered electronically such as telephony, television, and radio are being changed by digital technology and the Internet. These media can now be delivered from countless sources to equally countless destinations over the Internet. It is common to think of these media as being delivered in streaming modes (that is, packets delivered in real time), but this need not be the case for material that has been prerecorded. Users of iPods have already discovered that they can download music faster than they can listen to it.

With gigabit access to the Internet, one could download an hour's worth of conventional video in about 16 seconds. This fact certainly changes my understanding of "video on demand" from a streaming delivery to a file transfer. The latter is much easier on the Internet because one is not concerned about packet inter-arrival times (jitter), loss, or even orderly delivery because the packets can be reordered and retransmitted during the file transfer. I am told that about 10 hours of video are being uploaded to YouTube per second.

The battles over *Quality of Service* (QoS) are probably not over yet either. Services such as *Skype* and applications such as iChat from Apple demonstrate the feasibility of credible, real-time audio and video conferencing on the "best-efforts" public Internet. I have been surprised by the quality that is possible when both parties have reasonably high-capacity access to the Internet.

Technorati is said to be tracking on the order of 112 million blogs, and the *China Internet Network Information Center* (CNNIC) estimates 72 million Chinese blogs that are probably in addition to those tracked by Technorati. Adding to these are billions of Web pages and, perhaps even more significant, an unknown amount of information online in the form of large databases. The latter are not indexed in the same way that Web pages can be, but probably contain more information. Think about high-energy physics information, images from the Hubble and other telescopes, radio telescope data including the *Search for Extra-Terrestrial Intelligence* (SETI)^[2], and you quickly conclude that our modern society is awash in digital information.

It seems fair to ask how long accessibility of this information is likely to continue. By this question I do not mean that it may be lost from the Internet but, rather, that we may lose the ability to interpret it. I have already encountered such problems with image files whose formats are old and whose interpretation by newer software may not be possible. Similarly, I have ASCII text files from more than 20 years ago that I can still read, but I no longer have operating software that can interpret the formatting instructions to produce a nicely formatted page. I sometimes think of this problem as the "year 3000" problem: It is the year 3000 and I have just finished a Google search and found a PowerPoint 1997 file. Assuming I am running Windows 3000, it is a fair question whether the format of this file will still be interpretable. This problem would arise even if I were using open-source software. It seems unlikely that application software will last 1000 years in the normal course of events unless we deliberately take steps to preserve our ability to interpret digital content. Absent such actions, we will find ourselves awash in a sea of rotting bits whose meaning has long since been lost.

This problem is not trivial because questions will arise about intellectual property protection of the application, and even the operating system software involved. If a company goes out of business or asserts that it will no longer support a particular version of an application or operating system, do we need new regulations that require this software to be available on the public Internet in some way?

Even if we have skirted this problem in the past by rendering information into printed form, or microfilm, the complexity of digital objects is increasing. Consider spreadsheets or other complex objects that really cannot be fully “rendered” without the assistance of application software. So it will not be adequate simply to print or render information in other long-lived media formats. We really will need to preserve our ability to read and interpret bits.

The year 2008 also marks the tenth anniversary of a project that started at the U.S. Jet Propulsion Laboratory: *The Interplanetary Internet*. This effort began as a protocol design exercise to see what would have to change to make Internet-like capability available to manned and robotic spacecraft. The idea was to develop networking technology that would provide to the space exploration field the kind of rich and interoperable networking between spacecraft of any (Earth) origin that we enjoy between devices on the Internet.

The design team quickly recognized that the standard TCP/IP protocols would not overcome some of the long delays and disruptions to be expected in deep space communication. A new set of protocols evolved that could operate above the conventional Internet or on underlying transport protocols more suited to long delays and disruption. Called “delay and disruption tolerant networking”^[3, 4] or DTN, this suite of protocols is layered in the same abstract way as the Internet. The Interplanetary system could be thought of as a network of Internets, although it is not constrained to use conventional Internet protocols. The analog of IP is called the *Bundle Protocol*^[5], and this protocol can run above TCP or the *User Datagram Protocol* (UDP) or the new *Licklider Transport Protocol* (for deep space application). Ironically, the DTN protocol suite has also proven to be useful for terrestrial applications in which delay and disruption are common: tactical military communication and civilian mobile communication.

After 10 years of work, the DTN system will be tested onboard the Deep Impact mission platform late in 2008 as part of a program to qualify the new technology for use in future space missions. It is hoped that this protocol suite can be standardized for use by any of the world’s space agencies so that spacecraft from any country will be interoperable with spacecraft of other countries and available to support new missions if they are still operational and have completed their primary missions. Such a situation already exists on Mars, where the Rovers are using previously launched orbital satellites to relay information to Earth’s Deep Space Network using store-and-forward techniques like those common to the Internet.

The Internet has gone from dial-up to deep space in just the past 10 years. One can only begin to speculate about its application and condition 10 years hence. We will all have to keep our subscriptions to *The Internet Protocol Journal* to find out!

References

- [1] Cerf, V., “Looking Toward the Future,” *The Internet Protocol Journal*, Volume 10, No. 4, December 2007.
- [2] <http://www.seti.org>
- [3] <http://www.dtnrg.org/wiki>
- [4] V. Cerf, S. Burleigh, A. Hooke, L. Torgerson, R. Durst, K. Scott, K. Fall, and H. Weiss, “Delay-Tolerant Networking Architecture,” RFC 4838, April 2007.
- [5] Scott, K., and S. Burleigh, “Bundle Protocol Specification,” RFC 5050, November 2007.

VINTON G. CERF is vice president and chief Internet evangelist for Google. Cerf served as a senior vice president of MCI from 1994 through 2005. Widely known as one of the “Fathers of the Internet,” Cerf is the co-designer of the TCP/IP protocols and the architecture of the Internet. He received the U.S. National Medal of Technology in 1997 and the 2004 ACM Alan M. Turing award. In November 2005, he was awarded the Presidential Medal of Freedom. Cerf served as chairman of the board of the Internet Corporation for Assigned Names and Numbers (ICANN) from 2000 through 2007 and was founding president of the Internet Society. He is a Fellow of the IEEE, ACM, the American Association for the Advancement of Science, the American Academy of Arts and Sciences, the International Engineering Consortium, the Computer History Museum, and the National Academy of Engineering. He is an honorary Freeman of the City of London. Cerf holds a Bachelor of Science degree in Mathematics from Stanford University and Master of Science and Ph.D. degrees in Computer Science from UCLA. E-mail: vint@google.com

A Decade in the Life of the Internet

by Geoff Huston, APNIC

The evolutionary path of any technology can often take strange and unanticipated turns and twists. At some points simplicity and minimalism can be replaced by complexity and ornamentation, while at other times a dramatic cut-through exposes the core concepts of the technology and removes layers of superfluous additions. The technical evolution of the Internet appears to be no exception, and contains these same forms of unanticipated turns and twists.

This article presents a personal perspective of the evolution of the Internet over the last decade, highlighting my impressions of what has worked, what has not, and what has changed over this period. It has been an extraordinary decade for the Internet, encompassing a boom and a bust that would rate among history's best, a comprehensive restructuring of the communications industry, and a set of changes that have altered the way in which each of us now works and plays. And the Internet has even added a few new words to the language on the way.

Rather than offer a set of random observations, I will use the Internet Protocol model as a template, starting with the underlying transmission media, then looking at the internetwork layer, the transport layer, then applications and services, and, finally looking at the business of the Internet.

The Transmission Media Layer

It seems like it was in an entirely different lifetime, but the *Internet Service Provider* (ISP) business of 1998 was still centrally involved in the technology of dial-up modems. The state-of-the-art of modem speed had been continually refined from 9,600 bps to 14.4 kbps, to 28 kbps, to finally, 56 kbps, squeezing every last bit out the phase amplitude space contained in an analogue 3-KHz voice circuit. Modems were the bane of an ISP's life. They were capricious, constantly being superseded by the next technical refinement, unreliable, difficult for customers to use, and they were just slow. Almost everything else on the Internet was tailored to download reasonably quickly over a modem connection. Webpages were carefully tailored with compressed images, and plaintext was the dominant medium as a consequence.

Not all forms of Internet access were dial-up. ISDN was used in some places, but it was never cheap enough to take over as the ubiquitous access method. There were also access services based on *Frame Relay*, X.25, and various forms of digital data services. At the high end of the speed spectrum were T1 access circuits with 1.5-Mbps clocking, and T3 circuits clocked at 45 Mbps.

ISPs leased circuits from a telephony company (telco). In 1998 the ISP industry was undergoing a transition of its trunk IP infrastructure from T1 circuits to T3 circuits. It was not going to stop here, but squeezing even more capacity from the network was proving to be a challenge. Deployment of 622-Mbps IP circuits occurred, although many of these were constructed using 155-Mbps *Asynchronous Transfer Mode* (ATM) circuits using router load balancing to share the IP load over four of these circuits in parallel. Gigabit circuits were just beginning, and the initial tests of IP over 2.5-Gbps *Synchronous Digital Hierarchy* (SDH) circuits began in 1998.

In some ways 1998 was a pivotal year for IP transmission. Until this time IP was still just another application that was positioned as just another customer of the telco's switched-circuit infrastructure that was constructed primarily to support telephony. From the analogue voice circuits to the 64K digital circuit through to the trunk bearers, IP had been running on top of the voice network. By 1998 things were changing. The Internet had started to make ever larger demands on transmission capacity, and the factor accelerating further growth in the network was now not voice, but data. It made little sense to provision an ever larger voice-based switching infrastructure just to repackage it as IP, and by 1998 the industry was starting to consider just what an all-IP high-speed network would look like, from the photon all the way through to the application.

At the same time the fiber-optic systems were changing with the introduction of *Wavelength-Division Multiplexing* (WDM). Older fiber equipment with electro-optical repeaters and *Plesiochronous Digital Hierarchy* (PDH) multiplexers allowed a single fiber pair to carry around 560 Mbps of data. WDM allowed a fiber pair to carry multiple channels of data using different wavelengths, with each channel supporting a data rate of up to 10 Gbps. Channel capacity in a fiber strand is between 40 to 160 channels using *Dense WDM* (DWDM). Combined with the use of all-optical amplifiers, the most remarkable part of this entire evolution in fiber systems is that a Tbps cable system can be constructed today for much the same cost as a 560-Mbps cable system of the mid-1990s. The factor that accelerated deployment of these high-capacity fiber systems was never based on expansion of telephony, because the explosive growth of the industry was all about IP. So it came as no surprise that at the same time as the demand for IP transmission was increasing there was a shift in the transmission model, where instead of plugging routers into telco switching gear and using virtual point-to-point circuits for IP, we started to plug routers into wavelengths of the DWDM equipment and operate all-IP networks in the core of the Internet.

The evolution of access networks has seen a shift away from modems to numerous digital access methods, including DSL, cable modems, and high-speed wireless services. The copper pair of the telco network has proved surprisingly resilient, and DSL has achieved speeds of tens of megabits per second through this network, with the prospect of hundred-megabit systems appearing soon.

So, in terms of transmission, the last 10 years has seen the network migrate from an overlay system of kilobit-per-second access with multimegabit trunks operating as a customer of the telco switched network to a comprehensive IP network with access of megabits per second with multigigabit trunks, or a thousandfold increase in basic network capacity in that period.

The demand of the Internet for capacity continues, and we are now seeing work on standardizing 40- and 100-Gbps transmission systems in the IEEE; the prospect of terabit transmissions is now taking shape for the Internet.

The Internet Layer

If transmission has seen dramatic changes in the past decade, then what has happened at the IP layer over the same period?

The glib answer is “absolutely nothing!” But that answer would be ignoring a large amount of activity in this area. We have tried to change many parts of IP in the past decade, but, interestingly, none of the proposed changes has managed to gain any significant traction in the network, and IP today is largely no different from IP of a decade ago. *Mobility*^[1], *Multicast*^[2], and *IP Security* (IPSec)^[3] remain poised in the wings, still awaiting adoption by the Internet mainstream.

Quality of Service (QoS) was a “hot” topic in 1998, and it involved the search for a reasonable way for some packets to take the fast path while others took a more leisurely way through the network. We experimented with various forms of signaling, packet classifiers, queue-management algorithms, and interpretations of the *Type of Service* bits in the IPv4 packet header, and we explored the QoS architectures of *Integrated and Differentiated Services* in great detail. However, QoS never managed to achieve wide acceptance in mainstream Internet service environments. In this case the Internet took a simpler direction: In response to not enough network capacity, the alternate approach to installing additional mechanisms in the network—in the host protocol stack and even in the application in order to ration the capacity you have—is to simply expand the network to meet the total level of demand. So far the simple approach has prevailed in the network, and QoS remains largely unused^[4].

We have experimented with putting circuits back into the IP architecture in various ways, most notably with the *Multiprotocol Label Switching* (MPLS) technology^[5]. This technology used the label-swapping approach used in X.25, Frame Relay, and ATM virtual circuit switching systems; it created a collection of virtual paths from each network ingress to each network egress. The idea was that in the interior of the network you no longer needed to load up a complete routing table into each switching element, and instead of performing destination-address lookup you could perform a much smaller, and hopefully faster, label lookup.

This process did not eventuate, and switching packets using the 32-bit destination address continued to present much the same level of cost-efficiency at the hardware level as virtual circuit label switching. When you add the additional overhead of an additional level of indirection in terms of operational management of MPLS networks, MPLS became another technology that so far has not managed to achieve traction in mainstream Internet networks. However, MPLS is by no means a dormant technology, and one place where MPLS has enjoyed considerable deployment is in the corporate service sector where many *Virtual Private Networks*^[6] are constructed using MPLS as the core technology, steadily replacing a raft of traditional private data systems that used X.25, Frame Relay, ATM, *Switched Multimegabit Data Service* (SMDS), and switched Ethernet.

Of course one change at the IP level of the protocol stack that was intended in the past decade but has not occurred is *IP Version 6*^[7]. In 1998 we were forecasting that we would have consumed all the remaining unallocated IPv4 addresses by around 2008. We were saying at the time that, because we had completed the technical specification of IPv6, the next step was that of deployment and transition. There was no particular sense of urgency, and the comfortable expectation was that with a decade to go we did not need to raise any alarms. And this plan has worked, to some extent, in that today's popular desktop operating systems of Windows, MacOS, and UNIX all have IPv6 support. But other parts of this transition have been painfully slow. It was only a few months ago that the root of the *Domain Name System* (DNS) was able to answer queries using the IPv6 protocol as transport, and provide the IPv6 addresses of the root nameservers. Very few mainstream services are configured in a dual-stack fashion, and the prevailing view is still that the case for IPv6 deployment has not yet reached the necessary threshold. Usage measurements for IPv6 point to a level of deployment of around one-thousandth of the IPv4 network, and, perhaps more worrisome, this metric has not changed to any appreciable level in the past 4 years. So what about that projection of IPv4 unallocated pool exhaustion by 2008? How urgent is IPv6 now? The good news is that the *Internet Assigned Numbers Authority* (IANA) still has some 16 percent of the address space in its unallocated pool, so IPv4 address exhaustion is unlikely to occur this year. The bad news is that the global consumption rate of IP addresses is now at a level such that the remaining address pool can fuel the Internet for less than a further 3 years, and the exhaustion prediction is now sometime around 2010 to 2011.

So why have we not deployed IPv6 more seriously yet? And if we are not going to deploy IPv6, then what is the alternative? Of all the technical refinements to IP that have occurred, one that received little fanfare when it was first published has enjoyed massive deployment over the past decade, and that is the technology of *Network Address Translation* (NAT)^[8]. Today NAT devices are ubiquitous. It seems that every home access unit, every corporate firewall, every data center, and every service includes a NAT device.

One measure of the ubiquity of NATs is the transformation that has occurred in the application space. By 2008 applications have either adopted a strict client-server approach, where the client always initiates the network transaction, or were forced down a more complex path. Where there is some form of peer interaction, applications are now equipped with additional capabilities, including NAT behavior discovery, NAT binding management, application-level name spaces, and multiparty rendezvous mechanisms, all required to allow the application to function across NATs. So far we have managed to offload the problem of looming address scarcity in the Internet onto NATs, and the really significant change that has occurred in the past decade at the IP level is the default assumption about the semantics of an IP address. An IP address is no longer synonymous with the persistent identity of the remote party that anyone can use to initiate a communication, but a temporary token to allow a single transaction to complete. As a consequence, most Internet services have retreated into data centers and the business of hosting services has thrived. And the change that would have preserved the coherent end-to-end architecture of the Internet IP layer, namely IPv6, is still waiting for wide-scale deployment.

The next few years promise to be “interesting” in every form of meaning of the word. The exhaustion of the remaining IPv4 address pool is imminent, and if we are going to substitute IPv6 in place of IPv4, then we simply do not have enough time to achieve this substitution before the remaining IPv4 address pool is depleted. And although so far NATs have conveniently pushed the problem of increasing address scarcity off the network and over to the edge devices and onto applications, it is not clear that this approach can sustain an ever-growing Internet indefinitely. We have yet to understand just what a “carrier-grade NAT” might be, or whether it can even work in any useful manner. NATs were an accidental addition to the Internet, and their role in the coming years is unclear.

The early 1990s saw a flurry of activity in the routing space, and protocols were quickly developed and deployed. By 1998 the “standard” Internet environment involved the use of either *Intermediate System-to-Intermediate System* (IS-IS) or *Open Shortest Path First* (OSPF) as large-scale interior routing protocols and *Border Gateway Protocol 4* (BGP4) as the interdomain routing protocol^[9]. This picture has remained constant over the past decade. In some ways it is reassuring to see a technology that is capable of sustaining a quite dramatic growth rate, but perhaps that is not quite the complete picture.

We never quite completed the specification for the next interdomain routing protocol, and BGP4 is now showing signs of stress^[10]. The pool of *Autonomous System* (AS) numbers is forecast to run out early in 2011, and by then we need to have fielded a new variant of BGP that can operate with a much larger pool of AS numbers^[11].

Fortunately the technology development has been completed and an approach that allows incremental deployment has been devised, so this transition is not quite the traumatic transition that is associated with IPv6. But deployment is slow, and of the current level of adoption of the larger AS number set is, oddly enough, comparable to IPv6, at a level of around one-thousandth of the total AS number pool. The routing system has also been growing inexorably, and the capability of switching systems to cope with ever larger routing tables while at the same time offering continual improvements in cost-efficiencies is now looking less certain. So, once again we appear to be examining routing protocol theory and practice, and looking at alternate approaches to routing that can offer superior scaling properties to BGP for the future.

No listing of the major highlights in IP over the past decade would be complete without some mention of the perennial issue of *location* and *identity*.^[25] One of the original simplifications in the IP architecture was to place the semantics of identity, location, and forwarding into an IP address. Although that process has proved phenomenally effective in terms of simplicity of applications and simplicity of IP networks, it has posed some serious challenges with regard to mobility, routing, and network management. Each of these aspects of the Internet would benefit considerably if the Internet architecture allowed identity to be distinct from location. Numerous efforts have been directed at this problem over the past decade, particularly in IPv6, but so far we really have not arrived at an approach that feels truly comfortable in the context of IP.

So although it is possible to observe that not much has happened at the IP level in the past decade that is deployed in the Internet—and IP is still IP—there is still a considerable agenda to tackle at the Internet layer.

The Transport Layer

A decade ago, in 1998, the transport layer of the IP architecture consisted of the *User Datagram Protocol* (UDP) and TCP, and the network usage pattern was around 95-percent TCP and 5-percent UDP. Here, as well, not much has changed in the intervening 10 years.

We have developed two new transport protocols, the *Datagram Congestion Control Protocol* (DCCP) and the *Stream Control Transmission Protocol* (SCTP)^[12], which can be regarded as refinements of TCP to cover flow control for datagram streams in the case of DCCP and flow control over multiple reliable streams in the case of SCTP. However, in a world of transport-aware middleware that is the Internet today, the level of capability to actually deploy these new protocols in the public Internet is marginal at best.

TCP has proved to be remarkably resilient over the years, but as the capacity of the network increases the ability of TCP to continue to deliver ever faster data rates over distances that span the globe is becoming a significant concern. Recent times have seen much work to devise revised TCP flow-control algorithms that still share the network fairly with other concurrent TCP sessions, yet can ramp up to multigigabit-per-second data-transfer rates and sustain those rates over extended periods^[13]. At this stage much of this work is still in the area of research and experimentation, and TCP today as deployed on the Internet is much the same as TCP of a decade ago, with perhaps a couple of notable exceptions. The latest TCP stack from Microsoft in Vista uses dynamic tuning of the Receive window, and a larger inflation factor of the Send window in congestion avoidance where there is a large bandwidth delay product, and improved loss-recovery algorithms that are particularly useful in wireless environments. Linux now includes an implementation of *Binary Increase Congestion control* (BIC), which undertakes a binomial search to reestablish a sustainable send rate. Both of these approaches can improve the performance of TCP, particularly when sending the TCP session over long distances and trying to maintain high transfer speeds.

The Application and Service Layer

This area, unlike the transport layer, has seen quite profound changes over the past decade. A decade ago the Internet was on the cusp of portal mania, where *LookSmart* was the darling of the Internet boom and everyone were all trying to promote their own favorite “one stop shop” for all their Internet needs. We were still using various forms of hand-compiled directories, and navigation of the Internet was still the subject of various courses and books.

By 1998 *AltaVista* has made its debut, and change was already evident. This change, from directories and lists to active search, completely changed the Internet. These days we simply assume that we can type any query we have into a search engine and the search machinery will deliver a set of pointers to relevant documents. Each time this process occurs our expectations about the quality and utility of search engines are reinforced, and we have moved beyond swapping URLs as pointers and simply exchange search terms as an implicit reference to the material. Content is also changing as a result, because users no longer remain on a “site” and navigate around the site. Instead users are directing the search engines, and pulling the relevant page from the target site without reference to any other material.

Another area of profound change has been the rise of active collaboration over content, best typified in wikis. *Wikipedia* is perhaps the most cited example of user-created content, but almost every other aspect of content generation is also being introduced into the active user model, including *YouTube*, *Flickr*, *Joost*, and similar content.

Underlying these changes is another significant development, namely the changes in the content economy. In 1998 content providers and ISPs were competing for user revenue. Content providers were unable to make pay per view and other forms of direct financial relationship with users work in their favor, and were arguing that ISPs should fund content, because, after all, the only reason that users paid for Internet access was because of their perceived value of the content. ISPs, on the other hand, promoted the idea that content providers were enjoying a free ride across the ISP-funded infrastructure, and content providers should contribute to network costs. The model that has gained ascendancy as a result of this unresolved tension was that of advertised-funded content services, and this model has sustained a vastly richer, larger, and more compelling content environment.

At the same time the peer-to-peer network has emerged, and from its beginnings as a music-sharing subsystem, the distributed data model of content sharing now dominates the Internet with audio, video, and large data sets now using this form of content distribution and its associated highly effective transport architecture. Various measurements of Internet traffic have placed peer-to-peer content movement at between 40 and 80 percent of the overall traffic profile of the network.

In many ways applications and services have been the high frontier of innovation in the Internet in the past decade. An entire revolution in open interconnection of content elements is embraced under the generic term *Web 2.0*, and “content” is now a very malleable concept. It is no longer the case of “my computer, my applications, and my workspace” but an emerging model where not only the workspace for each user is held in the network, but where the applications themselves are part of the network, and all are accessed through a generic browser interface.

Any summary of the evolution of the application space over the last decade would not be complete without noting that whereas in 1998 the Internet was still an application that sat on top of the network infrastructure used to support the telephone network, by 2008 voice telephony was just another application layered on the infrastructure of the Internet, and the Internet had even managed to swallow the entire telephone number space into its DNS, using an approach called *ENUM*^[14].

The Business Layer

As much as the application environment of the Internet has been wildly erratic over the past decade, the business environment has been unpredictable as well, and the list of business winners and losers includes some of the historical giants of the telephone world as well as the Internet-bred new wave of entrants.

In 1998, despite the growing momentum of public awareness, the Internet was still largely a curiosity. It was an environment inhabited by geeks, game players, and academics, whose rites of initiation were quite arcane. As a part of the data networking sector, the Internet was just one further activity among many, and the level of attention from the mainstream telco sector was still relatively small. Most Internet users were customers of independent ISPs, and the business relationship between the ISP sector and the telco was tense and acrimonious. The ISPs were seen as opportunistic leeches on the telco industry; they ordered large banks of phone lines, but never made any calls; their customers did not hang up after 3 minutes, but kept their calls open for hours or even days at a time, and they kept ordering ever larger inventories of transmission capacity, yet had business plans that made the back of an envelope look professional by comparison. The telco was unwilling to make large long-term capital investments in additional infrastructure to pander to the extravagant demands of a wildcat set of Internet speculators and their fellow travelers. The telco, on the other hand was slow, expensive, inconsistent, ill-informed, and hostile to the ISP business. The telco wanted financial settlements and bit-level accounting, whereas the ISP industry appeared to manage quite well with a far simpler system of peering and tiering that avoided putting a value on individual packets or flows^[15]. This relationship was never going to last, and it resolved itself in ways that in retrospect were quite predictable. From the telco perspective it quickly became apparent that the only reason the telco was being pushed to install additional network capacity at ever increasing rates was the requirements of the ISP sector. From the ISP perspective the only way to grow at a rate that matched customer demand was to become one's own carrier and to take over infrastructure investment. And, in various ways, both outcomes occurred. Telcos bought ISPs, and ISPs became infrastructure carriers.

All this activity generated considerable investor interest, and the rapid value escalation of the ISP industry and then the entire Internet sector generated the levels of wild-eyed optimism that are associated only with an exceptional boom. By 2000 almost anything associated with the Internet, whether it was a simple portal, a new browser development, a search engine, or an ISP, attracted investor attention, and the valuations of Internet start-ups achieved dizzying heights. Of course one of the basic lessons of economic history is that every boom has an ensuing bust, and in 2001 the Internet bust happened. The bust was as inevitable and as brutal as the preceding boom was euphoric. But, like the railway boom and bust of the 1840s, when the wreckage was cleared away, what remained was a viable—and indeed a valuable—industry.

By 2003 the era of the independent retail ISP was effectively over. ISPs still exist, but those that are not competitive carriers tend to operate as IT business consultants that provide services to niche markets. Their earlier foray in to the mass market paved the way for the economies of scale that only the carrier industry could implement on the market.

But the grander aspirations of these larger players have not been met, and effective monopoly positions in many Internet access markets have not translated to effective control over the user's experience of the Internet, or anything even close to such control. The industry was already "unbundled," with intense competition occurring at every level of the market, including content, search, applications, and hosting. The efforts of the telco sector to translate their investment into mass-market Internet access into a more comprehensive control over content and its delivery in the Internet has been continually frustrated. The content world of the Internet has been reinvigorated by the successful introduction of advertiser-funded models of content generation and delivery, and this process has been coupled with the more recent innovations of turning back to the users themselves as the source of content, so that the content world is once again the focus of a second wave of optimism, bordering on euphoria.

And Now?

It has been a revolutionary decade for us all, and in the last 10 years the Internet has directly touched the lives of almost every person on this planet. Current estimates put the number of regular Internet users at 19 percent of the world's population.

Over this decade some of our expectations were achieved and then surpassed with apparent ease, whereas others remained elusive. And some things occurred that were entirely unanticipated. At the same time very little of the Internet we have today was confidently predicted in 1998, whereas many of the problems we saw in 1998 remain problems today.

What we have today is not the technical Internet we thought we were building a decade ago. It is not a coherent end-to-end network with clear signaling across commodity packet switching fabric, but a network that is replete with all forms of active middleware^[16], from NATs to firewalls^[17] and filters, including packet shapers, torrent detectors, *Voice over IP* (VoIP) blockers, and load balancers. It is neither a secure nor a safe network, but one that includes a continual barrage on end hosts in the form of more than a million different forms of viruses^[18], worms, and assorted malware^[19], as well as a barrage on users in the form of torrents of spam^[20]. The network is a host to a litany of hostile attacks, including gigabit traffic swamping attacks, redirection, inspection, passing off, and denial-of service attacks^[21]. The attacks are directed at links, routers^[22], the routing protocols^[23, 24], hosts, and applications. Our ability to effectively defend the network and its connected hosts continues to be, on the whole, ineffectual. Our level of interest in paying a premium to support highly secure systems still remains slight. But somehow we are not deterred by this situation. Somehow each of us has found a way to make our Internet work for us.

I am not sure that the next decade will bring the same level of intensity of structural change to the global communications sector, and perhaps that is a good thing given the collection of other challenges that are confronting us all in the coming decades. At the same time I think it would be good to believe that the past decade of development of the Internet has completely rewritten what it means to communicate, rewritten the way in which we can share our experience and knowledge, and, hopefully, rewritten the ways in which we can work together on these challenges.

References

The Internet Protocol Journal (IPJ) has published articles on all the major aspects of the technical evolution of the Internet over the past decade. To illustrate the extraordinary breadth of these articles, I have included as references here only articles that have been published in the IPJ.

- [1] Stallings, W., "Mobile IP," *IPJ*, Volume 4, No. 2, June 2001.
- [2] Handley, M., and Crowcroft, J., "Internet Multicast Today," *IPJ*, Volume 2, No. 4, December 1999.
- [3] Stallings, W., "IP Security," *IPJ*, Volume 3, No. 1, March 2000.
- [4] Huston, G., "QoS — Fact or Fiction?" *IPJ*, Volume 3, No. 1, March 2000.
- [5] Stallings, W., "MPLS," *IPJ*, Volume 4, No. 3, September 2001
- [6] Ferguson, P., and Huston, G., "What is a VPN?" *IPJ*, Volume 1, No. 1 & No. 2, June & September 1998.
- [7] Fink, R., "IPv6," *IPJ*, Volume 2, No. 1, March 1999.
- [8] Huston, G., "Anatomy: Inside Network Address Translators," *IPJ*, Volume 7, No. 3, September 2004.
- [9] Huston, G., "The BGP Routing Table," *IPJ*, Volume 4, No. 1, March 2001.
- [10] Huston, G., "Scaling inter-Domain Routing," *IPJ*, Volume 4, No. 4, December 2001.
- [11] Huston, G., "Exploring Autonomous System Numbers," *IPJ*, Volume 9, No. 1, March 2006.
- [12] Huston, G., "The Future for TCP," *IPJ*, Volume 3, No. 3, September 2000.
- [13] Huston, G., "Gigabit TCP," *IPJ*, Volume 9, No. 2, June 2006.

- [14] Huston, G., "ENUM," *IPJ*, Volume 5, No. 2, June 2002.
- [15] Huston, G., "Peering and Settlements," *IPJ*, Volume 2, No. 1 & No. 2, March & June 1999.
- [16] Huston, G., "The Middleware Muddle," *IPJ*, Volume 4, No. 2, June 2001.
- [17] Avolio, F., "Firewalls and Internet Security," *IPJ*, Volume 2, No. 2, June 1999.
- [18] Fraser, B., Rogers, L., and Pesante, L., "Was the Melissa Virus So Different?" *IPJ*, Volume 2, No. 2, June 1999.
- [19] Chen, T., "Virus Trends," *IPJ*, Volume 6, No. 3, September 2003.
- [20] Crocker, D., "Challenges in Anti-Spam Efforts," *IPJ*, Volume 8, No. 4, December 2005.
- [21] Patrikakis, C., Masikos, M., and Zouraraki, O., "Distributed Denial of Service Attacks," *IPJ*, Volume 7, No. 4, December 2004.
- [22] Lonvick, C., "Securing the Infrastructure," *IPJ*, Volume 3, No. 3, September 2000.
- [23] Kent, S., "Securing BGP: S-BGP," *IPJ*, Volume 6, No. 3, September 2003.
- [24] White, R., "Securing BGP: soBGP," *IPJ*, Volume 6, No. 3, September 2003.
- [25] Meyer, D., "The Locator Identifier Separation Protocol (LISP)," *IPJ*, Volume 11, No. 1, March 2008.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. The author of numerous Internet-related books, he is currently the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He was a member of the Internet Architecture Board from 1999 until 2005, and served on the Board of the Internet Society from 1992 until 2001. E-mail: gih@apnic.net

Mobile WiMAX

by Jarno Pinola and Kostas Pentikousis, VTT Technical Research Centre of Finland

One of the technologies that can lay the foundation for the next generation (fourth generation [4G]) of mobile broadband networks is popularly known as “WiMAX.” WiMAX, *Worldwide Interoperability for Microwave Access*, is designed to deliver wireless broadband bitrates, with *Quality of Service* (QoS) guarantees for different traffic classes, robust security, and mobility. This article provides an overview of mobile WiMAX, which is based on the wireless local and *Metropolitan-Area Network* (MAN) standards IEEE 802.16-2004^[1] and 802.16e-2005^[2]. We introduce WiMAX and focus on its mobile system profile and briefly review the role of the WiMAX Forum. We summarize the critical points of the WiMAX network reference model and present the salient characteristics of the PHY and MAC layers as specified in [1] and [2]. Then we address how mobile nodes enter a WiMAX network and explain the fundamentals of mobility support in WiMAX. Finally, we briefly compare WiMAX with *High-Speed Packet Access* (HSPA), another contender for 4G.

The Role of the WiMAX Forum

The WiMAX Forum is a nonprofit organization formed in 2001 to enhance the compatibility and interoperability of equipment based on the IEEE 802.16 family of standards. The IEEE 802.16 standards provide a large set of fundamentally different options for designing a wireless broadband system, including, for example, multiple options for *Physical* (PHY) layer implementation, *Media Access Control* (MAC) architecture, frequency bands, and duplexing. So many options lead to several possible system variants, which are all compatible with the IEEE standards. Although such multiplicity allows for deployment in very diverse environments, it may spell either solely vertical, single-vendor deployments or no deployment at all, because operators do not want to be locked in with any particular implementation. Thus, a major motivation for establishing the WiMAX Forum was to develop predefined system profiles for equipment manufacturers, which include a subset of the features included in the IEEE 802.16 standards. WiMAX Forum-certified products are guaranteed to be interoperable and to support wireless broadband services from fixed to fully mobile scenarios. The aim is to enable rapid market introduction of new standard-compliant WiMAX equipment and to promote the use of the technology in different sectors.

From IEEE 802.16 to Mobile WiMAX

The IEEE 802.16 standard was originally meant to specify a fixed wireless broadband access technique for point-to-point and point-to-multipoint links. During its development, however, it was decided that mobility support should also be considered.

The WiMAX Forum defines two system profiles based on [1] and [2], called *fixed* and *mobile* system profiles, respectively. Both include mandatory and optional PHY and MAC layer features that are required from all corresponding WiMAX-certified products. Because [1] and [2] specify only the PHY and MAC layers, an end-to-end architecture specification was deemed necessary in order to enable fast growth in manufactured quantities, market share, and interoperability. In response, the WiMAX Forum established the *Network Working Group* (NWG) with the aim of developing an end-to-end network reference model architecture based on IP supporting both fixed and mobile WiMAX (refer to [3] and [4]).

In short, according to the NWG reference model, a WiMAX network is partitioned into three independent architectural components: the user equipment (also referred to as *Customer Premises Equipment* [CPE]), the *Radio Access Network* (RAN, based on IEEE 802.16), and the network providing IP connectivity with the rest of the Internet. Clearly, this model allows a single operator to freely mix and match offerings from different manufacturers for these three parts, at least after interoperable equipment becomes readily available. Furthermore, in principle, each of these components of an operational network can be deployed and managed by different service providers. This scenario makes the network architecture flexible, eases network operation and maintenance, can increase competition under certain conditions, and is conducive to new business models. For example, municipalities can venture jointly with local or national network operators to deploy WiMAX in suburban and rural areas.

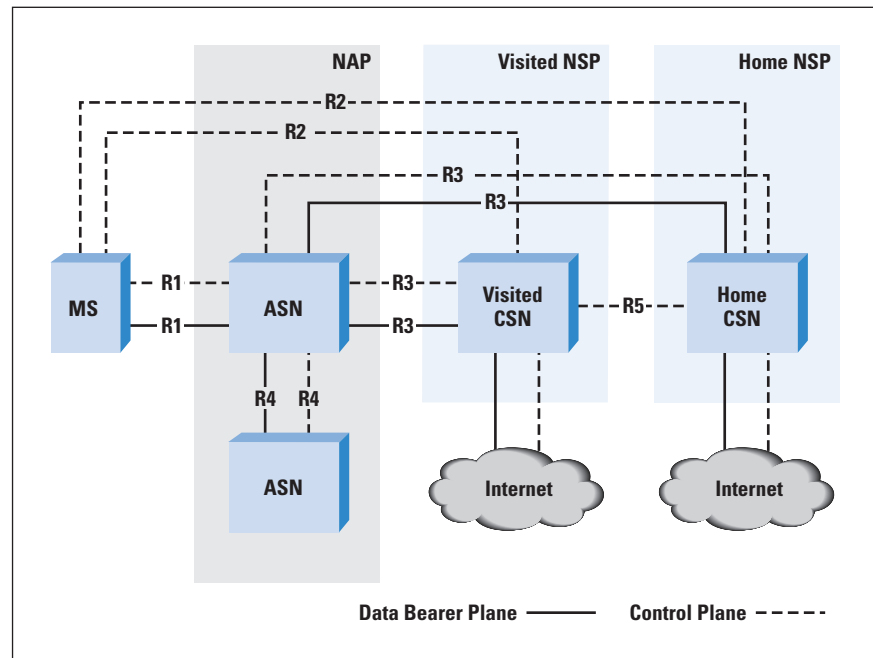
In contrast with earlier wireless data networks^[5], IP is fundamental in a WiMAX network. Indeed, IP currently plays a dominant role in the present state of the telecommunications industry. The premise is that by embracing IP, service providers and equipment manufacturers will face fewer problems when introducing WiMAX into their networks and product portfolios. Moreover, protocols standardized by the *Internet Engineering Task Force* (IETF) are preferred over proprietary solutions and are adopted as extensively as possible in the reference model.

Mobile WiMAX Network Reference Model

The WiMAX Forum NWG network reference model defines three basic architectural entities: the *Mobile Station* (MS), the *Access Service Network* (ASN), and the *Connectivity Service Network* (CSN). The role of the MS is to provide user access to the WiMAX network. The ASN is the Radio Access Network and is formed by numerous *Base Stations* (BSs) and *ASN Gateways* (ASN-GWs), managed by a *Network Access Provider* (NAP). CSN is the network entity providing IP connectivity to the WiMAX radio equipment, including all the IP core network functions required for internetworking with the rest of the world. CSNs are maintained by *Network Service Providers* (NSPs).

The ASN and CSN are further broken up into smaller functional entities, which communicate with each other using standardized interfaces called *reference points*. These reference points guarantee that a certain set of protocols and procedures are always supported and can function irrespective of the underlying hardware. The currently defined reference points are used for different control and management purposes, as well as for data bearing between the network entities. Figure 1 illustrates the network reference model and the main reference points.

Figure 1: WiMAX Forum NWG
Network Reference Model



The reference points are defined as follows in [3]: Reference point R1 consists of protocols and procedures compliant to [1], [2], and [6]. R1 implements the specifications of the air interface between the MS and the BS. R2, an interface between the MS and a CSN, is used solely for management purposes, including mobility management. R3 serves the same purpose between an ASN and a CSN, and R4 is used for micromobility management between two ASNs. R5 enables interworking between two CSNs for macromobility management.

In addition to reference points R1–R5, another three intra-ASN reference points are defined (not illustrated in Figure 1). R6, which consists of a set of control- and bearer-plane protocols for BS and ASN-GW communication, controls the data path and MS mobility events between these two ASN entities. R7 is an optional set of protocols used for coordinating R6 functions. Finally, R8 consists of bearer-plane protocols that enable data transfer between the base stations involved in a handover (also called *handoff*).

With respect to mobility, the reference model considers two different scenarios called *ASN-anchored mobility* and *CSN-anchored mobility*. ASN-anchored mobility (or intra-ASN mobility, or micromobility) management is employed when MS handovers occur from one BS to another, and both are controlled by the same ASN-GW. On the other hand, CSN-anchored mobility (or inter-ASN mobility, or macromobility) management is employed when MS movement dictates a handover from the currently serving BS to another one that is in a different subnetwork, controlled by a different ASN-GW. In the ASN-anchored case, handovers are managed solely by the MS and the ASN. In the CSN-anchored case, both ASN and CSN entities are engaged in mobility management.

Typically, ASN-anchored mobility procedures take precedence and CSN-anchored mobility management is employed only if necessary. Because ASN-anchored mobility takes place inside a single ASN, it does not change the MS network layer (IP) configuration. Three different functions are specified for ASN-anchored mobility management, all considered peer-to-peer interactions between different architectural entities:

- The *handoff* (HO) function controls the handover decision operation and handover signaling. The HO function supports mobile- and network-initiated handovers and, additionally, it may support *Fast Base Station Switching* (FBSS) or *Macro Diversity Handover* (MDHO)^[2].
- The *Data Path* (DP) function manages the data path setup and data packet transmission between two functional entities.
- The context function addresses the exchanges required in order to retrieve or set up any state in the network elements.

On the other hand, when MS movement necessitates CSN-anchored mobility management, the MS IP layer configuration changes as a result of the handover. In this case, mobility management is based on *Mobile IPv4* (MIPv4)^[7] or *Mobile IPv6* (MIPv6)^[8], if the MS supports it. Alternatively, the reference model adopts *Proxy MIP* (PMIP)^[9] to handle the handover. In PMIP, the MIP function is moved from the MS to a network instance called a *PMIPv4 client*, which takes care of all MIP signaling on behalf of the MS. Support for PMIP is specified only for MIPv4 in [3] and [4]. Note that in a handover from one ASN to another, MIP is used to complement ASN-anchored mobility management. The latter is still necessary to control the link-layer handover procedures. That is, after the micromobility handover is successfully completed, MIP independently takes care of the macromobility handover, that is, establishes communication paths between the new ASN-GW and the CSN. CSN-anchored mobility handovers are always network-initiated.

By embracing IETF protocols and providing an end-to-end architecture with independent functional entities, the WiMAX Forum NWG network reference model provides a clear framework for the application developers to work in. The model provides only operational requirements and does not prescribe particular technical solutions to realize them, allowing for proprietary yet standards-compliant implementations and enabling technical competition between different manufacturers.

Before examining mobility support in WiMAX, we review the basics of the IEEE 802.16 PHY and MAC layers.

OFDM and OFDMA

IEEE 802.16 and thus WiMAX adopted *Orthogonal Frequency Division Multiplexing* (OFDM), a multicarrier modulation scheme, as its PHY layer. In OFDM, the available bandwidth is divided into several parallel orthogonal subcarriers with lower bandwidth. A wideband channel is defined as a group of adjacent narrowband channels: a high-bitrate data stream is divided into these subcarriers and multiple narrowband data streams are transmitted over the air. Because the data symbol duration is inversely proportional to bitrate, the transmitted symbol duration is increased and the level of *Inter-Symbol Interference* (ISI) can be reduced. ISI is caused by multipath propagation in the wireless communication medium, where the transmitted data symbols can arrive at the receiver through different propagation routes because of reflections from buildings in urban areas and from hills and trees in rural areas. OFDM also uses guard intervals between successive data symbols and cyclic prefixes in order to decrease the effect of ISI even more.

One reason for the wide adoption of OFDM in modern broadband communication systems is its hardware implementation simplicity. OFDM signals can be formed and processed using *Inverse Fast Fourier Transform* (IFFT) and *Fast Fourier Transform* (FFT), at the transmitter and receiver, respectively, and both transforms can be implemented directly in hardware for higher performance. OFDM bodes well for mobile broadband systems through frequency diversity and adaptivity in both modulation and channel coding. By using *Adaptive Modulation and Coding* (AMC), the end-to-end quality deterioration due to the excess delays and deep fading conditions caused by mobility can be prevented, or at least diminished.

OFDM can also be used as a multiaccess scheme by having subcarriers grouped into subchannels, which can be assigned to different users contending for the data link. Each subchannel can contain a different number of subcarriers, and by altering the subcarrier group sizes and observing the channel conditions, it is possible to use differentiation in the channel allocation for different users.

This technique of using OFDM as a multiaccess scheme is called *Orthogonal Frequency Division Multiple Access* (OFDMA). Mobile WiMAX uses OFDMA as its PHY layer instead of plain OFDM, and subchannelization to both uplink and downlink transmissions is possible.

In OFDMA, the subcarriers assigned to subchannels can be either concurrent or taken from different regions of the total bandwidth. Both of these allocation schemes have advantages. When subcarriers assigned to one subchannel are distributed over the available bandwidth, frequency diversity can be attained. In mobile systems this diversity is advantageous because it can be used to make the transmission link more resistant against fast fading. A subchannelization scheme based on dispersed subcarrier allocation to subchannels, called *Partial Usage of Subcarriers* (PUSC), is mandatory in all mobile WiMAX implementations.

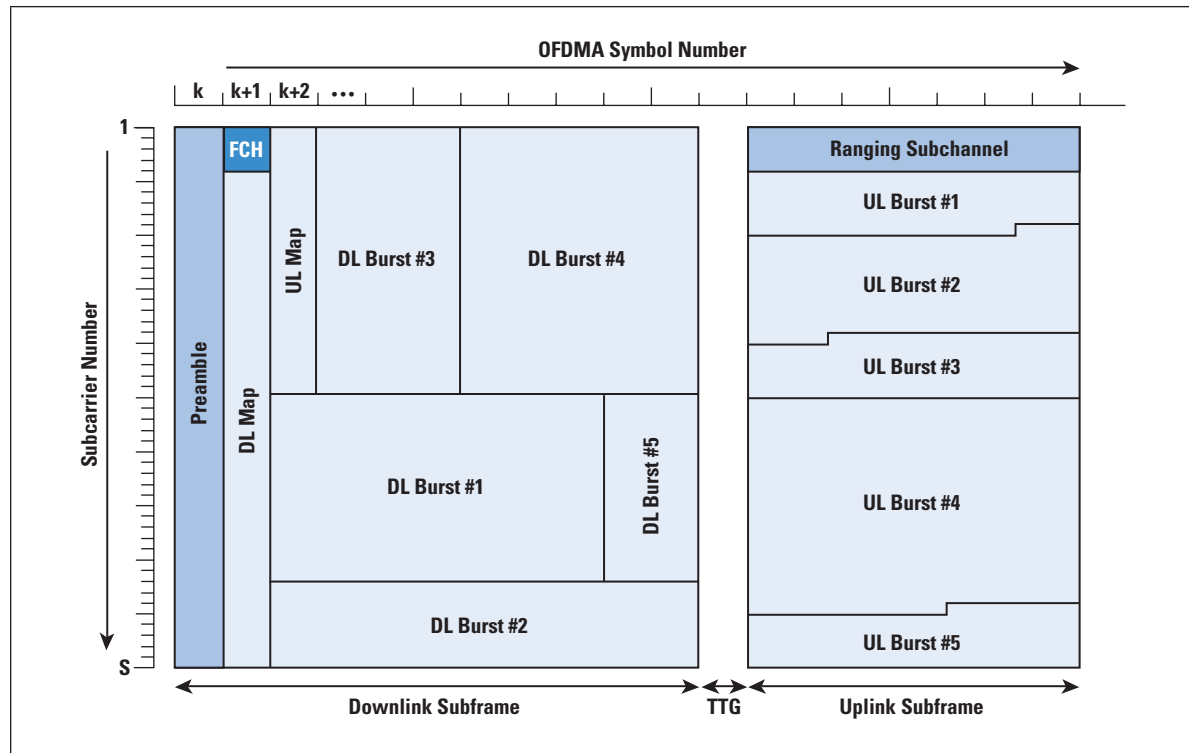
WiMAX systems can use *Time-Division Duplexing* (TDD) or *Frequency-Division Duplexing* (FDD) when allocating air interface resources to users. In TDD, the uplink and downlink transmissions are done over the same carrier frequencies and the separation between the transmission directions is done by assigning time slots, in which the transmission to one direction at a time is scheduled. In FDD, uplink and downlink transmissions are done simultaneously over different carrier frequencies.

Commonly used in mobile WiMAX equipment, TDD allows more flexible sharing of the available bandwidth between the uplink and downlink transmissions. On balance, TDD requires synchronization between multiple adjacent base stations so that transmissions in neighboring cells do not interfere with each other. A TDD frame (Figure 2) is divided into two subframes: first comes a downlink frame and after a short guard interval, called the *Transmit/Receive Transition Gap* (TTG), an uplink frame follows in the same frequency band. Each downlink subframe starts with a preamble, which is used for synchronization and channel estimation. To enhance tolerance against mobility-inflicted channel impairments, WiMAX allows optional support for a more frequent preamble repetition during transmission. In the uplink, short preambles, also called *midambles*, can be used after 8, 16, or 32 OFDM symbols, and in the downlink, short preambles in front of every data burst can be used. After the preamble comes a *Frame Control Header* (FCH), which consists of uplink and downlink *Media Access Protocol* (MAP) messages, which inform users about their transmission parameters.

Flexible data multiplexing from different users into one OFDM or OFDMA frame is also supported, as illustrated in Figure 2. Both uplink and downlink subframes can include data bursts of different types from multiple users, and they can be of variable length.

A small portion of the uplink subframe is reserved for transmission parameter adjustment and bandwidth request purposes. Moreover, small amounts of user data can be sent in this portion of the uplink subframe. The total OFDM frame size can range between 2.5 and 20 ms, but the initially supported frame size in present WiMAX equipment is 5 ms.

Figure 2: An example of a WiMAX OFDMA Frame

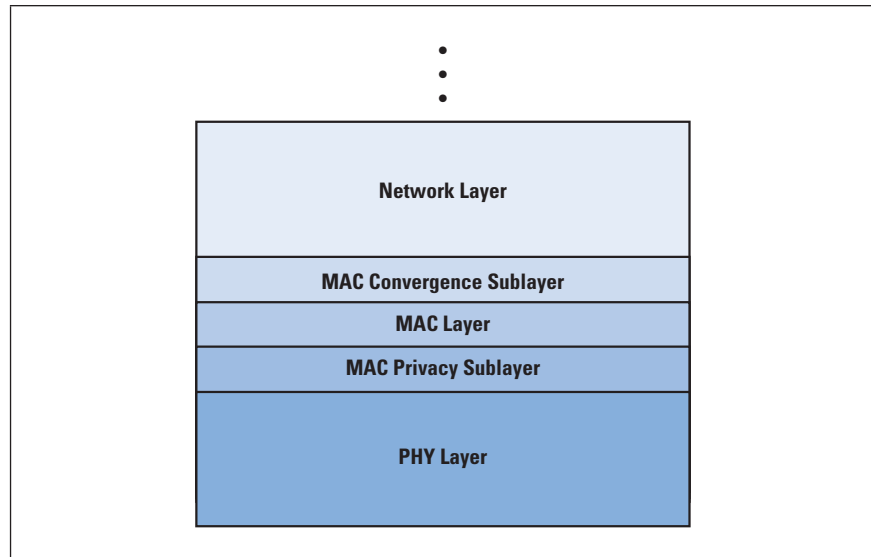


Media Access Control

The MAC layer is primarily an adaptation layer between the PHY layer and the upper layers. Its most important task, when transmitting data, is to receive *MAC Service Data Units* (MSDUs) from the layer above, aggregate and encapsulate them into *MAC Protocol Data Units* (MPDUs), and pass them down to the OFDM or OFDMA PHY layer for transmission. When data is received, the MAC layer takes MPDUs from the PHY layer, decapsulates and reorganizes them into MSDUs, and passes them on to the upper-layer protocols.

An additional layer between the MAC and upper protocol layers called the *Convergence Sublayer* (CS) is also defined in [1] and [2] and illustrated in Figure 3. For the upper layers, CS functions as an interface to the MAC layer. Even though in principle a CS is presented for a variety of different protocols, currently [3] and [4] support CS only for IP and Ethernet. Other protocols can, of course, use these CSs through encapsulation. The CS may also support upper-protocol header compression.

Figure 3: WiMAX Protocol Stack



Similarly with the PHY layer, shown in Figure 3, the MAC layer allows flexible allocation of transmission capacity to different users. Variably sized MPDUs from different flows can be included into one data burst before being handed over to the PHY layer for transmission. Multiple small MSDUs can be aggregated into one MPDU and, conversely, one big MSDU can be fragmented into multiple small ones in order to further enhance system performance. For example, by bundling up several MPDUs or MSDUs, the PHY and MAC layer header overheads, respectively, can be reduced.

It is important to remember that the BS MAC layer manages bandwidth allocation for both uplink and downlink transmissions. The BS assigns bandwidth for the downlink transmission according to incoming network traffic. For the uplink transmission, bandwidth is allocated based on the requests received from the MS. Because basically all connections are controlled by the BS, QoS can be efficiently implemented into WiMAX equipment. Currently, the MAC layer of a mobile WiMAX BS should include support for five different QoS classes, briefly summarized in Table 1.

Table 1: Mobile WiMAX QoS Classes

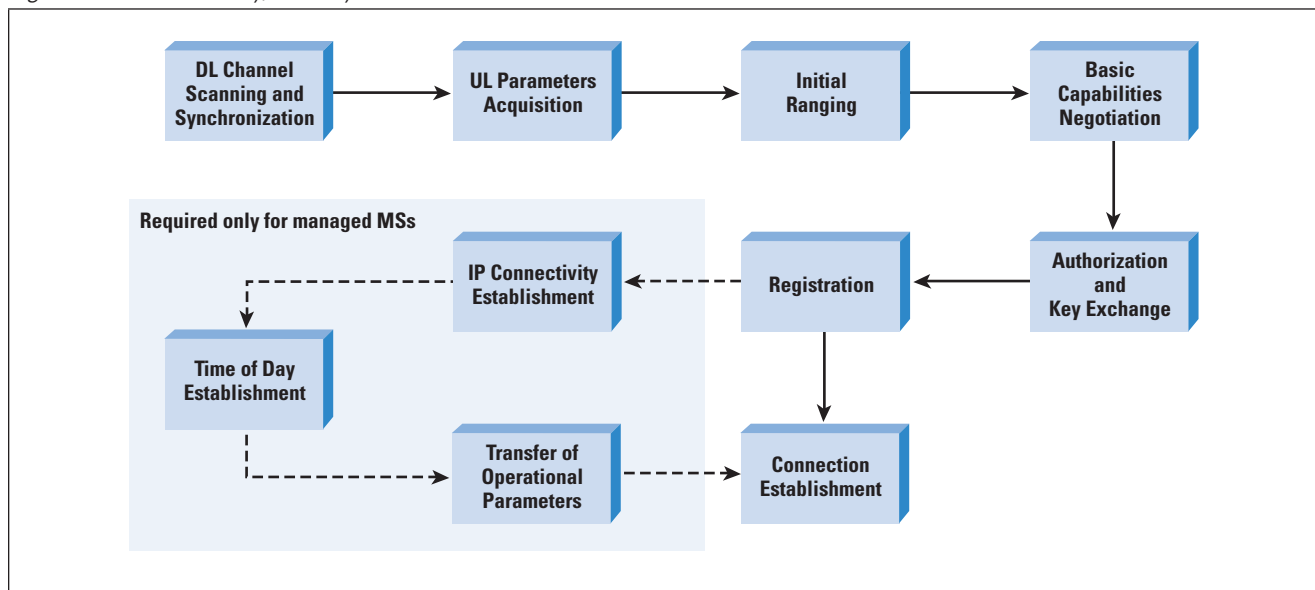
QoS Class	Supported Service	Example Application
Unsolicited Grant Services (UGS)	Latency- and jitter-sensitive applications with fixed-size data packets at Constant Bitrate (CBR)	Voice over IP (VoIP) without silence suppression
Real-Time Variable Rate (RT-VR)	Real-time applications with variable-size data packet bursts	Video and audio streaming
Non-Real-Time Polling Services (nrtPS)	Delay-tolerant applications with variable-size data packets and guaranteed bitrate demands	File transfers
Extended Real-Time Variable Rate (ERT-VR)	Real-time applications with Variable Bitrate (VBR) data streams and guaranteed bitrate and delay demands	VoIP with silence suppression
Best Effort (BE)	Data streams with no minimum service-level demands	Web browsing, instant messaging, and data transfer

Prior to any data transmission over a WiMAX link, the MS and the BS must form a unidirectional connection between their respective MAC layers. A unique identifier, called *Connection Identifier* (CID), is assigned to each uplink and downlink connection pair. The CID serves as a temporary address for the transmitted data packets over the WiMAX link. Another identifier, called *Service Flow Identifier* (SFID), is assigned by the BS to unidirectional packet flows with the same QoS parameters, that is, service flows. The BS also handles the mapping of SFIDs to CIDs in the QoS control process. Note that the MAC layer incorporates sophisticated power-management techniques and robust, state-of-the-art security features, but these features are out of scope for this article.

Network Entry and Reentry

Figure 4 illustrates the basic steps that every MS must go through when entering or reentering a WiMAX network. First, a MS scans the downlink channel and synchronizes with the BS, after which the MS acquires the transmit parameters for the uplink transmission from the BS *Uplink Channel Descriptor* (UCD) message and performs initial ranging, hence acquiring the correct timing offset and power adjustments. A MS extracts an initial ranging-interval time slot from an uplink MAP message. If a MS cannot complete the initial ranging successfully, it must start scanning for a new downlink channel.

Figure 4: Network Entry/Reentry Procedure



The basic capabilities negotiation process starts when the MS sends a message containing its capabilities to the BS; the BS responds with a message containing the capabilities it has in common with the MS. If *Privacy Key Management* (PKM) is enabled at both the MS and the BS, the next step is to perform the authorization and key-exchange procedure, so that the MS can register with the network. The BS sends back a registration response message that contains the secondary management CID, if the MS is managed.

After a managed MS obtains this secondary management CID, it becomes “manageable.” The successful reception of the registration response message is a prerequisite for any MS in order to be able to transmit to and receive from the network.

When a managed MS enters the network, the next step is to establish IP connectivity by using the assigned secondary management connection and by either invoking the *Dynamic Host Configuration Protocol* (DHCP)^[10] or DHCPv6^[11], or using the IPv6 stateless address autoconfiguration^[12], depending on the information provided by the BS registration response message. If the MS uses MIPv4 or MIPv6, it can secure its address by using the secondary management connection with MIP. The establishment of IP connectivity and time of day, as well as the transfer of the operational parameters, are needed only for managed MSs. These parameters can be managed with IP management messages through a secondary management connection, for example, by using the DHCP, *Trivial File Transfer Protocol* (TFTP)^[13], or *Simple Network Management Protocol* (SNMP)^[14]. These additional steps during network entry are necessary for the operation of the IP management protocols.

If DHCP is used to establish IP connectivity, a managed MS must also establish the time of day so that the management system can time-stamp certain events. Both the MS and the BS must be set at the same time of day, with an accuracy of the nearest second. The time of day is retrieved using the secondary management connection with the *Time Protocol*^[15]. The current time is formed by combining the time retrieved from the server with the time offset extracted from the DHCP reply message. Although the time of day is not needed for the registration to complete successfully, it is required in order to keep the connection operational. Finally, the managed MS must acquire its operational parameters with TFTP.

After a managed MS has obtained its operational parameters, or after an unmanaged MS has registered with the network, the MS preprovisioned service-flow connections are established.

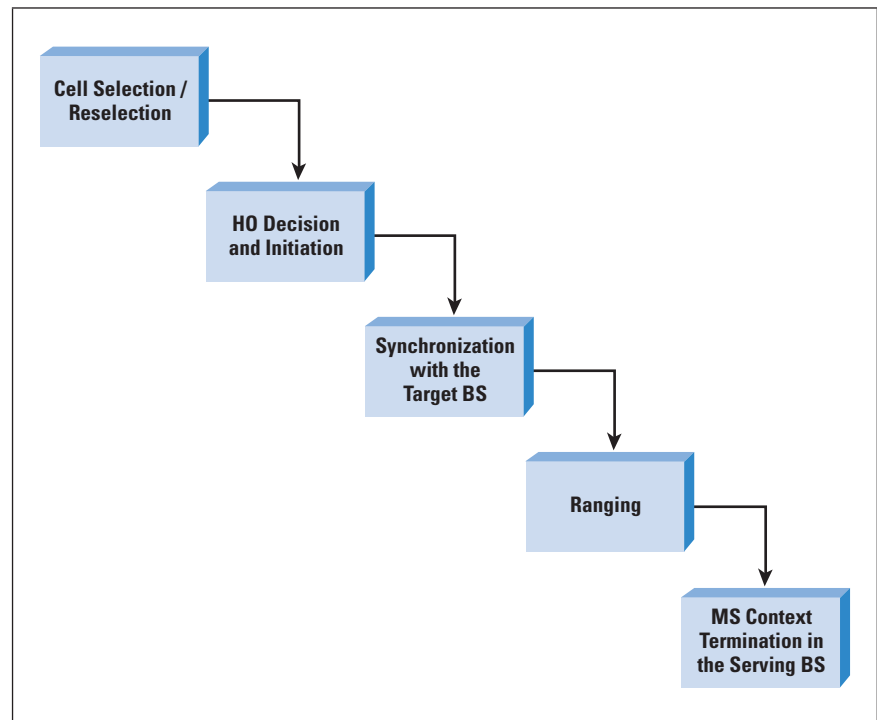
Mobility Support

As discussed previously, IEEE 802.16e introduced mobility support, defining an OFDMA PHY layer and signaling mechanisms to enable location and mobility management, paving the way for mobile WiMAX. The WiMAX Forum details four mobility scenarios in addition to the fixed WiMAX scenario. In the nomadic and portable mobility scenarios, the point of attachment of a fixed *Subscriber Station* (SS) can change. The simple mobility scenario allows MSs to roam within the coverage area with speeds up to 60 km/h, but handovers may cause connection interruptions of up to 1 second. In the so-called *full-mobility scenario*, the MS speed can be as much as 120 km/h, and transparent handovers are supported. This last scenario is what many might consider as the real mobile WiMAX scenario, but all five scenarios are “standards-compliant.”

Although three different types of handovers are defined in [2], *Hard Handover* (HHO), *Macro Diversity Handover* (MDHO), and *Fast Base Station Switching* (FBSS), only HHO is mandatory for all mobile WiMAX equipment. This type of handover is often referred to as a *break-before-make handover*: first, the MS disconnects from the serving BS and then connects to the target BS. Because of the short disconnection period, packets may be lost; HHO is less sophisticated than either MDHO or FBSS and may be inappropriate for some applications. The MS must also register with the target BS and reauthenticate with the network, typically meaning further delays before actual data exchange can (re)start. If multiple handover types are supported and enabled, the BS decides which type should take precedence over the other. MDHO and FBSS are enabled or disabled during the registration of the MS with the BS.

Figure 5 illustrates the five stages of a successful HHO in mobile WiMAX. The first stage is to select the target BS cell based on information about the network topology surrounding the serving BS through periodically broadcasted neighbor advertisements. The advertisements include the same information on the serving BS neighbors that the *Downlink Channel Descriptor* (DCD) and *Uplink Channel Descriptor* (UCD) messages of the neighboring BSs would include. For example, a neighbor advertisement message includes channel information of the neighboring BSs so that the MS can synchronize with them and perform scanning operations to evaluate their suitability as potential targets for a HO.

Figure 5: The Five Phases of a Successful HHO



The second phase is to make the actual decision to initiate the handover procedure, when a certain network (say, congestion in the serving cell requires load balancing) or channel condition threshold (for example, low received *Signal-to-Interference + Noise Ratio* [SINR] in the current cell) is crossed. The actual decision to start the message exchange for the MS to migrate from the radio interface of the serving BS to the radio interface of another BS can be made by the MS, BS, or the network. In the third phase, the MS synchronizes with the downlink transmission of the target BS and obtains the transmission parameters for the downlink and the uplink. The time consumed to perform the synchronization procedure depends on the amount of information the MS received about the target BS in the neighbor advertisement messages prior to the handover. The average synchronization latency without previously acquired information about the target BS ranges from two to three frame cycles, or approximately 4 to 40 ms depending on the OFDMA frame duration used in the system. The more extensive the channel parameter list received in the neighbor advertisement messages prior to the handover, the shorter the time to achieve the synchronization.

After synchronizing, the MS and the target BS initiate the ranging procedure. During this fourth step in HHO, MS and BS exchange the required information so that the MS can reenter the network. The target BS can request information about the MS from the (previously) serving BS and other network entities. Again, the more information made available to the target BS, the shorter the time to reenter the network, because the target BS may skip some steps from the network (re)entry procedure described earlier. In short, sharing context information before the actual handover optimizes the handover procedure and decreases its latency. In the last step of a HHO, the MS context at the serving BS is terminated and resources are released.

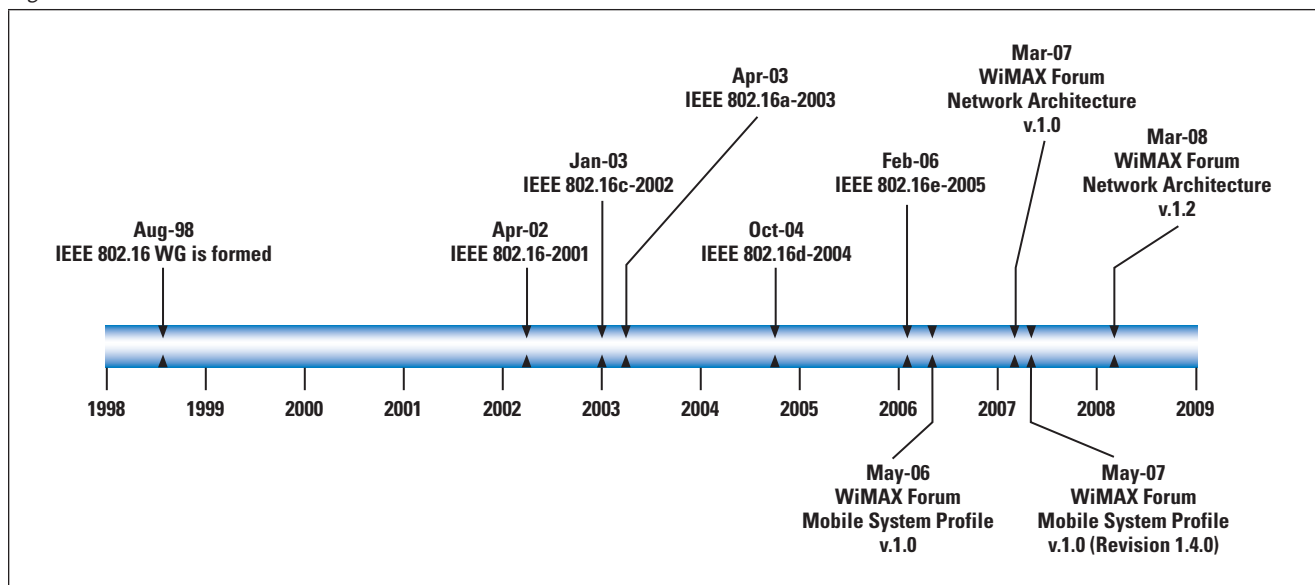
If MDHO and FBSS are supported, the following stages, in addition to those already described in the HHO procedure, must be performed: (a) decision to enable MDHO or FBSS, (b) diversity set update, and (c) anchor BS selection. In macrodiversity communications the MS maintains a connection to one or more serving BSs simultaneously, enabling soft or make-before-break handovers. In [2], the transition of the MS from the air interface of one or more serving BSs to the air interface of one or more target BSs is referred to as a MDHO. The MS and the BS both maintain a list called the *diversity set*, which includes all serving BSs involved in the MDHO communication. The MS maintains both uplink and downlink unicast connections to all the BSs in the diversity set, and one of the serving BSs is defined as the anchor BS. Note that all BSs involved in the diversity set use the same set of CIDs for the connections established between the MS and the serving BSs.

In FBSS, the MS transmits to and receives data from a single serving BS during any frame period. The BS, to which the MS has the connection to at any given frame, is called the *anchor BS*. The MS maintains a diversity set, which includes all active BSs in its range, and can change its anchor BS on a frame-by-frame basis, based on certain criteria. The transition from the serving anchor BS to the target anchor BS in FBSS is done without invocation of the normal handover procedure, and only the anchor BS update procedure is needed. After all, the MS has collected all required information about all BSs during the diversity set update ranging procedures.

Mobile WiMAX vs. HSPA

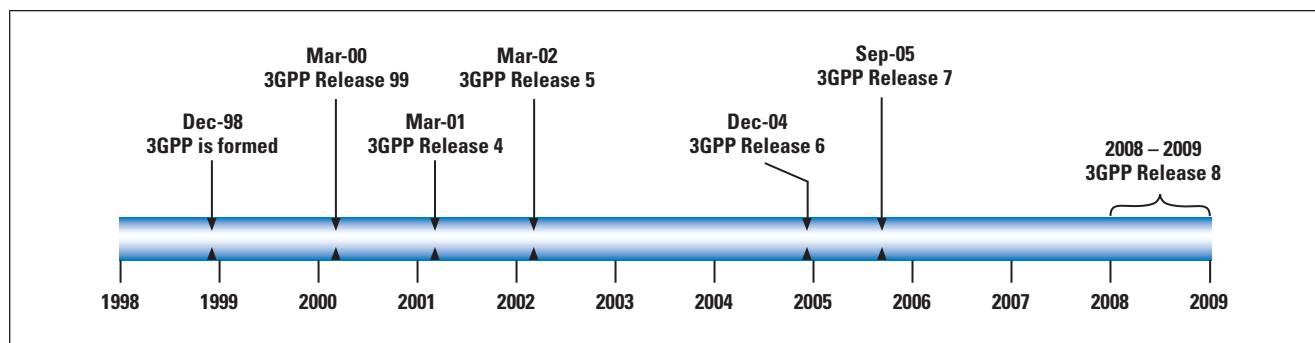
Mobile WiMAX and *High-Speed Packet Access* (HSPA) are expected to be the two major contestants in the rapidly growing wireless broadband market. The two, however, come from different origins. Figure 6 summarizes the evolution toward mobile WiMAX. It all started with the establishment in August 1998 of the IEEE 802.16 working group, which published its first standard (IEEE 802.16-2001) in April 2002. This first version defines a single carrier system operating in the 10- to 66-GHz frequency band and only under *line-of-sight* (LOS) conditions. The IEEE 802.16c-2002 amendment detailed system profiles for the original standard based on the 10- to 66-GHz frequency band. IEEE 802.16a-2003 introduced support for 2- to 11-GHz frequencies and *non-line of sight* (NLOS) operation, and adopted the use of OFDM and OFDMA. IEEE 802.16d-2004^[1] consolidated all these previous versions and amendments in a single document, and further enhanced the system. Fixed WiMAX is based on IEEE 802.16d-2004, [3], and [4]. Mobile WiMAX is based on the IEEE 802.16e-2005 amendment^[2], which introduced mobility support, as well on [3] and [4].

Figure 6: The Road Toward Mobile WiMAX



HSPA is a set of technological enhancements to the already widely deployed *Wideband Code Division Multiple Access* (WCDMA) cellular networks defined by the *Third Generation Partnership Project* (3GPP). Figure 7 illustrates the WCDMA specification evolution. The origins of HSPA can be traced in the foundation of 3GPP in December 1998. The original aim of 3GPP was to develop a third-generation WCDMA system, and in the process, HSPA was introduced. In March 2000, Release 99, the original standard specifying the WCDMA system, was published. A year later, the first enhancements were published in Release 4, which introduced, among others, an IP-based core network. Release 5 introduced *High-Speed Downlink Packet Access* (HSDPA) and defined the *3GPP IP Multimedia Subsystem* (IMS). *High-Speed Uplink Packet Access* (HSUPA) and some further improvements to HSDPA were defined in Release 6 (December 2004). Release 7 further enhanced QoS support and defined mechanisms to decrease network latency. Release 8 is expected to be published in 2008, and it will include specifications for the next step, called *3GPP Long-Term Evolution* (LTE). LTE is meant to deliver maximum cell throughputs an order of magnitude larger than HSPA.

Figure 7: The Evolution of the 3GPP WCDMA Standard



Mobile WiMAX evolved out of a broadband wireless LAN/MAN technology, and vendors currently report that it can deliver maximum cell capacities of 46 and 7 Mbps in downlink and uplink transmissions, respectively. However, mobility management is a later addition and, according to Maravedis, by September 2007 only 12 percent of all deployed *Customer Premises Equipment* (CPE) was IEEE 802.16e-2005-compliant^[16]. On the other hand, HSPA is based on a solid foundation of mobility management techniques with wide deployment in cellular networks around the globe, but can currently deliver maximum cell throughputs of only 14.4 and 5.8 Mbps in downlink and uplink transmissions, respectively.

Either commercial or trial networks of both technologies have already been implemented all over the world. However, according to the *Global Mobile Suppliers Association* (GSA), HSPA networks have yet to be deployed in China and India, both of which are large and rapidly growing market areas for wireless communications. According to Maravedis, both India and China have at least WiMAX trial deployments in place.

As mentioned already, the vast majority of current WiMAX deployments do not support mobility. Up to now, fixed WiMAX has been used mainly for last-mile broadband connectivity for sparsely populated rural areas. The largest commercial IEEE 802.16e-2005-compliant system is currently the *Wireless Broadband (WiBro)*^[17] network in South Korea, which supports simple mobility up to 60 km/h. Even though WiMAX and WiBro are both based on the same standards, WiBro was developed by the South Korean telecommunications industry before the WiMAX Forum adopted mobility support for its system profiles. WiMAX and WiBro are often cited as separate technologies, even though cooperation is in place in order to assure interoperability between the two.

Summary

In this article we presented an overview of mobile WiMAX, a much-heralded technology for next-generation mobile broadband networks; mobile WiMAX is an intricate system. We introduced WiMAX and the role of the WiMAX Forum, and summarized the important points of the WiMAX network reference model and the PHY and MAC layers. We addressed mobility support, but not the security aspects. Finally, we briefly compared WiMAX with HSPA, presenting their respective evolutions and illustrating their worldwide deployments. We hope that this article will serve as a valuable primer, and we highly recommend that those interested in the mobile WiMAX technology check the bibliography.

Bibliography

- [1] IEEE 802.16 Working Group, “IEEE Standard for Local and Metropolitan Area Networks, Part 16: Air Interface for Fixed Broadband Wireless Access Systems,” IEEE Standard 802.16-2004, October 2004.
- [2] IEEE 802.16 Working Group, “IEEE Standard for Local and Metropolitan Area Networks, Part 16: Air Interface for Fixed and Mobile Broadband Wireless Access Systems, Amendment 2: Physical and Medium Access Control Layers for Combined Fixed and Mobile Operation in Licensed Bands,” IEEE Standard 802.16e-2005, February 2006.
- [3] WiMAX Forum Network Working Group, “WiMAX Forum Network Architecture—Stage 2: Architecture Tenets, Reference Model and Reference Points—Release 1, Version 1.2,” WiMAX Forum, January 2008.
- [4] WiMAX Forum Network Working Group, “WiMAX Forum Network Architecture—Stage 3: Detailed Protocols and Procedures—Release 1, Version 1.2,” WiMAX Forum, January 2008.

- [5] K. Pentikousis, "Wireless Data Networks," *The Internet Protocol Journal*, Volume 8, No. 1, March 2005, pp. 6–14.
- [6] IEEE 802.16 Working Group, "IEEE Standard for Local and Metropolitan Area Networks. Part 16: Air Interface for Fixed and Mobile Broadband Wireless Access Systems. Amendment 3: Management Plane Procedures and Services," IEEE Standard 802.16g-2007, December 2007.
- [7] C. Perkins (Ed.), "IP Mobility Support for IPv4," RFC 3344, August 2002.
- [8] D. Johnson, C. Perkins, and J. Arkko, "Mobility Support in IPv6," RFC 3775, June 2004.
- [9] K. Leung, G. Domemety, P. Yegani, and K. Chowdhury, "WiMAX Forum/3GPP2 Proxy Mobile IPv4," Internet-Draft, Work in Progress.
- [10] R. Droms, "Dynamic Host Configuration Protocol," RFC 2131, March 1997.
- [11] R. Droms (Ed.), J. Bound, B. Volz, T. Lemon, C. Perkins, and M. Carney "Dynamic Host Configuration Protocol for IPv6 (DHCPv6)," RFC 3315, July 2003.
- [12] S. Thomson and T. Narten, "IPv6 Stateless Address Autoconfiguration," RFC 2462, December 1998.
- [13] K. Sollins, "The TFTP Protocol (Revision 2)," RFC 1350, July 1992.
- [14] J. Case, M. Fedor, M. Schoffstall, and J. Davin "A Simple Network Management Protocol (SNMP)," RFC 1157, May 1990.
- [15] J. Postel and K. Harrenstien, "Time Protocol," RFC 868, May 1983.
- [16] K. Pentikousis, J. Pinola, E. Piri, F. Fitzek, T. Nissilä, and I. Harjula, "Empirical Evaluation of VoIP Aggregation over a Fixed WiMAX Testbed," Proceedings of The 4th International Conference on Testbeds and Research Infrastructures for the Development of Networks & Communities (TRIDENTCOM), 18–20 March, 2008, Innsbruck, Austria.
- [17] Telecommunications Technology Association, "Specifications for 2.3GHz Band Portable Internet (WiBro™) Service," TTA Standard TTAS.KO-06.0082/R1, December 2005.

JARNO PINOLA received his M.Sc. from the University of Oulu, Oulu, Finland, in Spring 2008. During his studies, he specialized in telecommunication systems and wrote his Master's Thesis on mobility management issues in wireless broadband systems. Currently he is working as a Research Scientist at VTT Technical Research Centre of Finland in Oulu, Finland. He can be contacted via e-mail at: [**jarno.pinola@vtt.fi**](mailto:jarno.pinola@vtt.fi)

KOSTAS PENTIKOUSIS studied computer science at Aristotle University of Thessaloniki, Thessaloniki, Greece (B.Sc. 1996), and Stony Brook University, Stony Brook, New York, USA (M.Sc. 2000, Ph.D. 2004). He is a tenured Senior Research Scientist at VTT Technical Research Centre of Finland, in Oulu, Finland. He has published internationally in several areas, including mobile computing (mobility triggers, multiaccess, media-independent handovers, and energy consumption); transport protocols; applications; network traffic measurements and analysis; and simulation and modeling. Visit [**http://ipv6.willab.fi/kostas**](http://ipv6.willab.fi/kostas) for more information and contact details.

IDNs

The DNS protocol is 8-bit clean (“Internationalizing the Domain Name System,” IPJ, Volume 11, No. 1, March 2008), even if some DNS clients and servers are not. The hardest thing about changing any Internet protocol is coordinating clients and servers during the transition.

And yet, with the DNS, no transition is needed to support UTF-8 domain names. If you want to publish a UTF-8 domain name, then run a name server that supports UTF-8. If you want to be able to access domain names in your own language, switch to DNS software that supports it. Implementations that are 8-bit clean are already available; ordinary market mechanisms will handle the rest.

Punycode is a gross hack that makes my stomach roil. You know it, I know it, any engineer will agree with you, so how did it get through the IETF?

The argument for where to stop internationalization does not spread to `protocol://` because it’s “gobble-de-gook” in English, too. Dots are a completely arbitrary character used to separate the hierarchy. There’s plenty of space at the top for UTF-8 names.

The real problem with IDN is homoglyphs.

—Russ Nelson,
nelson@crynwr.com

The author responds:

It would certainly make more sense in terms of design elegance and minimalism within the DNS if the label that was stored in the DNS was precisely the same label that was used in the interface between applications and the DNS client software. There is something rather clumsy about the approach that stores an encoded version of a canonical version of the label value, and relies on the application being capable of performing the *stringprep* and encoding functions in consistent and uniform ways. The resultant limitations on what can actually sit in DNS labels on a language-by-language basis are, in part, an outcome of the potential indeterminism of this canonicalization function.

But indeterminism is not a tolerable outcome of the DNS. The DNS is not a guessing game, and inconsistencies in the mapped transforms that are provided by the DNS trigger intolerable insecurities in the networked environment. So the *nameprep* profiles and the related restrictions on allowable Unicode code points are unavoidable if we want to avoid this indeterminism in the DNS.

So if *nameprep* is required in any case, then what we are left with to consider is the decision to use the Punycode *ASCII Compatible Encoding* (ACE) to map Unicode labels into the *Letter-Digit-Hyphen* (LDH) subset of ASCII. But is the Punycode ACE really that much of a problem? Within the overall IDN framework the Punycode algorithm is not so complex that the risk of incorrect implementations is significant, the algorithm is not processor-intensive, and the outcome does not inflate the encoded labels to an impossible length. The advantage of Punycode is that the DNS servers do not require modification, and the clients that manipulate IDNs required additional *nameprep* functions in any case, so Punycode was evidently intended to be the least-impact approach that spared DNS servers from a potential requirement for modification.

To me, this solution appears to be a design tradeoff, in so far as the ACE approach circumvents the observed problem of non-8-bit clean DNS servers sitting within the deployed DNS, and does not in and of itself demand novel roles and functions on the part of the clients of the DNS in addition to what was already necessitated by the IDN *nameprep* function. However, at the same time it creates an annoying inconsistency in the overall framework of the design of the DNS, where certain labels in the DNS are intended to trigger a Punycode transform into an equivalent Unicode string while other labels are meant to be used without further transforms applied.

My judgment of the short-term path of least risk sits with the ACE approach as adopted for IDNs, but at the same time I agree with Russ' discomfort that the path that preserves the long-term essential broad utility and function of the DNS through consistency of design and application sits in an 8-bit clean DNS without the adornment of any form of an ACE.

And, yes, I agree with Russ that the most significant problem with IDNs is homoglyphs, because of continued reliance of an underlying approach of "appearance is everything" in terms of the integrity of the DNS as an identity framework.

—Geoff Huston,
gih@apnic.net

More IDNs

The LDH restriction referred to in "Internationalizing the Domain Name System" (IPJ, Volume 11, No. 1, March 2008) was relaxed in RFC 1123^[1] to allow a host name to begin with either a letter or a digit.

—Andrew Friedman

[1] R. Braden, Editor, "Requirements for Internet Hosts—Application and Support," RFC 1123, October 1989.

The author responds:

My thanks to Andrew for pointing this out. It has been commonly recounted that this relaxation of the LDH convention was associated with the successful registration of the DNS name **3com.com** and that the RFC paperwork was revised following this registration. Since then the most visible set of names that used this “liberal” revision of LDH with names that have leading digits were telephone number mapping name sets, including the venerable **tpc.int** domain of the early 1990s and, more recently, ENUM. As for names with leading hyphens, I don’t believe that we are at the point of allowing Morse code into the DNS yet, but I’m sure that someone somewhere is working on it!

—Geoff
(-- . . --- .-. .-.)

We want to hear from You

Your feedback is important to us. Please send your comments and suggestions to **ipj@cisco.com**. And don’t forget to visit our Website at **<http://www.cisco.com/ipj>** where you can read or download back issues, update and renew your subscription, and find articles using our index files. We also encourage you to participate in our online forum at **<http://ipjforum.org>**

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Fragments

OLSR stands for *Optimized Link State Routing Protocol*.

DUMBO

The *Digital Ubiquitous Mobile Broadband OLSR* (DUMBO) project deploys mobile wireless networks on an ad hoc basis for emergency conditions, such as after a natural disaster when a fixed network infrastructure is not available.

A *Mobile ad hoc Network* (MANET) consists of mobile nodes that automatically cooperate to support the exchange of information through wireless medium. Since the MANET does not rely on fixed telecommunication infrastructure, it is suitable for emergency situations and can be set up in a short amount of time. Using lightweight portable mobile nodes, MANET coverage can penetrate deep into areas not easily accessible by roads or into areas where the telecommunication infrastructure has been destroyed.

DUMBO allows streaming video, *Voice over IP* (VoIP) and short messages to be simultaneously transmitted from a number of mobile laptops to a central command center, or to the other rescuers at the same or different disaster sites. The DUMBO command center has a face recognition module that identifies potential matches between unknown victims' face photos taken from the field and a collection of stored known face images. In addition, sensors can be deployed to measure environmental data such as temperature and humidity. Data from the sensors can be sent to the command center which analyzes or passes it on to the other mobile nodes. The command center can be located either in the disaster area or anywhere with Internet access. DUMBO technology is currently being deployed in cyclone-ravaged Burma. See <http://www.interlab.ait.ac.th/dumbo/> and <http://www.relief.asia/>

Upcoming Events

The *Internet Engineering Task Force* (IETF) will meet in Dublin, Ireland, July 27 – August 1 and in Minneapolis, Minnesota, November 16 – 21, see <http://www.ietf.org/>

APNIC, the *Asia Pacific Network Information Centre*, will hold its Open Policy meeting in Christchurch, New Zealand, August 25 – 29, see <http://www.apnic.net/meetings/26/>

[Ed.: I will be organizing a pipe organ demonstration event on August 26 as part of the opening reception for APNIC 26, see <http://organdemo.info>]

The *North American Network Operators' Group* (NANOG) will meet in Los Angeles, California, October 12 – 14. Immediately following the NANOG meeting, the *American Registry for Internet Numbers* (ARIN) will meet in the same location, October 15 – 17. See <http://nanog.org> and <http://arin.net>

The *Internet Corporation for Assigned Names and Numbers* (ICANN) will meet in Paris, France, June 22 – 26, and in Cairo, Egypt, November 2 – 7. See <http://icann.org>

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2008 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA
--

The Internet Protocol Journal

September 2008

Volume 11, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
GMPLS and the Optical Internet	2
IPv4 Address Exhaustion.....	19
Letters to the Editor.....	37
Book Reviews	39
Fragments	46

FROM THE EDITOR

If you are reading the printed version of this journal you will notice a subtle change in the paper. This issue is printed on an uncoated stock, specifically Exact® Offset Opaque White 60#, a recycled paper made by Wausau Paper Corporation. This paper is slightly thinner, and thus lighter, than the paper we have been using. It is also less reflective and easier to write notes on. We invite your feedback on this paper as we experiment with various solutions to reduce our carbon footprint. As always, send your comments to: ipj@cisco.com

This journal has a long history of covering existing and emerging technologies that form part of the underlying infrastructure for both the global Internet and private enterprise networks. Recent articles have focused on wireless systems such as WiMAX, and we have other articles on wireless technologies in the pipeline. This time, however, we look at *optical networking*, specifically *Generalized Multiprotocol Label Switching* (GMPLS) as a technology for next-generation internets. The article is by Francesco Palmieri.

The topic of IP Version 4 address exhaustion has been discussed in several articles in this journal, and is currently being heavily debated in the *Regional Internet Registries* (RIRs). As we approach the inevitable date when the IPv4 address pool “runs out,” we are returning to this topic with several articles. The first of these articles is included in this issue. Geoff Huston sets the stage by reviewing some of the history and answering the basic question of “why” we find ourselves at a point in history where the IPv4 addresses will run out before we have deployed any significant amount of IPv6 systems. In future issues we will follow Geoff’s introduction with several other perspectives on this situation.

Once again, let me remind you to visit our Website at <http://www.cisco.com/ipj>, where you can renew and update your subscription, download back issues, and find additional resources such as our online forum at <http://ipjforum.org>

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

GMPLS Control Plane Services in the Next-Generation Optical Internet

by Francesco Palmieri, Federico II University of Napoli, Italy

One of the major concerns in the Internet-based information society today is the tremendous demand for more and more bandwidth. Optical communication technology has the potential for meeting the emerging needs of obtaining information at much faster yet more reliable rates because of its potentially limitless capabilities—huge bandwidth (nearly 50 terabits per second^[1]), low signal distortion, low power requirement, and low cost. The challenge is to turn the promise of optical networking into reality to meet our Internet communication demands for the next decade. With the deployment of *Dense Wavelength Division Multiplexing* (DWDM) technology, a new and very crucial milestone is being reached in network evolution. The speed and capacity of such wavelength switched networks—with hundreds of channels per fiber strand—seem to be more than adequate to satisfy the medium to long term connectivity demands. In this scenario, carriers need powerful, commercially viable and scalable devices and control plane technologies that can dynamically manage traffic demands and balance the network load on the various fiber links, wavelengths, and switching nodes so that none of these components is over- or underused.

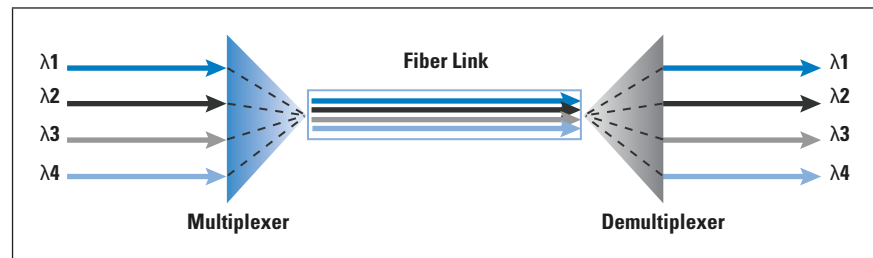
This process of adaptively mapping traffic flows onto the physical topology of a network and allocating resources to these flows—usually referred to as *traffic engineering*—is one of the most difficult tasks facing Internet backbone providers today. *Generalized Multiprotocol Label Switching* (GMPLS) is the most promising technology. GMPLS will play a critical role in future IP pure optical networks by providing the necessary bridges between the IP and optical layers to deliver effective traffic-engineering features and allow for interoperable and scalable parallel growth in the IP and photonic dimension. The GMPLS control plane technology, when fully available in next-generation optical switching devices, will support all the needed traffic-engineering functions and enable a variety of protection and restoration capabilities, while simplifying the integration of new photonic switches and existing label switching routers.

Wavelength Division Multiplexing

Traditional *Electronic Time-Division Multiplexed* (ETDM) networks use an electrical signal form to switch traffic along routes and restore signal strength. These networks do not fully exploit the bandwidth available on optical fibers because only a single frequency (wavelength or *lambda*) of light is used on each fiber to transmit data signals that can be modulated at a maximum bit rate of the order of 40 Gbps. The high bandwidth of optical fibers can be better used through WDM technology by which distinct data signals may share an optical fiber, provided they are transmitted on carriers having different wavelengths^[2].

In more detail, the optical transmission spectrum is divided into numerous nonoverlapping wavelengths, with each wavelength supporting a single communication channel. Each channel, which can be viewed as a *light path*, is transmitted at a different wavelength (or frequency). Multiple wavelengths are multiplexed into a single optical fiber and multiple light-path data is transmitted as shown in Figure 1.

Figure 1: WDM Functional Model



Dense WDM (DWDM), an evolution of WDM referring essentially to the closer spacing of channels, is the current favorite multiplexing technology for long-haul communications in modern optical networks. Hence, all the major carriers today devote significant effort to developing and applying DWDM technology in their business.

All-optical networks employing the concept of WDM and wavelength routing are thought to be the transport networks for the future^[3]. In such networks, two adjacent nodes are connected by one or multiple fibers, each carrying multiple wavelengths or channels. Each node consists of a dynamically configurable optical switch that supports fiber switching and wavelength switching; that is, the data on a specified input fiber and wavelength can be switched to a specified output fiber on the same wavelength^[4]. In order to transfer data between source–destination node pairs, a light path needs to be established by allocating the same wavelength throughout the route of the transmitted data. Benefiting from the development of all-optical amplifiers, light paths can span more than one fiber link and remain entirely optical from end to end. It has been demonstrated that the introduction of wavelength-routing networks not only offers the advantages of higher transmission capacity and routing node throughput, but also satisfies the growing demand for protocol transparency and simplified operation and management^{[3] [5]}.

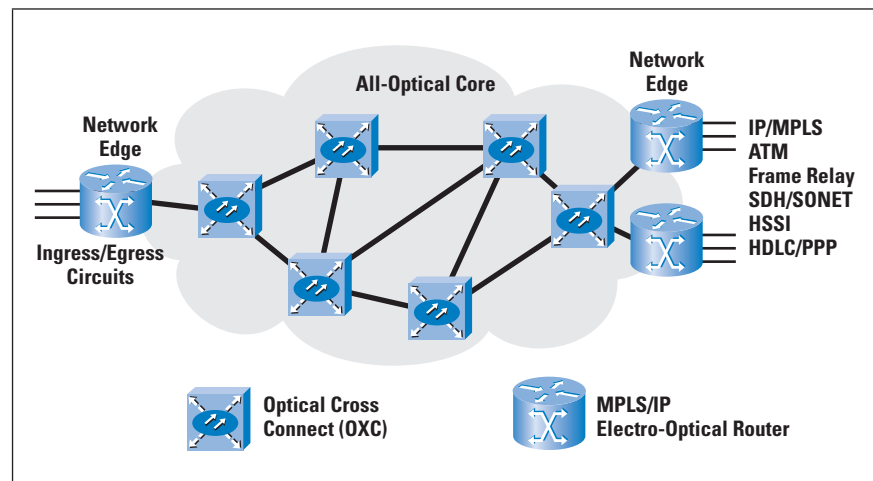
Optical Transport Backbones

The modern Internet transport infrastructure can be physically seen as a very complex mesh of variously interconnected optical or traditional ETDM subnetworks, where each subnetwork consists of several heterogeneous routing and switching devices built by the same or different vendor and operating according to the same control plane protocols and policies. With these very different types of devices, all the forwarding decisions will be based on a combination of packet or cell, timeslot, wavelengths, or physical ports, depending on the position (edge or core) and role (intermediate or termination or gateway node) of the switching devices in the network layout.

In particular, WDM-switched optical subnetworks are typically used as backbone infrastructures to interconnect a large number of different IP as well as other packet networks such as SDH, ATM, and Frame Relay.

New optical devices such as DWDM multiplexers, *Add/Drop Multiplexers* (ADM), and *Optical Cross-Connects* (OXC) are making possible an intelligent all-optical core where packets are routed through the network without leaving the optical domain. The optical network and the surrounding IP networks are independent of each other, and an edge IP router interacts with its ingress switching node only over a well-defined *User-Network Interface* (UNI). Clearly, the optical network is responsible for setting up light paths between the edge IP routers. A light path can be either switched or permanent. Switched light paths are established in real time using proper signaling procedures, and they may last for a short or a long period of time. Permanent light paths are set up administratively by subscription, and they typically last for a very long time. An edge IP router requests a switched light path from its ingress optical switching device using a proper signaling protocol over the UNI. See Figure 2.

Figure 2: The Optical Transport Infrastructure

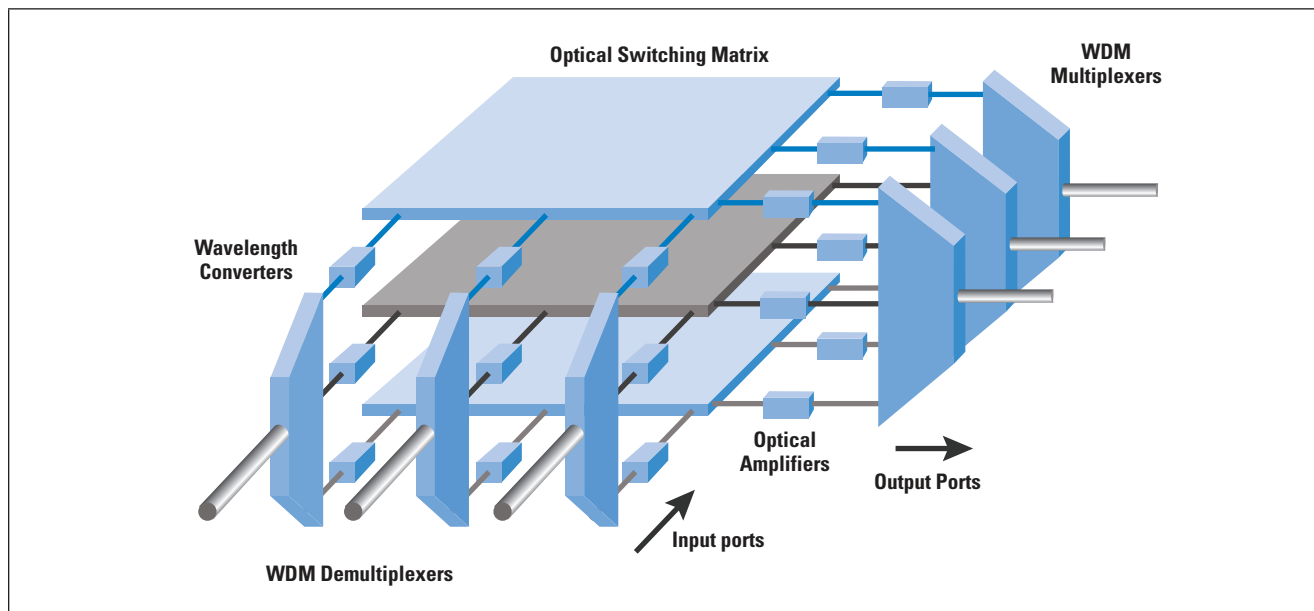


The key concept to guarantee desirable speeds and correct functional behavior in these networks is to maintain the signal in pure optical form, thereby avoiding the prohibitive overhead of conversion to and from electrical form. Such a network would be “optical transparent” in the sense that it would be able to transport client signals with any format and with a wide range of bit rates (at least from about 10 Mbps to more than 10 Gbps). In particular, transparent OXCs, used to selectively switch wavelengths between their input and output ports, are likely to emerge as the preferred option for switching multigigabit or even terabit data streams, because any slow electronic per-packet processing is avoided.

Transparent Optical Switching Nodes

Transparent OXC systems are expected to be the cornerstone of the photonic layer, offering carriers more dynamic and flexible options in building network topologies with enhanced performance and scalability. The development of large and flexible transparent OXCs, now enabled by a new generation of optical components such as optical amplifiers, tunable lasers, and wavelength filters, is still a significant challenge^[1]. Their architecture makes use of optical switching fabrics, wavelength multiplexers and demultiplexers, and transparent wavelength converters, which eliminate the need for optoelectronic transponders. A simple and linear architectural model for an optical transparent OXC is shown in Figure 3.

Figure 3: OXC Architectural Model



Here, the WDM demultiplexers separate incoming grouped wavelengths from input ports into individual lambdas. A sufficiently large low-loss connectivity and compact-design, all-optical switching fabric can be realized by using the reflection of light and *Micro-Electromechanical Systems* (MEMS) technology, now widely available on the market. This multilayer switching fabric driven by a micro-machined electrical actuator redirects, according to the control plane instructions, each wavelength into appropriate output ports passing through optical amplifiers, typically *Erbium-Doped Fiber Amplifiers* (EDFAs), which boost the signal power in line without the need for any optoelectronic conversion to cope with the effects of light dispersion and attenuation on long distances. The WDM multiplexer then groups the wavelengths from the above multiple layers of cross-connects. Furthermore, the wavelength that arrives into an OXC can be directly passed to the optical switching fabric, to be switched to the appropriate output fiber or previously converted, based on the control plane instructions, to another particular wavelength with the use of a tunable wavelength converter (without being transformed to electricity) if the former output wavelength is not available.

This architecture is transparent; that is, the optical signal does not need to be transformed to electricity at all, implying that this architecture can support any protocol and any data rate. Hence, possible upgrades in the wavelength transport capacity can be accommodated at no extra cost. Furthermore, this architecture decreases the cost because it involves the use of fewer devices than the other architectures. In addition, transparent wavelength conversion eliminates constraints on conversions. In this way the real switching capacity of the OXC is increased, leading to cost reduction. First-generation OXCs require manual configuration. Clearly, an automatic switching capability allowing optical nodes to dynamically modify the network topology based on changing traffic demand is highly desirable.

Automatically Switched Optical Networks

For automatically switched networks, where network nodes may directly initiate or terminate new connections or perform wavelength-level switching in the network, sophisticated and flexible control functions are needed.

The *control plane* supports connection management by clients and also provides protection and restoration services. The control plane of an optical network is also responsible for tracking the network topology and for notifying the state of the network resources. Two families of protocols achieve this task:

- *Routing protocols* are specifically responsible for the reliable advertisement of the optical network topology and the available bandwidth resources within and between network domains. In particular, some areas are relevant within this context: the bundling of links with equivalent or logically bundled characteristics, the definition of the routing areas in an optical domain, the rich specifications of an optical link resource as opposed to a typical advertisement of the up or down interface of IP networks, and the advertisement of the shared risk group (optical fibers flowing in the same cable or duct) to which an optical connection belongs.
- *Signaling protocols* are responsible for provisioning, maintaining, and deleting connections. Optical networks are characterized by connection-oriented paradigms that require a resource reservation protocol. State-of-the-art control plane technologies operating on traditional IP-based networks focus on soft-state protocols that require periodic refresh throughout the participating nodes. In optical networks, where the data plane is separated from the control plane, a possible solution is also to adopt a hard state reservation protocol without periodic refresh to limit the effect caused on the data plane by a failure in the control plane. Furthermore, redundant, generalized label binding is encouraged to reserve protection paths in the mesh network.

Data transport is the most obvious task and the main purpose of an optical network *data plane*. It provides uni- or bidirectional information transport (transmission and switching) between users, detects faults, and monitors signal quality. More specifically, the data plane performs, under the directions of the control plane, data routing to the appropriate ports; channel adds and drops to external, older networks (using the edge interfaces); and label or lambda swapping through an array of WDM demultiplexers, wavelength converters, OXCs, optical amplifiers, and multiplexers.

An important concern that must be addressed in designing an optical network is the cross effect of the failure of a data or control plane. Failures of the data plane are usually addressed by the control plane itself by rerouting the disrupted flows at the appropriate level. The control plane must then advertise quickly the new network state to the neighboring nodes to avoid the presence of stale information in the link databases. A failure of the IP-based control plane usually significantly affects the data plane.

Traffic Engineering in Optical Networks

Traffic engineering should be viewed as assistance to the routing and switching infrastructure that provides additional information used in forwarding traffic along alternate paths across the network, trying to optimize service delivery throughout the network by improving its balanced usage and avoiding congestion caused by uneven traffic distribution. Traffic engineering is required in the modern Internet mainly because the current dynamic routing protocols always use the shortest paths to forward traffic. This practice, obviously, conserves network resources, but it causes some of them to be overused while the other resources remain underused. Furthermore, the routing protocols mentioned earlier never account for specific traffic flow requirements such as bandwidth and *Quality of Service* (QoS) needs. Practitioners in the field often assert that traffic engineering essentially signifies the ability to place traffic where the capacity exists to accommodate it—whereas network engineering denotes the ability to install capacity where the traffic exists.

When a traffic-engineering application implements the right set of features, it should provide precise control over the placement of traffic flows within a routing and switching domain, gaining better network use and realizing a more manageable network. A traffic-engineering solution suitable for transparent optical networks always consists of numerous basic functional components; for example:

- *Traffic monitoring, analysis, and aggregation*—This function collects traffic statistics from the network elements; for example, the OXCs. Then the statistics are analyzed or aggregated to prepare for the traffic engineering and network reconfiguration related to decision making.

- *Bandwidth demand projection*—Bandwidth demand projection estimates the bandwidth requirements in the near future based on past and present measurements and the characteristics of the traffic arrival processes. The bandwidth projections are used for subsequent allocation.
- *Reconfiguration trigger*—This variable consists of a set of policies that decide when a network-level reconfiguration is performed. This decision is based on traffic measurements, bandwidth predictions, and operational areas; for example, to suppress the influence of transitional factors and reserve adequate time for the network to converge.
- *Topology design*—Topology design provides a network topology based on the traffic measurements and predictions. Conceptually this process can be considered as optimizing a graph (that is, OXC connected by light paths at the WDM layer) for specific objectives (for example, maximizing throughput), subject to certain constraints (for example, nodal degree or interface capacity), for a given load matrix (that is, traffic load applied to the network.) This area is, in general, a NP-hard problem. Because reconfiguration is regularly triggered by continually changing traffic patterns, an optimized solution may not be stable. It may be more practical to develop heuristics that place more emphasis on factors such as fast convergence, and less on ongoing traffic, rather than on optimality.
- *Topology migration*—Topology migration consists of algorithms to coordinate the network migration from an old topology to a new one. Because WDM reconfiguration deals with large-capacity channels, changing allocation of channel resources in this coarse granularity significantly affects a large number of end-user flows. Traffic flows have to adapt to the light-path changes at and after each migration step. These effects can potentially spread over the routing pattern of the network, in turn possibly affecting more user flows.

Traditionally, all provisioning and engineering in optical networks has required manual planning and configuration, resulting in setup times of days or even weeks and a marked reluctance among network managers to de-provision resources in case doing so would affect other services. In the last few years, during which control protocols have been deployed to dynamically provide traffic engineering and provisioning or management assistance in optical networks, the control protocols have been proprietary and have greatly suffered from interoperability problems. Consequently, a new standardized control plane framework, supporting evolutionary traffic-engineering features, is needed for automatically switched optical transport networks to foster the expedited development and deployment of a new class of versatile optical switches that specifically address the optical transport needs of the Internet.

The important remaining challenge to be addressed in developing a dynamically reconfigurable optical network is that of controlling the optical resources, especially under distributed control where the network elements exchange information among themselves in a standardized multivendor environment. Performance and reliability requirements make this challenge of paramount importance to photonic networks. Beyond eliminating proprietary “islands of deployment,” this common control plane enables independent innovation curves within each product class, and faster service deployment with end-to-end provisioning using a single set of semantics.

The GMPLS Paradigm

GMPLS, the emerging paradigm for the design of control planes for OXCs, aims to address and solve all the challenges mentioned previously, trying to automatically and dynamically configure any kind of network element. It was proposed shortly after *Multiprotocol Label Switching* (MPLS) to extend its packet control plane to encompass time division (for example, for SONET/SDH), wavelength (for optical lambdas) and spatial switching (for example, for incoming port or fiber to outgoing port or fiber). Nongeneralized MPLS overlays a packet-switched IP network to facilitate traffic engineering and allow resources to be reserved and routes predetermined. It provides virtual links or tunnels through the network to connect nodes that lie at the edge of the network. For packets injected into the ingress of an established tunnel, normal IP routing procedures are suspended; instead the packets are label-switched so that they automatically follow the tunnel to its egress.

With the success of MPLS in packet-switched IP networks, optical network providers have accelerated a process to generalize the applicability of MPLS to cover all-optical networks as well. The premise of GMPLS is that the idea of a label can be generalized to be anything that is sufficient to identify a traffic flow. For example, in an optical fiber whose bandwidth is divided into wavelengths, the whole of one wavelength could be allocated to a requested flow. The *Label Switch Routers* (LSRs) at either end of the fiber simply have to agree on which frequency to use. From a control plane perspective, an LSR bases its functions on a table that maintains relations between incoming label or port and outgoing label or port. It should be noted that in the case of the OXC, the table that maintains the relations is not a software entity but it is implemented in a more straightforward way, for example, by appropriately configuring the micro-mirrors of the optical switching fabric.

There are several constraints in reusing the GMPLS control plane. These constraints arise from the fact that LSRs and OXCs use different data technologies. More specifically, LSRs manipulate packets that bear an explicit label, and OXCs manipulate wavelengths that bear the label implicitly; that is, the label value is implicit in the fact that the data is being transported within the agreed frequency band.

Furthermore, because the analogy of a label in the OXC is a wavelength or an optical channel, there are no equivalent concepts of label merging nor label push and pop operations in the optical domain, and label swapping can be realized through wavelength conversion. The transparency and multiprotocol properties of such a control plane approach would allow an OXC to route optical channel trails carrying various types of digital payloads (including IP, ATM, SDH, etc.) coherently and uniformly.

GMPLS Control Plane Functions and Services

GMPLS focuses mainly on the control plane services that perform connection management for the data plane (the actual forwarding logic) for both packet-switched interfaces and non-packet-switched interfaces. The GMPLS control plane essentially facilitates four basic functions:

- *Routing control*—Provides the routing capability, traffic engineering, and topology discovery
- *Resource discovery*—A mechanism to keep track of the system resource availability such as bandwidth, multiplexing capability, and ports
- *Connection management*—Provides end-to-end service provisioning for different services, including connection creation, modification, status query, and deletion
- *Connection restoration*—Implements an additional level of protection to the networks by establishing for each connection one or more presignaled backup paths and enabling very fast switching in case of failure between them.

The fundamental service offered by the GMPLS control plane is dynamic end-to-end connection provisioning. The operators need only to specify the connection parameters and send them to the ingress node. The network control plane then determines the optical paths across the network according to the parameters that the user provides and signals the corresponding nodes to establish the connection. The whole procedure can be done within seconds instead of hours. The other important service is bandwidth on demand, which extends the ease of provisioning even further by allowing the client devices that connect to the optical network to request the connection setup in real time as needed. In order to establish a connection that will be used to transfer data between a source–destination node pair, a light path needs to be established by allocating, in presence of the so-called *continuity constraint*, the same wavelength throughout the route of the transmitted data or selecting the proper wavelength conversion-capable nodes across the path. In fact, if the wavelength continuity constraint is not fully enforced, some wavelength conversion-capable nodes can be placed in the network to reduce the overall blocking probability in case of wavelength resource exhaustion on some nodes. Light paths can span more than one fiber link and remain entirely optical from end to end.

However, according to the mandatory clash constraint, two light paths traversing the same fiber link cannot share the same wavelength on that link. That is, each wavelength on a given fiber is not a sharable resource between light paths.

In general, if there are multiple feasible wavelengths (lambdas) between a source node and a destination node, then a Wavelength Assignment algorithm is required to select a wavelength for a given light path. The wavelength selection can be performed either after an optical route has been determined (in the so-called *decoupled approach*), or in parallel with finding a route. In the latter case, we refer to the coupled approach, in which the entire job is accomplished by a single *Routing and Wavelength Assignment* (RWA) algorithm. When light paths are established and taken down dynamically, routing and wavelength assignment decisions must be made as connection requests arrive to the network. It is possible that, for a given connection request, there may be insufficient network resources to set up a light path, in which case the connection request is blocked. The connection may also be blocked if there is no common wavelength available on all the links along the chosen route. Thus, the objective in the dynamic situation is to choose a route and a wavelength that maximizes the probability of setting up a given connection, while at the same time attempting to minimize the blocking for future connections.

In addition, because the quality of an optical signal degrades as it travels through several optical components and fiber segments, the deployment of “long-distance” light paths may require signal regeneration at strategic locations in a nationwide or global WDM network. As a result, the algorithms performing routing and wavelength assignment, virtual-topology embedding, wavelength conversion, etc. must also be mindful of the locations of the sparse signal regenerators in the network. Such regenerators, which are placed at select locations in the network, “clean up” the optical WDM signal either entirely in the optical domain or through an optoelectronic conversion followed by an electro-optic conversion. Thus the signal from the source travels through the network as far as possible before its quality drops below a certain threshold, thereby requiring it to be regenerated at an intermediate node. The same signal could be regenerated several times in the network before it reaches the destination.

Furthermore, in current multilayer transport networks the bandwidth demanded by traffic typically is orders of magnitude lower than the capacity of lambda links, and the number of available wavelengths per fiber is limited and costly. Hence, it is not worth assigning exclusive end-to-end light paths to these demands, so a better sub-lambda granularity is required. Thus, to increase the throughput of a network with a limited number of lambdas per fiber, *traffic grooming* is required in certain nodes, typically those on the network edge.

The GMPLS control plane ensures traffic-grooming capability on edge nodes by operating on a two-layer model; that is, an underlying pure optical wavelength routed network and an “optoelectronic” time-division multiplexed layer built over it. In the wavelength routed layer, operating exclusively at lambda granularity, when a transparent light path connects two physically adjacent or distant nodes, these nodes will seem adjacent for the upper layer. The upper layer can perform multiplexing of different traffic streams into a single wavelength-based light path through simultaneous time and space switching. Similarly it can demultiplex different traffic streams of a single lambda path. It can also perform remultiplexing: some of the demands demultiplexed can be again multiplexed into some other wavelength paths and handled together along it. This is due to the “generalized” and hence multilayer nature of the GMPLS control plane.

The electronic layer is clearly required for multiplexing packets coming from different ports. This upper electronic layer can be a classical or “next-generation” technology, such as IP/MPLS, but it can also be based on any other networking technology (that is SDH/SONET, ATM, Ethernet, etc.). However, the technology of the upper layer must be unique for all traffic streams that have to be demultiplexed and then multiplexed again, because the network cannot directly multiplex, for example, ATM cells with Ethernet frames.

Another service that gives greatest flexibility to users in handling their own virtual network topologies on the transport core is the *Optical Virtual Private Network* (OVPN), which allows users to have full network resource control of a defined partition of the carrier optical network. Although users have full network resource control of that portion of the network, the OVPN is just a logical network partition and the end users still do not have access and visibility to the carrier’s networks. This service can save the carrier’s operation resources by allowing end users to perform circuit provisioning and setup procedures.

GMPLS Interfaces

GMPLS encompasses control plane signaling for multiple interface types. The diversity of controlling not only switched packets and cells but also TDM network traffic and optical network components makes GMPLS flexible enough to position itself in the direct migration path from electronic to all-optical network switching. The five main interface types supported by GMPLS follow:

- *Packet Switching Capable* (PSC)—These interfaces recognize packet boundaries and can forward packets based on the IP header or a standard MPLS “shim” header.
- *Layer 2 Switch-Capable* (L2SC)—These interfaces recognize frame and cell headers and can forward data based on the content of the frame or cell header (for example, an ATM LSR that forwards data based on its *Virtual Path Identifier/Virtual Circuit Identifier* (VPI/VCI) value, or Ethernet bridges that forward the data based on the MAC header).

- *Time-Division Multiplexing-Capable* (TDMC)—These interfaces forward the data based on the time slot in a repeating cycle (for example, SDH cross-connect or ADM, interfaces implementing the Digital Wrapper G.709, and *Plesichronous Digital Hierarchy* [PDH] interfaces).
- *Lambda Switch-Capable* (LSC)—These interfaces are for wavelength-based MPLS control of optical devices and wavelength switching devices, such as *optical ADMs* (OADMs) and OXCs, operating at the granularity of the single wavelength or group of wavelengths (waveband). These interfaces forward the optical signal from an incoming optical wavelength to an outgoing optical wavelength. Traffic is forwarded based upon wavelength or waveband.
- *Fiber-Switch-Capable* (FSC)—These interfaces forward the signal from one or more incoming fibers to one or more outgoing fibers for spatial control of interface selection, automated patch panels, and physical fiber switching systems. Traffic is forwarded based on port, fiber, or interface.

These supported interfaces are hierarchal in structure and controlled simultaneously by GMPLS.

Generalized Label

GMPLS defines several new forms of label—the *generalized label* objects. These objects include the generalized label request, the generalized label, the explicit label control, and the protection flag. The generalized label can be used to represent timeslots, wavelengths, wavebands, or space-division multiplexed positions.

With plain MPLS labels embedded in the cell or packet structure for in-band control plane signaling, with the different kinds of interfaces supported by GMPLS it is impossible to embed label-specific information, in terms of fiber port or wavelength switching, into the traffic packet structure. Consequentially, new “virtual” labels have been added to the MPLS label structure. These virtual labels comprise specific indicators that represent wavelengths, fiber bundles, or fiber ports and are distributed to GMPLS nodes through out-of-band GMPLS signaling. GMPLS out-of-band signaling causes a control-channel separation problem.

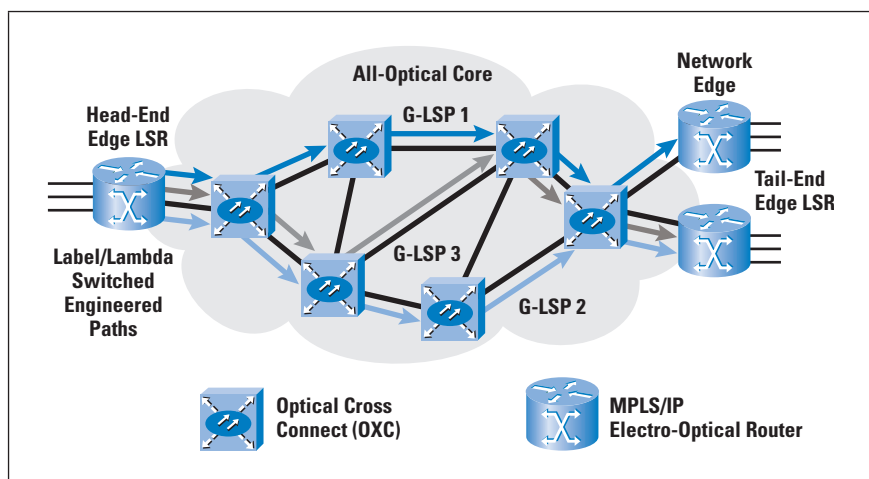
With MPLS, the control information is found in the label, which is directly attached to the data payload. However, when you send the control information out of band, the label is separated from the data that it is attempting to control. GMPLS provides a means for identifying explicit data channels. Having the ability to identify data channels allows the control message to be associated with a particular data flow, whether it is a wavelength, fiber, or fiber bundle.

Generalized Label-Switched Paths

The handling of *label-switched paths* (LSPs) under GMPLS differs from that of MPLS. MPLS does not provide for bidirectional LSPs. Each direction LSP has to be established in turn. Under GMPLS, the LSP can be established bidirectionally. The traffic-engineering requirements for the bidirectional LSP are the same in both directions, and it is established for both directions through only one signaling message, allowing for reductions in latency-related setup time. In the optical environment, OXC translates label assignments into corresponding wavelength assignments and sets up *generalized LSPs* (G-LSPs) using their local control interfaces to the other switching devices. Subsequent to G-LSP setup, no explicit label or lambda lookup or processing operations are performed by the OXC nodes.

GMPLS supports traffic engineering by allowing the node at the network ingress to specify the route that a G-LSP will take by using explicit light-path routing. An explicit route is specified by the ingress as a sequence of hops and wavelengths that must be used to reach the egress, which is different from the hop-by-hop routing that is usually associated with PSC networks.

Figure 4: G-LSPs Ensuring Traffic Engineering



GMPLS also maintains the capability already available with MPLS to nest G-LSPs. Nested G-LSPs make possible the building of a forwarding hierarchy. At the top of this hierarchy are nodes that have FSC interfaces, followed by nodes that have LSC interfaces, followed by nodes that have TDMC interfaces, and followed by nodes with PSC interfaces. Nesting of G-LSPs between interface types increases flexibility in service definition and makes it possible for service providers operating a GMPLS network to deliver both bundled and unbundled services.

Because the deployment of DWDM equipment makes feasible the creation a large number of individual connections between two adjacent nodes, another very useful feature of bundling is the ability to simultaneously handle multiple adjacent links. Link bundling treats the traffic of these links as a single link.

In order for the adjacent links to be bundled, they must be on the same GMPLS segment, they must be of the same type, and they must have the same traffic-engineering requirements. These requirements reduce the amount of link advertisements that need to be maintained throughout the network, thereby increasing the control plane scalability. Just as in MPLS label stacking, GMPLS labels only contain information about a single level of hierarchy. The difference for GMPLS is that this hierarchy can be fiber-, wavelength-, timeslot-, packet- or cell-based.

For instance, if a connection is desired from one PSC interface to another PSC interface, and the traffic traverses physically separate fibers, a unique LSP has to be established for each level in turn. First, the FSC LSP, then the LSC LSP, then the TDMC LSP, and finally the PSC LSP have to be established through GMPLS signaling.

Signaling and Routing Protocols

In order to set up a light path, a signaling protocol is also required to exchange control information among nodes, to distribute labels, and to reserve resources along the path. In our case, the signaling protocol is closely integrated with the routing and wavelength assignment protocols. Suitable GMPLS signaling protocols for the GMPLS control plane include *Resource Reservation Protocol* (RSVP) and *Constraint-Based Label Distribution Protocol* (CR-LDP). Any of the objects that are defined within the GMPLS specification can be carried within the message of either of these signaling protocols that are responsible for all the connection management actions such as setup, modify, or remove the G-LSPs. Clearly, support for provisioning and restoration of end-to-end optical trails within a photonic network consisting of heterogeneous networking elements imposes new requirements for these signaling protocols. Specifically, optical trails require small setup latency (especially for restoration purposes), support for bidirectional trails, rapid failure detection and notification, and fast intelligent trail restoration.

Both RSVP and CR-LDP can be used to reserve a single wavelength for a light path if the wavelength is known in advance. These protocols can also be modified to incorporate wavelength selection functions into the reservation process^[7]. In RSVP, signaling takes place between the source and destination nodes. The signaling messages may contain information such as QoS requirements for the carried traffic and label requests for assigning labels at intermediate nodes that reserve the appropriate resources for the path. CR-LDP uses TCP sessions between nodes in order to provide a hop-by-hop reliable distribution of control messages, indicating the route and the required traffic parameters for the route. Each intermediate node reserves the required resources, allocates a label, and sets up its forwarding table before backward signaling to the previous node.

To correctly perform resource reservation, allocation, and topology discovery on the available optical link resources, each node needs to maintain a representation of the state of each link in the network. The link state includes the total number of active channels, the number of allocated channels, and the number of channels reserved for light-path restoration. Additional parameters can be associated with allocated channels; for example, some light paths can be preemptable or have associated hold priorities. When the local inventory is constructed, the node engages in a routing protocol to distribute and maintain the topology and resource information. Standard IP routing protocols, such as *Open Shortest Path First* (OSPF) or *Intermediate System-to-Intermediate System* (IS-IS) with GMPLS Traffic Engineering extensions, can be used to reliably propagate the information.

The extensions to OSPF and IS-IS add additional information about links and nodes into the link-state database. Such information includes the type of LSPs that can be established across a given link (for example, packet forwarding, SONET/SDH trails, wavelengths, or fibers), as well as the current unused bandwidth, the maximum size of G-LSP that can be established, and the administrative groups supported. This information allows the node computing the explicit route for an LSP to do so more intelligently. Furthermore, any switching node cooperating in the GMPLS control plane will maintain a per-interface or per-fiber *Wavelength Forwarding Information Base* (WFIB) because lambdas and channels (labels) are specific to a particular interface or fiber, and the same lambda or channel (label) could be used concurrently on multiple interfaces or fibers.

Link Management Protocol

GMPLS also uses the *Link Management Protocol* (LMP) to communicate proper cross-connect information between the network elements. LMP runs between adjacent systems for link provisioning and fault isolation. It can be used for any type of network element, particularly in natively photonic switches. LMP automatically generates and maintains associations between links and labels for use in label swapping^[6]. Automating the labeling process simplifies management and avoids the errors associated with manual label assignment. LMP provides control-channel management, link-connectivity verification, link-property correlation, and fault isolation. Control-channel management establishes and maintains connectivity between adjacent nodes using a keepalive protocol. Link verification verifies the physical connectivity between nodes, thereby detecting loss of connections and misrouting of cable connections. Fault isolation pinpoints failures in both electronic and optical links without regard to the data format traversing the link.

In order for these link bundles to be handled accordingly, GMPLS needed a method to manage the links between adjacent nodes. LMP was developed to address several link-specific problems that surfaced when generalizing the MPLS protocol across different interface types. The main responsibilities of the LMP follow:

- *Control-Channel Management*—Establishment of a control channel is critical to GMPLS signaling. The maintenance of the control channel between adjacent nodes must be able to exchange information related to LSP establishment.
- *Link-Property Correlation*—When link bundling occurs, GMPLS requires a way to verify that all traffic-engineering requirements are similar between links of adjacent nodes. Link-property correlation performs the verification and the aggregation of such links.
- *Link-Connectivity Verification*—This feature is used by GMPLS to verify the connectivity between data links when the control channel is separate from each data link.
- *Fault Management*—Fault management helps the network isolate faults down to the individual link.

Although LMP assumes the messages are IP encoded, it does not dictate the actual transport mechanism used for the control channel. However, the control channel must terminate on the same two nodes that the bearer channels span. Therefore, this protocol can be implemented on any OXC, regardless of the internal switching fabric. A requirement for LMP is that each link has an associated bidirectional control channel and that free bearer channels must be opaque (that is, able to be terminated); however, when a bearer channel is allocated, it may become transparent. Note that this requirement is trivial for optical cross-connects with electronic switching planes, but is an added restriction for photonic switches.

Conclusion

Innovations in the field of optical components will take advantage of the introduction of all-optical networking in all areas of information transport and will offer system designers the opportunity to create new solutions that will allow smooth evolution of all telecommunication networks. A new class of versatile IP-addressable optical switching devices is emerging, operating according to a common GMPLS-based control plane to support full-featured traffic engineering in modern optical transparent infrastructures.

The main advantage of this approach is that it is based on already existing and widely deployed protocols while simplifying network management and engineering tasks that can be performed in a unified way in both the data and the optical domains. Furthermore, it offers a function framework that can accommodate future expectations concerning the way networks will work and the way services will be provided to clients. Thus we envision a horizontal network, harmonized by a common GMPLS-based control plane, where all network elements work as peers to dynamically establish optical paths through the network.

This new photonic internetwork will make it possible to provision high bandwidth in tenths of seconds, and enable new revenue-generating services and dramatic cost savings for service providers.

In the same way that digital communication technologies changed the twentieth century into the “electronic century,” the optical technologies discussed in this article will make the next century “the photonic century.” All winning strategies must rely on such GMPLS-based photonic infrastructures—an environment in which innovations work at the speed of light.

For Further Reading

- [1] B. E. A. Saleh and M. C. Teich, *Fundamentals of Photonics*, John Wiley & Sons Inc., 1991.
- [2] P. Raghavan and E. Upfal, “Efficient Routing in All-Optical Networks,” *Proceedings of ACM STOC’94*, 1994.
- [3] B. Mukherjee, *Optical Communication Networks*, McGraw-Hill, 1997.
- [4] A. Mokhtar and M. Azizoglu, “Adaptive Wavelength Routing in All-Optical Networks,” *IEEE/ACM Transactions on Networking*, vol. 6, pp. 197–206, April 1998.
- [5] E. Karasan and S. Banerjee, “Performance of WDM Transport Networks,” *IEEE Journal on Selected Areas in Communications*, vol. 16, pp. 1081–1096, September 1998.
- [6] A. Banerjee, J. Drake, J. Lang, B. Turner, K. Kompella, and Y. Rekhter, “Generalized Multiprotocol Label Switching: An Overview of Routing and Management Enhancements,” *IEEE Communications Magazine*, January 2001.
- [7] A. Banerjee, J. Drake, J. Lang, B. Turner, D. O. Awduche, L. Berger, K. Kompella, and Y. Rekhter, “Generalized Multiprotocol Label Switching: An Overview of Signalling Enhancements and Recovery Techniques,” *IEEE Communications Magazine*, July 2001.

FRANCESCO PALMIERI holds two Computer Science degrees from Salerno University, Italy. Since 1997, he has led the network management and operation centre of the Federico II University, in Napoli, Italy. He has been closely involved with the development of the Internet in Italy in the last few years, particularly within the academic and research sector, and is actually a member of the Technical Scientific Committee and of the Computer Emergency Response Team of the Italian NREN GARR. He worked for several international companies on a variety of networking-related projects concerned with nationwide communication systems, network management, transport protocols, and IP networking. He is an active researcher in the fields of high-performance, evolutionary networking, and network security. He regularly publishes in leading technical journals and conferences and gives invited talks and keynote speeches. E-Mail: Francesco.Palmieri@unina.it

The Changing Foundation of the Internet: Confronting IPv4 Address Exhaustion

by Geoff Huston, APNIC

Throughout its relatively brief history, the Internet has continually challenged our preconceptions about networking and communications architectures. For example, the concepts that the network itself has no role in management of its own resources, and that resource allocation is the result of interaction between competing end-to-end data flows, were certainly novel innovations, and for many they have been very confrontational. The approach of designing a network that is unaware of services and service provisioning and is not attuned to any particular service whatsoever—leaving the role of service support to end-to-end overlays—was again a radical concept in network design. The Internet has never represented the conservative option for this industry, and has managed to define a path that continues to present significant challenges.

From such a perspective it should not be surprising that the next phase of the Internet story—that of the transition of the underlying version of the IP protocol from IPv4 to IPv6—refuses to follow the intended script. Where we are now, in late 2008, with IPv4 unallocated address pool exhaustion looming within the next 18 to 36 months, and IPv6 still largely not deployed in the public Internet, is a situation that was entirely unanticipated and, even in hindsight, entirely surprising.

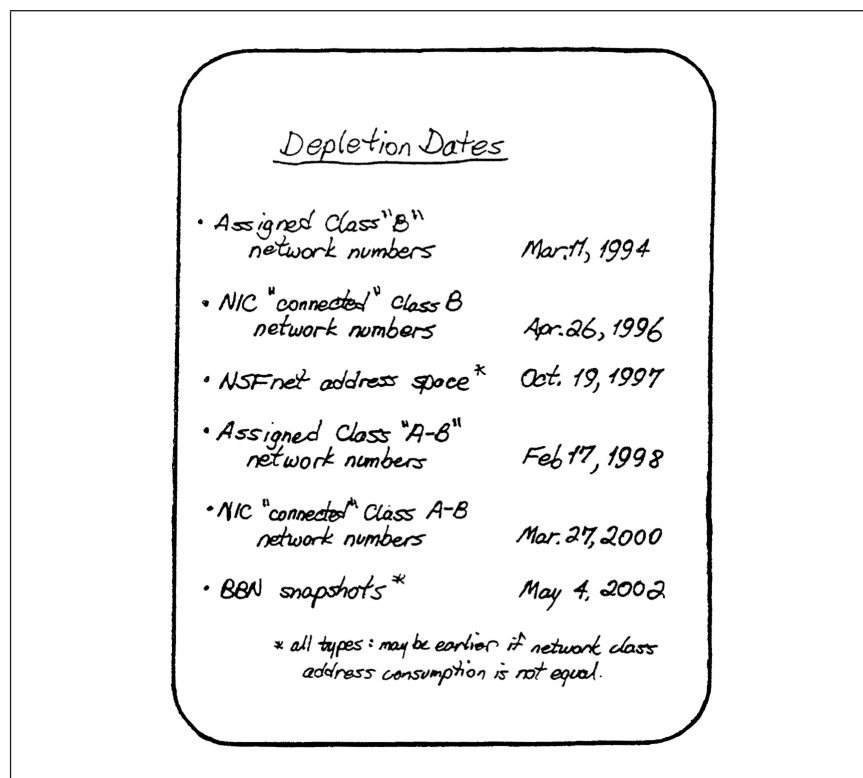
The topic examined here is *why* this situation has arisen, and in examining this question we analyze the options available to the Internet to resolve the problem of IPv4 address exhaustion. We examine the timing of the IPv4 address exhaustion and the nature of the intended transition to IPv6. We consider the shortfalls in the implementation of this transition, and identify their underlying causes. And finally, we consider the options available at this stage and identify some likely consequences of such options.

When?

This question was first asked on the TCP/IP list in November 1988, and the responses included foreshadowing a new version of IP with longer addresses and undertaking an exercise to reclaim unused addresses^[1]. The exercise of measuring the rate of consumption of IPv4 addresses has been undertaken many times in the past two decades, with estimates of exhaustion ranging from the late 1990s to beyond 2030. One of the earliest exercises in predicting IPv4 address exhaustion was undertaken by Frank Solensky and presented at IETF 18 in August 1990. His findings are reproduced in Figure 1.

At that time the concern was primarily the rate of consumption of Class B network addresses (or of /16 prefixes from the address block 128.0.0.0/2, to use current terminology). Only 16,384 such Class B network addresses were within the class-based IPv4 address plan, and the rate of consumption was such that the Class B networks would be fully consumed within 4 years, or by 1994. The prediction was strongly influenced by a significant number of international research networks connecting to the Internet in the late 1980s, with the rapid influx of new connections to the Internet creating a surge in demand for Class B networks.

Figure 1: Report on IPv4 Address Depletion^[2]



Successive predictions were made in the context of the *Internet Engineering Task Force* (IETF) in the *Address Lifetime Expectancy* (ALE) Working Group, where the predictive model was refined from an exponential growth model to a logistical saturation function, attempting to predict the level at which all address demands would be met.

The predictive technique described here is broadly similar, using a statistical fit of historical data concerning address consumption into a mathematical model, and then using this model to predict future address consumption rates and thereby predict the exhaustion date of the address pool.

The predictive technique models the IP address distribution framework. Within this framework the pool of unallocated /8 address blocks is distributed by the *Internet Assigned Numbers Authority* (IANA) to the five *Regional Internet Registries* (RIRs). (A “/8 address block” refers to a block of addresses where the first 8 bits of the address values are constant. In IPv4 a /8 address block corresponds to 16,777,216 individual addresses.) Within the framework of the prevailing address distribution policies, each RIR can request a further address allocation from IANA when the remaining RIR-managed unallocated address pool falls below a level required to meet the next 9 months of allocation activity. The amount allocated is the number of /8 address blocks required to augment the RIR’s local address pool to meet the anticipated needs of the regional registry for the next 18 months. However, in practice, the RIRs currently request a maximum of 2 /8 address blocks in any single transaction, and do so when the RIR-managed address pool falls below a threshold of the equivalent of 2 /8 address blocks.

As of August 2008 some 39 /8 address blocks are left in IANA’s unallocated address pool. A predictive exercise has been undertaken using a statistical modeling of historical address consumption rates, using data gathered from the RIRs’ records of address allocations and the time series of the total span of address space announced in the Internet interdomain default-free routing table as basic inputs to the model. The predictive technique is based on a least-squares best fit of a linear function applied to the first-order differential of a smoothed copy of the address consumption data series, as applied to the most recent 1,000 days’ data.

The linear function, which is a best fit to the first-order differential of the data series, is integrated to provide a quadratic time-series function to match the original data series. The projection model is further modified by analyzing the day-of-year variations from the smoothed data model, averaged across the past 3 years, and applying this daily variation to the projection data to account for the level of seasonal variations in the total address consumption rate that has been observed in the historical data. The anticipated rate of consumption of addresses from this central pool of unallocated IPv4 addresses is expected to be about 15 /8s in 2009, and slightly more in 2010.

RIR behaviors are modeled using the current RIR operational practices and associated address policies, which are used to predict the times when each RIR will be allocated a further 2 /8s from IANA. This RIR consumption model, in turn, allows the IANA address pool to be modeled.

This anticipated rate of increasing address consumption will see the remaining unallocated addresses that are held by IANA reach the point of exhaustion in February 2011. The most active RIRs are anticipated to exhaust their locally managed unallocated address pools in the months following the time of IANA exhaustion.

The assumptions behind this form of prediction follow:

- The current policy framework relating to the distribution of addresses will continue to apply without any further alteration through to complete exhaustion of the unallocated address pool.
- The demand curves will remain consistent, meaning that there will be no forms of disruption to demand, such as a panic rush on the remaining addresses or some introduced externality that affects total address demand.
- The level of return of addresses to the unallocated address pool will not vary significantly from existing levels of address return.

Although the statistical model is based on a complete data set of address allocations and a detailed hourly snapshot of the address span advertised in the Internet routing table, a considerable level of uncertainty is still associated with this prediction.

First, the behavior of the *Internet Service Provider* (ISP) industry and the other entities that are the direct recipients of RIR address allocations and assignments are not ignorant of the impending exhaustion condition, and there is some level of expectation of some form of last-minute rush or panic on the part of such address applicants when exhaustion of this address pool is imminent. The predictive model described here does not include such a last-minute acceleration of demand.

The second factor is the skewed distribution of addresses in this model. From 1 January 2007 until 20 July 2008, 10,402 allocation or assignments transactions were recorded in the RIRs' daily statistics files. These transactions accounted for a total of 324,022,704 individual IPv4 addresses, or the equivalent of 19.3 /8s. Precisely one-half of this address space was allocated or assigned in just 107 such transactions.

In other words, some 1 percent of the recipients of address space in the past 18 months have received some 50 percent of all the allocated address space. The reason why this distribution is relevant here is that this predictive exercise assumes that although individual actions are hard to predict with any certainty, the aggregate outcome of many individuals' actions assumes a much greater level of predictability.

This observation about aggregate behavior does not apply in this situation, however, and the predictive exercise is very sensitive to the individual actions of a very small number of recipients of address space because of this skewed distribution of allocations. Any change in the motivations of these larger-sized actors that results in an acceleration of demand for IPv4 will significantly affect the predictions of the longevity of the remaining unallocated IPv4 address pool.

The third factor is that this model assumes that the policy framework remains unaltered, and that all unallocated addresses are allocated or assigned under the current policy framework, rather than under a policy regime that is substantially different from today's framework. The related assumption here is that the cost of obtaining and holding addresses remains unchanged, and that the perceptions of future scarcity of addresses do not affect the policy framework of address distribution of the remaining unallocated IPv4 addresses.

Given this potential for variation within this set of assumptions, a more accurate summary of the current expectations of address consumption would be that the exhaustion of the IANA unallocated IPv4 address pool will occur sometime between July 2009 and July 2011, and that the first RIR will exhaust all its usable address space within 3 to 12 months from that date, or between October 2009 and July 2012.^[3]

What Next?

Apart from the exact date of exhaustion that is predicted by this modeling exercise, none of the information relating to exhaustion of the unallocated IPv4 address pool should be viewed as particularly novel information. The IETF *Routing and Addressing* (ROAD) study of 1991 recognized that the IPv4 address space was always going to be completely consumed at some point in the future of the Internet^[4].

Such predictions of the potential for exhaustion of the IPv4 address space were the primary motivation for the adoption of *Classless Inter-Domain Routing* (CIDR) in the *Border Gateway Protocol* (BGP), and the corresponding revision of the address allocation policies to craft a more exact match between planned network size and the allocated address block. These predictions also motivated the protracted design exercise of what was to become the IPv6 protocol across the 1990s within the IETF. The prospect of address scarcity engendered a conservative attitude to address management that, in turn, was a contributory factor in accelerating the widespread use of *Network Address Translation* (NAT)^[5] in the Internet during the past decade. By any reasonable metric this industry has had ample time to study this problem, ample time to devise various strategies, and ample time to make plans and execute them.

And this reality has been true for the adoption of classless address allocations, the adoption of CIDR in BGP, and the extremely widespread use of NAT. But all of these measures were short-term, whereas the longer-term measure, that of the transition to IPv6, was what was intended to come after IPv4. But IPv6 has not been the subject of widespread adoption so far, while the time of anticipated exhaustion of IPv4 has been drawing closer. Given almost two decades of advance warning of IPv4 address exhaustion, and a decade since the first stable implementations of IPv6 were released, we could reasonably expect that this industry—and each actor within this industry—is aware of the problem and the need for a stable and scalable long-term solution as represented by IPv6. We could reasonably anticipate that the industry has already planned the actions it will take with respect to IPv6 transition, and is aware of the triggers that will invoke such actions, and approximately when they will occur.

However, such an expectation appears to be ill-founded when considering the broad extent of the actors in this industry, and there is little in the way of a common commitment as to what will happen after IPv4 address exhaustion, nor even any coherent view of plans that industry actors are making in this area.

This lack of planning makes the exercise of predicting the actions within this industry following address exhaustion somewhat challenging, so instead of immediately describing future scenarios, it may be useful to first describe the original plan for the response of the Internet to IPv4 address exhaustion.

What Was Intended?

The original plan, devised in the early 1990s by the IETF to address the IPv4 address shortfall, was the adoption of CIDR as a short-term measure to slow down the consumption of IPv4 addresses by reducing the inefficiency of the address plan, and the longer-term plan of the specification of a new version of the Internet Protocol that would allow for adoption well before the IPv4 address pool was exhausted.

The industry also adopted the use of NAT as an additional measure to increase the efficiency of address use, although the IETF did not strongly support this protocol. For many years the IETF did not undertake the standardization of NAT behaviors, presumably because NAT was not consistent with the IETF's advocacy of end-to-end coherence of the Internet at the IP level of the protocol stack.

Over the 1990s the IETF undertook the exercise of the specification of a successor IP protocol to Version 4, and the IETF's view of the longer-term response was refined to be advocacy of the adoption of the IPv6 protocol and the use of this protocol as the replacement for IPv4 across all parts of the network.

In terms of what has happened in the past 15 years, the adoption of CIDR was extremely effective, and most parts of the network were transitioned to use CIDR within 2 years, with the transition declared to be complete by the IETF in June 1996. And, as noted already, NAT has been adopted across many, if not most, parts of the network. The most common point of deployment of NAT has not been at an internal point of demarcation between provider networks, but at the administrative boundary between the local customer network and the ISP, so that the common configuration of *Customer Premises Equipment* (CPE) includes NAT functions. Customers effectively own and operate NAT devices as a commonplace aspect of today's deployed Internet.

CIDR and NAT have been around for more than a decade now, and the address consumption rates have been held at very conservative levels in that period, particularly so when considering that the bulk of the population of the Internet was added well after the advent of CIDR and NAT.

The longer-term measure—the transition to IPv6—has not proved to be as effective in terms of adoption in the Internet.

There was never going to be a “flag-day” transition where, in a single day, simultaneously across all parts of every network the IP protocol changed to using IPv6 instead of IPv4. The Internet is too decentralized, too large, too disparate, and too critical for such actions to be orchestrated, let alone completed with any chance of success. A flag day, or any such form of coordinated switchover, was never a realistic option for the Internet.

If there was no possibility of a single, coordinated switchover to IPv6, the problem is that there was never going to be an effective piecemeal switchover either. In other words, there was never going to be a switchover where host by host, and network by network, IPv6 is substituted for IPv4 on a piecemeal and essentially uncoordinated basis. The problem here is that IPv6 is not “backward-compatible” with IPv4. When a host uses IPv6 exclusively, then that host has no direct connectivity to any part of the IPv4 network. If an IPv6-only host is connected to an IPv4-only network, then the host is effectively isolated. This situation does not bode well for a piecemeal switchover, where individual components of the network are switched over from IPv4 to IPv6 on a piecemeal basis. Each host that switches over to IPv6 essentially disconnects itself from the IPv4 Internet at that point.

Given this inability to support backward compatibility, what was planned for the transition to IPv6 was a “dual-stack” transition. Rather than switching over from IPv4 to IPv6 in one operation on both hosts and networks, a two-step process has been proposed: first switching from IPv4 only to a “dual-stack” mode of operation that supports both IPv4 and IPv6 simultaneously, and second—and at a much later date—switching from dual-stack IPv4 and IPv6 to IPv6 only.

During the transition more and more hosts are configured with dual stack. The idea is that dual-stack hosts prefer to use IPv6 to communicate with other dual-stack hosts, and revert to use IPv4 only when an IPv6-based end-to-end conversation is not possible. As more and more of the Internet converts to dual stack, it is anticipated that use of IPv4 will decline, until support for IPv4 is no longer necessary. In this dual-stack transition scenario, no single flag day is required and the dual-stack deployment can be undertaken in a piecemeal fashion. There is no requirement to coordinate hosts with networks, and as dual-stack capability is supported in networks the attached dual-stack hosts can use IPv6. This scenario still makes some optimistic assumptions, particularly relating to the achievement of universal deployment of dual stack, at which point no IPv4 functions are used, and support for IPv4 can be terminated. Knowing when this point is reached is unclear, of course, but in principle there is no particular timetable for the duration of the dual-stack phase of operation.

There are always variations, and in this case it is not necessarily that each host must operate in dual-stack mode for such a transition. A variant of the NAT approach can perform a rudimentary form of protocol translation, where a *Protocol-Translating NAT* (or NAT-PT^[6]) essentially transforms an incoming IPv4 packet to an outgoing IPv6 packet, and conversely, using algorithmic binding patterns to map between IPv4 and IPv6 addresses. Although this process relieves the IPv6-only host of some additional complexity of operation at the expense of some added complexity in *Domain Name System* (DNS) transformations and service fragility, the essential property still remains that in order to speak to an IPv4-only remote host, the combination of the local IPv6 host and the NAT-PT have to generate an equivalent IPv4 packet. In this case the complexity of the dual stack is now replaced by complexity in a shared state across the IPv6 host and the NAT-PT unit. Of course this solution does not necessarily operate correctly in the context of all potential application interactions, and concerns with the integrity of operation of NAT-PT devices are significant, a factor that motivated the IETF to deprecate the existing NAT-PT specification^[7]. On the other hand, the lack of any practical alternatives has led the IETF to subsequently reopen this work, and once again look at specifying the standard behavior of such devices^[8].

The detailed progress of a dual-stack transition is somewhat uncertain, because it involves the individual judgment of many actors as to when it may be appropriate to discontinue all support for IPv4 and rely solely on IPv6 for all connectivity requirements. However, one factor is constant in this envisaged transition scenario, and whether it is dual stack in hosts or dual stack through NAT-PT, or various combinations thereof, the requirement that there are sufficient IPv4 addresses to span the addressing needs of the entire Internet across the complete duration of the dual-stack transition process is consistent.

Under this dual-stack regime every new host on the Internet is envisaged to need access to both IPv6 and IPv4 addresses in order to converse with any other host using IPv6 or IPv4. Of course this approach works as long as there is a continuing supply of IPv4 addresses, implying that the envisioned timing of the transition was meant to have been completed by the time that IPv4 address exhaustion happens.

If this transition were to commence in earnest at the present time, in late 2008, and take an optimistic 5 years to complete, then at the current address consumption rate we will require a further 90 to 100 /8 address blocks to span this 5-year period. A more conservative estimate of a 10-year transition will require a further 200 to 250 /8 address blocks, or the entire IPv4 address space again, assuming that we will use IPv4 addresses in the future in precisely the same manner as we have used them in the past and with precisely the same level of usage efficiency as we have managed to date.

Clearly, waiting for the time of IPv4 unallocated address pool exhaustion to act as the signal to industry to commence the deployment of IPv6 in a dual-stack transition framework is a totally flawed implementation of the original dual-stack transition plan.

Either the entire process of dual-stack transition will need to be undertaken across a far faster time span than has been envisaged, or the manner of use of IPv4 addresses, and, in particular their usage efficiency in the context of dual-stack transition support, will need to differ markedly from the current manner of address use. Numerous forms of response may be required, posing some challenging questions because there is no agreed precise picture of what markedly different and significantly more efficient form of address use is required here. To paraphrase the situation, it is clear that we need to do “something” differently, and do so as a matter of some urgency, but we have no clear agreement on what that something is that we should be doing differently. This situation obviously is not an optimal one.

What was intended as a transition mechanism for IPv6 is still the only feasible approach that we are aware of, but the forthcoming exhaustion of the unallocated IPv4 address pool now calls for novel forms of use of IPv4 addresses within this transitional framework, and these novel forms may well entail the deployment of various forms of address translation technologies that we have not yet defined, let alone standardized. The transition may also call for scaling capabilities from the interdomain routing system that also head into unknown areas of technology and deployment feasibility.

Why?

At this point it may be useful to consider how and why this situation has arisen.

If the industry needed an abundant supply of IPv4 addresses to underpin the entire duration of the dual-stack transition to IPv6, then why didn't the industry follow the lead of the IETF and commence this transition while there was still an abundant supply of IPv4 addresses on hand? If network operators, service providers, equipment vendors, component suppliers, application developers, and every other part of the Internet supply chain were aware of the need to commence a transition to IPv6 well before effective exhaustion of the remaining pool of IPv4 addresses, then why didn't the industry make a move earlier? Why was the only clear signal for a change in Internet operation to commence a dual-stack transition to IPv6 one that has been activated too late to be useful for the industry to act on efficiently?

One possible reason may lie in a perception of the technical immaturity of IPv6 as compared to IPv4. It is certainly the case that many network operators in the Internet are highly risk-adverse and tend to operate their networks in a mainstream path of technologies rather than constantly using leading-edge advance releases of hardware and software solutions. Does IPv6 represent some form of unacceptable technical risk of failure that has prevented its adoption? This reasoning does not appear to be valid in terms of either observed testing or observation of perceptions about the technical capability of IPv6. The IPv6 protocol is functionally complete and internally consistent, and it can be used in almost all contexts where IPv4 is used today. IPv6 works as a platform for all forms of transport protocols, and is fully functional as an internetwork layer protocol that is functionally equivalent to IPv4. IPv6 NAT exists, *Dynamic Host Configuration Protocol Version 6* (DHCPv6) provides dynamic host configuration for IPv6 nodes, and the DNS can be completely equipped with IPv6 resource records and operate using IPv6 transport for queries and responses.

Perhaps the only notable difference between the two protocols is the ability to perform host scans in IPv6, where probe packets are sent to successive addresses. In IPv6 the address density is extremely low because the low-order 64-bit interface address of each host is more or less unique, and within a single network the various interface addresses are not clustered sequentially in the number space. The only known use of address probing to date has been in various forms of hostile attack tools, so the lack of such a capability in IPv6 is generally seen as a feature rather than an impediment. IPv6 deployment has been undertaken in a small scale for many years, and although the size of the deployed IPv6 base remains small, the level of experience gained with the technology functions has been significant. It is possible to draw the conclusion that IPv6 is technically capable and this capability has been broadly tested in almost every scenario except that of universal use across the Internet.

It also does not appear that the reason was a lack of information or awareness of IPv6. The efforts to promote IPv6 adoption have been under way in earnest for almost a decade now. All regions and many of the larger economies have instigated programs to promote the adoption of IPv6 and have provided information to local industry actors of the need to commence a dual-stack transition to IPv6 as soon as possible. In many cases these promotional programs have enjoyed broad support from both public and industry funding sources. The coverage of these promotional efforts has been widespread in industry press reports. Indeed, perhaps the only criticism of this effort is possibly too much promotion, with a possible result that the effectiveness of the message has been diluted through constant repetition.

A more likely area to examine in terms of possible reasons why industry has not engaged in dual-stack transition deployment is that of the business landscape of the Internet. The Internet can be viewed as a product of the wave of progressive deregulation in the telecommunications sector in the 1980s and early 1990s. New players in the deregulated industry searching for a competitive edge to unseat the dominant position of the traditional incumbents found the Internet as their competitive lever. The result was perhaps unexpected, because it was not one that replaced one vertically integrated operator with a collection of similarly structured operators whose primary means of competition was in terms of price efficiency across an otherwise undifferentiated service market, as we saw in the mobile telephony industry. In the case of the Internet, the result was not one that attempted to impose convergence on this industry, but one that stressed divergence at all levels, accompanied by branching role specialization at every level in the protocol stack and at every point in the supply chain process. In the framework of the Internet, consumers are exposed to all parts of the supply process, and do not rely on an integrator to package and supply a single, all-embracing solution. Consumers make independent purchases of their platform technology, their software, their applications, their access provider, and their means of advertising their own capabilities to provide goods and services to others, all as independent decisions, all as a result of this direct exposure to the consumer of every element in the supply chain.

What we have today is an industry structure that is highly diverse, broadly distributed, strongly competitive, and intensely focused on meeting specific customer needs in a price-sensitive market, operating on a quarter-by-quarter basis. Bundling and vertical integration of services has been placed under intense competitive pressure, and each part of the network has been exposed to specialized competition in its right. For consumers this situation has generated significant benefits. For the same benchmark price of around US\$15 to US\$30 per month, or its effective equivalent in purchasing power of a local currency, today's Internet user enjoys multimegabit-per-second access to a richly populated world of goods and services.

The price of this industry restructure has been a certain loss of breadth and depth of the supply side of the market. If consumers do not value a service, or even a particular element of a service, then there is no benefit in incurring marginal additional cost in providing the service. In other words, if the need for a service is not immediate, then it is not provided. For all service providers right through the supply side the focus is on current customer needs, and this focus on current needs, as distinct from continued support of old products or anticipatory support of possible new products, excludes all other considerations.

Why is this change in the form of communications industry operation an important factor in the adoption of IPv6? The relevant question in this context is that of placing IPv6 deployment and dual-stack transition into a viable business model. IPv6 was never intended to be a technology visible to the end user. It offers no additional functions to the end user, nor any direct cost savings to the customer or the supplier. Current customers of ISPs do not need IPv6 today, and neither current nor future customers are aware that they may need it tomorrow. For end users of Internet services, e-mail is e-mail and Web-based delivery of services is just the Web. Nothing will change that perspective in an IPv6 world, so in that respect customers do not have a particular requirement for IPv6, as opposed to a generic requirement for IP access, and will not value such an IPv6-based access service today in addition to an existing IPv4 service. For an existing customer IPv6 and dual stack simply offer no visible value. So if the existing customer base places no value on the deployment of IPv6 and dual stack, then the industry has little incentive to commit to the expenditure to provide it.

Any IPv6 deployment across an existing network is essentially an unfunded expenditure exercise that erodes the revenue margins of the existing IPv4-based product. And as long as sufficient IPv4 address space remains to cover the immediate future needs, looking at this situation on the basis of a quarter-by-quarter business cycle, then the decision to commit to additional expenditure and lower product margins to meet the needs of future customers using IPv6 and dual-stack deployments is a decision that can comfortably be deferred for another quarter. This business structure of today's Internet appears to represent the major reason why the industry has been incapable of making moves on dual-stack transition within a reasonable timeframe as it relates to the timeframe of IPv4 address pool exhaustion.

What of the strident calls for IPv6 deployment? Surely there is substance to the arguments to deploy IPv6 as a contingency plan for the established service providers in the face of impending IPv4 address exhaustion, and if that is the case, why have service providers discounted the value of such contingency motivations? The problem to date is that IPv4 address exhaustion is now not a novel message, and, so far, NAT usage has neutralized the urgency of the message.

The NAT protocol is well-understood, it appears to work reliably, applications work with it, and it has influenced the application environment to such an extent that now no popular application can be fielded unless it can operate across this protocol. For conventional client-server applications, NAT represents no particular problem. For peer-to-peer-based applications, the rendezvous problem with NAT has been addressed through application gateways and rendezvous servers. Even the variability of NAT behavior is not a service provider liability, and it is left to applications to load additional functions to detect specific NAT behavior and make appropriate adjustments to the behavior of the application.

The conventional industry understanding to date is that NAT can work acceptably well within the application and service environment. In addition, NAT usage for an ISP represents an externalized cost, because it is essentially funded and operated by the customer and not the ISP. The service provider's perspective is that considering that this protocol has been so effective in externalizing the costs of IPv4 address scarcity from the ISP for the past 5 years, surely it will continue to be effective for the next quarter. To date the costs of IPv4 address scarcity have been passed to the customer in the form of NAT-equipped CPE devices and to the application in the form of higher complexity in certain forms of application rendezvous. ISPs have not had to absorb these costs into their own costs of operation. From this perspective, IPv6 does not offer any marginal benefits to ISPs. For an ISP today, NATs are purchased and operated by customers as part of their CPE equipment. To say that IPv6 will eliminate NATs and reduce the complexities and vulnerabilities in the NAT service model is not directly relevant to the ISP.

The more general observation is that, for the service provider industry currently, IPv6 has all the negative properties of revenue margin erosion with no immediate positive benefits. This observation lies at the heart of why the service provider industry has been so resistant to the call for widespread deployment of IPv6 services to date.

It appears that the current situation is not the outcome of a lack of information about IPv6, nor a lack of information about the forthcoming exhaustion of the IPv4 unallocated address pool. Nor is it the outcome of concerns over technical shortfalls or uncertainties in IPv6, because there is no evidence of any such technical shortcomings in IPv6 that prevent its deployment in any meaningful fashion. A more likely explanation for the current situation is an inability of a highly competitive deregulated industry to be in a position to factor longer-term requirements into short-term business logistics.

What Next?

Now we consider some questions relating to IPv4 address exhaustion. Will the exhaustion of the current framework that supplies IP addresses to service providers cause all further demand for addresses to cease at that point?

Or will exhaustion increase the demand for addresses in response to various forms of panic and hoarding behaviors in addition to continued demand from growth?

The size and value of the installed base of the Internet using IPv4 is now very much larger than the size and value of incremental growth of the network. In address terms the routed Internet currently (as of 14 August 2008) spans 1,893,725,831 IPv4 addresses, or the equivalent of 112.2 /8 address blocks. Some 12 months ago the routed Internet spanned 1,741,837,080 IPv4 addresses, or the equivalent of 103.8 /8 address blocks, representing a net annual growth of 10 percent in terms of advertised address space.

These facts lead to the observation that, even in the hypothetical scenario where all further growth of the Internet is forced to use IPv6 exclusively while the installed base still uses IPv4, it is highly unlikely that the core value of the Internet will shift away from its predominate IPv4 installed base in the short term.

Moving away from the hypothetical scenario, the implication is that the relative size and value of new Internet deployments will be such that these new deployments may not have sufficient critical mass by virtue of their volume and value as to be in a position to force the installed base to underwrite the incremental cost to deploy IPv6 and convert the existing network assets to dual-stack operation in this timeframe. The corollary of this observation is that new Internet network deployments will need to communicate with a significantly larger and valuable IPv4-only network, at least initially. The fact that IPv6 is not backward-compatible with IPv4 further implies that hosts in these new deployments will need to cause IPv4 packets with public addresses in their packet headers to be sent and received, either by direct deployment of dual stack or by proxies in the form of protocol-translating NATs. In either case the new network will require some form of access to public IPv4 addresses. In other words, after exhaustion of the unallocated address pools, new network deployments will continue to need to use IPv4 addresses.

From this observation it appears highly likely that the demand for IPv4 addresses will continue at rates comparable to current rates across the IPv4 unallocated address pool and after it is exhausted. The exhaustion of the current framework of supply of IPv4 addresses will not trigger an abrupt cessation of demand for IPv4 addresses, and this event will not cause the deployment of IPv6-only networks, at least in the short term of the initial years following IPv4 address pool exhaustion. It is therefore possible to indicate that immediately following this exhaustion event there will be a continuing market need for IPv4 addresses for deployment in new networks.

Although a conventional view is that this market need is likely to occur in a scenario of dual-stacked environments, where the hosts are configured with both IPv4 and IPv6, and the networks are configured to also support the host operation of both protocols, it is also conceivable to envisage the use of deployments where hosts are configured in an IPv6-only mode and network equipment undertakes a protocol-translating NAT function. In either case the common observation is that we apparently will have a continuing need for IPv4 addresses well after the event of IPv4 unallocated pool exhaustion, and IPv6 alone is no longer a sufficient response to this problem.

How?

If demand continues, then what is the source of supply in an environment where the current supply channel, namely the unallocated pool of addresses, is exhausted? The options for the supply of such IPv4 addresses are limited.

In the case of established network operators, some IPv4 addresses may be recovered through the more intensive use of NAT in existing networks. A typical scenario of current deployment for ISPs involves the use of private address space in the customer's network and NAT performed at the interface between the customer network and the service provider infrastructure (the CPE). One option for increasing the IPv4 address usage efficiency could involve the use of a second level of NAT within the service provider's network, or the so-called "carrier-grade" NAT option^[9]. This option has some attraction in terms of increasing the port density use of public IPv4 addresses, by effectively sharing the port address space of the public IPv4 address across multiple CPE NAT devices, allowing the same number of public IPv4 addresses to be used across a larger number of end-customer networks.

The potential drawback of this approach is that of added complexity in NAT behavior for applications, given that an application may have to traverse multiple NATs, and the behavior of the compound NAT scenario becomes in effect the behavior of the most conservative of the NATs in the path in terms of binding times and access. Another potential drawback is that some applications have started to use multiple simultaneous transport sessions in order to improve the performance of the download of multipart objects. For single-level CPE NATs with more than 60,000 ports to be used for the customer network, this application behavior had little effect, but the presence of a carrier NAT servicing a large number of CPE NATs may well restrict the number of available ports per connection, in turn affecting the utility of various forms of applications that operate in this highly parallel mode. Allowing for a peak simultaneous demand level of 500 ports per customer provides a potential use factor of some 100 customers per IP address.

Given a large enough common address pool, this factor may be further improved by statistical multiplexing by a factor of 2 or 3, allowing for between 200 and 300 customers per NAT address. Of course such approximations are very coarse, and the engineering requirement to achieve such a high level of NAT usage would be significant. Variations on this engineering approach are possible in terms of the internal engineering of the ISP network and the control interface between the CPE NATs and the ISP equipment, but the maximal ratio of 200 to 300 customers per public IP address appears to be a reasonable upper bound without unduly affecting application behaviors.

Another option is based on the observation that, of the currently allocated addresses, some 42 percent of them, or the equivalent of some 49 /8 address blocks, are not advertised in the interdomain routing table, and are presumed to be either used in purely private contexts, or currently unused. This pool of addresses could also be used as a supply stream for future address requirements, and although it may be overly optimistic to assume that the entirety of this unadvertised address space could be used in the public Internet, it is possible to speculate that a significant amount of this address pool could be used in such a manner, given the appropriate incentives. Speculating even further, if this address pool were used in the context of intensive carrier-grade NATs with an achieved average deployment level of, say, 10 customers per address, an address pool of 40 /8s would be capable of sustaining some 7 billion customer attachments.

Of course, no such recovery option exists for new entrants, and in the absence of any other supply option, this situation will act as an effective barrier to entry into the ISP market. In cases where the barriers to entry effectively shut out new entrants, there is a strong trend for the incumbents to form cartels or monopolies and extract monopoly rentals from their clients. However, it is unlikely that the lack of supply will be absolute, and a more likely scenario is that addresses will change hands in exchange for money. Or, in other words, it is likely that such a situation will encourage the emergence of markets in addresses. Existing holders of addresses have the option to monetize all or part of their held assets, and new entrants, and others, have the option to bid against each other for the right to use these addresses. In such an open market, the most efficient usage application would tend to be able to offer the highest bid, in an environment dominated by scarcity tending to provide strong incentives for deployment scenarios that offer high levels of address usage efficiency.

It would therefore appear that options are available to this industry to increase the usage efficiency of deployed address space, and thereby generate pools of available addresses for new network deployments. However, the motive for so doing will probably not be phrased in terms of altruism or alignment to some perception of the common good. Such motives sit uncomfortably within the commercial world of the deregulated communications sector.

Nor will it be phrased in terms of regulatory impositions. It will take many years to halt and reverse the ponderous process of public policy and its expression in terms of regulatory measures, and the “common-good” objective here transcends the borders of regulatory regimes. This consideration tends to leave this argument with one remaining mechanism that will motivate the industry to significantly increase the address usage efficiency: monetizing addresses and exposing the costs of scarcity of addresses to the address users. The corollary of this approach is the use of markets to perform the address distribution function, creating a natural pricing function based on levels of address supply and demand.

References

- [1] TCP/IP Mailing List, Message Thread: “Running out of Internet Addresses,” November 1988.
http://www-mice.cs.ucl.ac.uk/multimedia/misc/tcp_ip/8813.mm.www/index.html#121
- [2] F. Solenksy, “Internet Growth,” Steering Group Report, p. 61, Proceedings of the 18th IETF Meeting, August 1990.
<http://www.ietf.org/proceedings/prior29/IETF18.pdf>
- [3] G. Huston, “The IPv4 Internet Report,” August 2008,
<http://ipv4.potaroo.net>
- [4] P. Gross and P. Almquist, “IESG Deliberations on Routing and Addressing,” RFC 1380, November 1992.
- [5] K. Egevang and P. Francis, “The IP Network Address Translator (NAT),” RFC 1631, May 1994.
- [6] G. Tsirtsis and P. Srisuresh, “Network Address Translation – Protocol Translation (NAT-PT),” RFC 2766, February 2000.
- [7] C. Aoun and E. Davies, “Reasons to Move the Network Address Translator – Protocol Translator (NAT-PT) to Historic Status,” RFC 4966, July 2007.
- [8] M. Bagnulo, P. Matthews, and I. van Beijnum, “NAT64/DNS64: Network Address and Protocol Translation from IPv6 Clients to IPv4 Servers,” Internet Draft, work in progress, **draft-bagnulo-behave-nat64-00.txt**, June 2008.
- [9] T. Nishitani and S. Miyakawa, “Carrier Grade Network Address Translator (NAT) Behavioral Requirements for Unicast UDP, TCP and ICMP,” Internet Draft, work in progress, **draft-nishitani-cgn-00.txt**, July 2008.

- [10] Olaf Maennel, Randy Bush, Luca Cittadini, Steven M. Bellovin, “A Better Approach than Carrier-Grade-NAT,”
<http://rip.psg.com/~randy/080820.alt-to-cgn.pdf>
- [11] William Lehr, Tom Vest, Eliot Lear, “Running on Empty: The Challenge of Managing Internet Addresses,” to be presented at the 36th Research Conference on Communication, Information and Internet Policy (TPRC), on 27 September 2008.
http://eyeconomics.com/backstage/References_files/Lehr-Vest-Lear-TPRC2008-080915.pdf
- [12] Hain, Tony, “A Pragmatic Report on IPv4 Address Space Consumption,” *The Internet Protocol Journal*, Volume 8, No. 3, September 2005
- [13] <http://icann.org/en/announcements/proposal-ipv4-report-29nov07.htm>
(See also “Fragments” on page 46.)

GEOFF HUSTON is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He graduated from the Australian National University with a B.Sc. and M.Sc. in Computer Science. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005; he served on the Board of Trustees of the Internet Society from 1992 to 2001.
E-mail: gih@apnic.net

Letters to the Editor

I sincerely congratulate you for Geoff Huston's excellent article in *The Internet Protocol Journal*, June 2008, on the "Decade of Internet Evolution." The article shows an amazing insight into the Internet as it has recently evolved and deserves as wide an audience as possible.

The only comment I could make is that though Huston hints about separating the IP address from the host name, he does not explicitly mention the *Host Identity Protocol* (HIP)^[1]. Previous issues of the Journal have this omission as well.

Note: As we struggle in the IETF and everywhere else in the industry with NAT traversal, mobility, and multihoming, we see countless approaches for each application layer protocol separately. HIP seems to fulfill the promise of solving these problems comprehensively.

Thanks for the privilege to continue reading the Journal; keep such papers coming.

—Henry Sinnreich, Adobe Systems, Inc.
hsinnrei@adobe.com

- [1] R. Moskowitz, P. Nikander, P. Jokela, Ed., and T. Henderson, "Host Identity Protocol," RFC 5201, April 2008. See also: <http://www.ietf.org/html.charters/hip-charter.html>

The author responds:

Thank you for your generous comments.

At some point I was toying (dangerously!) with writing an article that attempted to predict the next 10 years, looking at what appears to be important today and what that could mean in the future. There is no doubt that the tight binding of identity and location is one of the assumptions that has made the Internet both simple and effective for the past decade. But where we sit today, in a world dominated by scale, mobility, a dense mesh of interconnectivity, highly capable end devices, dense middleware, and a panoply of specialized requirements, we need to look forward to methods that allow separation of identity and location. Now this separation could be at the level of the Internet Protocol itself, as in HIP or *Site Multihoming by IPv6 Intermediation* (SHIM6); or at the level of the transport session, as exemplified at present by the *Stream Control Transmission Protocol* (SCTP); or even at the application level, where the various offerings related to *Voice over IP* (VoIP) and *Peer-to-Peer* (P2P) have been working at the level of multiparty application rendezvous and application identity that sit on top of an adaptive platform of dynamic discovery of the characteristics of the underlying transport subsystem.

Each approach appears to offer some significant leverage in scaling the network in diverse ways, while at the same time presenting us with some fascinating insights into possible architectures that could address our needs in the next decade. No doubt the next 10 years will present us with some quite novel challenges with the imminent exhaustion of the unallocated IPv4 address pool and the associated observation that the schedule for the update of IPv6 has proceeded so slowly that we will be forced to be remarkably inventive with IPv4. HIP may well be a central part of such invention, but, more generally, I have no doubt that we will examine more generally how we can devise refinements to the networking model that preserve useful notions of identity across a rather fluid sea of shared location tokens.

Regards,

—*Geoff Huston, APNIC*
gih@apnic.net

Ten Years of IPJ

We received many congratulatory messages in response to our June 2008 Anniversary Issue. The following are some quotes from our readers:

“Compliments and congratulations for the tenth anniversary of this great Journal. It is great because it is making us realize the synergy between what has been and what is to come.”

—*John Okewole, Lagos, Nigeria*

“This week I received the June 2008 issue of IPJ. I have been a subscriber for several years and it has been a great pleasure to find great contents in IPJ, such as the current issue that brings reviews on Internet evolution. I would like to send my congratulations to the IPJ team for 10 years of publication and my best wishes for future success.”

—*Frederico Fari, Belo Horizonte, Brazil*

“I think that IPJ is a great journal. I hope you will not be forced to give up the paper edition because is a beautiful one (and it allows me to read during the evening hours when all computers and children in the house are shut down :-)”

—*Andrea Montefusco, Rome, Italy*

Book Reviews

Two Books on Cyber Law

Code and Other Laws of Cyberspace

Code and Other Laws of Cyberspace, by Lawrence Lessig, Basic Books, 1999, ISBN 0-465-03913-8. <http://code-is-law.org/>

Code 2.0

Code 2.0, by Lawrence Lessig, Basic Books, 2006, ISBN-10: 0-465-03914-6, ISBN 13: 978-0-465-03914-2. <http://codev2.cc/>

First published in 1999, then Harvard Law School Professor Lawrence Lessig's cautionary tale about the inescapable influence of certain material features of the built Internet has since become a foundational "Internet studies" text in universities and laws schools around the world. Lessig, who now occupies an endowed chair at Stanford Law School, makes a series of troubling observations about the Internet, his chosen sector of focus since setting aside his mid-1990s work on legal and institutional development in post-Soviet societies.

Lessig's key findings from that previous work are that rules matter—especially the sort of rules embodied in "constitutions" and other foundational institutions; that rules are artifacts of contingent human intent and design; and that rules can be changed. Being a "classical liberal" on the model of John Stuart Mill, Lessig advocates the sort of rules that afford maximum liberty for individuals against a triumvirate of coercive influences, including not only governments but also market power and oppressive social mores.

Now however, a fourth challenge to personal liberty has been exposed by the advent of the Internet—or rather, of *cyberspace*, which Lessig describes as the lived experience of participants in the rich application space that has been built atop the Internet. This new constraining factor is "architecture," which Lessig defines as "the built environment," or "the way the world is," that is, the cumulative result of all of the contingent historical events and decisions that have shaped the material circumstances confronting Internet users (or *cyberspace denizens*) today. *Code* is Lessig's term for the instruction sets (that is, programs, applications, etc.) that are the building blocks of the architecture of cyberspace; it is the stuff that emerges from the decision making of a relatively few (the *code writers*), which accretes over time into the less-malleable architecture that shapes the everyday choices and possibilities of everyone else whom the Internet or cyberspace touches.

New Code Means New Power(s)

According to Lessig, the code that defines cyberspace—which he calls "West Coast Code"—demands particular attention, both because of its omnipresence and because of how it differs from the other, more familiar factors that can impinge on individual liberty.

Like the canons of law (also known as “East Coast Code”), code is basically a collection of rules written with human goals and objectives in mind. However, in its effects code more closely resembles the laws of nature, because it requires neither the awareness nor the consent of its subjects in order to be effective. Although this claim sounds suspiciously like a variant, or perhaps an illustration of Arthur C. Clarke’s *Third Law of Prediction* (which states that any sufficiently advanced technology will be indistinguishable from the supernatural), there is purpose behind Lessig’s observation. The self-enforcing character of code is doubly problematic in the case of cyberspace, he suggests, because unlike the law, code affords no appeal, no recourse, and no formal, institutional review and interpretation of the kind that lawyers and judges exercise in legal matters. Without such expert oversight, code might come to be used as a tool to subvert individual liberties or public values, for either commercial or political gain, without anyone’s being the wiser. In fact, he implies, the lack of transparency of code almost invites such abuses.

At this point some might be tempted to dismiss Lessig’s program as just “sour grapes” from a high-profile industry spokesman sensing this erosion of the traditional prominence and centrality of his profession in a new code-centric world. Lessig believes passionately in the exercise of law and judicial review as master tools for keeping other important forces—government power, market power, and social norms—broadly aligned with “important public values.” He extols the relationships among the rule of law, democracy, and politics, the latter of which invests law with legitimacy to raise or lower the cost of particular individual actions (for example, by taxing, criminalizing, valorizing, or subsidizing them) to encourage conformity with publicly chosen goals and values. He observes that “architecture is a kind of law” and that “code codifies values, and yet, oddly, most people speak as if code were just a question of engineering.” It takes no great leap of imagination to conclude that code too should be subject to the same kind of legal and judicial oversight that keeps the rest of society running smoothly. Eliminating any doubt, Lessig asserts that:

Technology is plastic. It can be remade to do things differently. We should expect—and demand—that it can be made to reflect any set of values that we think important. The burden should be on the technologists to show us why that demand can’t be met.

However, such a dismissal would indeed be too easy, for Lessig also expresses misgivings about the professionalization and segregation of “constitutional thinking” within the legal sector. “Constitutional thought has been the domain of lawyers and judges for too long,” Lessig writes, and as a result everyone else has grown less comfortable—and also less competent—in engaging in fruitful conversation about fundamental, “constitutional” values.

And yet Lessig suggests that this skill has also atrophied within the legal community, as more and more jurists have embraced an “originalist” interpretive philosophy that holds that the U.S. Constitution provides no guidance for how to resolve conflicts between old values—what Lessig calls *latent ambiguities*—or how to address wholly novel concerns raised by technologies such as the Internet. Originalists (Lessig mentions U.S. Supreme Court Justice Antonin Scalia) assert that in such cases the only recourse is the political and legislative processes—where, one assumes, limited experience with both technology and constitutional debate make the prospects for success even dimmer. Lessig writes that “We (legal scholars) have been trapped by a mode of reasoning that pretends that all the important questions have already been answered,” but that “the constitutional discourse of our present Congress is far below the level at which it must be to address the questions about constitutional values that will be raised by cyberspace.”

Diagnosis from a Distance

Lessig is without question eminently qualified to make such observations about his home-turf legal and political spheres. However, it is less clear that his blanket charge of deliberative incompetence is equally valid across the full range of Internet and cyberspace stakeholders. Neither is it clear that the architecture of cyberspace is as uniquely problematic as he suggests, compared to the architecture of other, more familiar domains. Finally, Lessig’s own admittedly limited technical expertise may lead him to misapprehend the boundary between cyberspace and the Internet, and to underestimate the radicalness of his proposed cyberspace fix.

Taking these ideas in reverse order, Lessig’s conception of the structural and functional distinction between the Internet and cyberspace merits closer scrutiny. As explained later, Lessig advocates profound technical changes to bring the functions of code under the rule of law (or laws, because Lessig wishes to accommodate subsidiary jurisdictions as well as sovereign differences in law). However, he envisions this intervention affecting only the “code” domain, not the “Internet’s core protocols”:

When I speak about regulating the code, I’m not talking about changing these core TCP/IP protocols...In my view these components of the network are fixed. If you required them to be different, you’d break the Internet. Thus rather than imagining the government changing the core, the question I want to consider is how the government might either (1) complement the core with technology that adds regulability, or (2) regulate applications that connect to the core.

Lessig's specific ideas for achieving this function while preserving the core are not fully detailed in this context until *Code 2.0* (2006), which Lessig describes as an update rather than a full rewrite, albeit one with new relevance to match a "radically different time." The central idea involves the introduction of an "identity layer" that permits authoritative in-band querying and signaling of the jurisdiction(s) to which every would-be Internet user is subject. The deployment of this system would be accompanied by the development of a comprehensive distributed database of Internet usage restrictions mandated by every legally recognized jurisdiction around the world. Together, these components would operate as a kind of "domain interdiction system" that would automatically black-hole all Internet resource queries that are legally impermissible to individuals based on their jurisdiction(s) of origin, regardless of their actual location.

This proposal is clearly vulnerable to criticism of many kinds—technical, ethical, practical, etc.—and to be fair Lessig anticipates and preemptively responds to several of the most obvious ones. Space limitations preclude any review of those arguments here, but it is impossible to resist a few short observations. First, it is not clear why Lessig imagines that his proposed system would be anything less than a fundamental intervention in the core function and protocols of the Internet. Today several different high-profile technical developments that could plausibly be described as changing TCP/IP are under way, but they (hopefully) will not break the Internet. At the same time, TCP/IP is not the only technology that is essential to the Internet "core." The system that Lessig advocates is clearly inspired by the *Domain Name System* (DNS), it would of necessity be similarly global and ubiquitous in scope and scale, and it would likely function by selectively blocking some DNS responses based on the initiator's identity. Although some once regarded the DNS as a mere application (for example, shortly after it was invented), few today would categorize it as anything other than a core protocol. Also, given the degree to which any implementation of the proposed identity system would preempt many "normative" features that are associated with the Internet core (for example, the principles behind the *end-to-end* arguments), it is unclear what would remain "unbroken" therein that might still warrant any special consideration or separate treatment. We can only hope that Lessig's optimism on this question is justified, because looming developments in certain wireless standards as well as in the management of IP addressing may provide for more concrete—and less revisable—answers in the very near future.

Objects in View May Be Closer Than They Appear

Then there is the question of how much code really makes the architecture of cyberspace different from the architecture of other domains. Many of Lessig's claims on this point date back to the first version of the book, when Internet exceptionalism was still new enough for deflationary counterarguments to seem provocative.

Although the revolutionary potential of the Internet continues to inspire many (this reviewer included), the past decade of booms, busts, compromises, and indictments have done much to temper that faith. It is not that Lessig's concerns about the opaque nature of cyberspace architecture, about the substantial influence that code writers and network owners command, and about the vulnerability of the whole system to a crisis-induced authoritarian turn aren't reasonably well-founded. But they are equally apropos to most other important spheres of life. The phrase "possession is nine-tenths of the law" has multiple meanings, and was coined many decades before the Internet was invented. The inexplicability of many current "real-world" legislative and judicial outcomes without recourse to some cynical theory of unacknowledged interests and unobservable influence certainly raises many questions about the architecture of the space beyond cyberspace. And Lessig's warnings about national security fears precipitating a sudden loss of freedoms (taken from Jonathan Zittrain's *Z-Theory*) now seem prophetic—albeit less for the Internet than for the earliest and largest host society of the Internet. One might observe that Lessig is guilty of his own kind of exceptionalism—one that, ironically, may obscure the degree to which constitutional challenges in the real and virtual worlds are more or less the same. In fact, Lessig's subsequent shift of priorities from code to intellectual property law recently ended with a return to his original home turf of law and politics—perhaps in belated recognition that sometimes, even when you have a good story, East Coast Code is still the only durable recourse.

Finally, there is the question of constitutional acumen. This question is the critical one for Lessig (he uses some form of the term *constitution* more than 250 times in the main text), because for him the term evokes nothing less than "an architecture... a way of life that structures and constrains social and legal power, to the end of protecting fundamental values." In this sense, he adds, constitutions are built rather than found. Moreover, they have been built in different (albeit sometimes overlapping) places by different institutions and societies, many with quite different conceptions of which fundamental values to uphold. From whence will the architecture of values of cyberspace emerge? Who will be its authors? Lessig never quite gives a final answer, even for his own home jurisdiction, but he does help to winnow out several likely suspects. As noted previously, he invests little faith in the current U.S. legislative branch. He also has reservations about many members of his own legal profession, although the need to preserve backward compatibility with the primary U.S. Constitution and to reconcile newly revealed "latent ambiguities" therein obviously recommends some legal training at the very least. Government and industry represent the most likely perpetrators of liberty-undermining code, Lessig claims, so he looks for no help from those quarters.

In the end Lessig provides some oblique advice for judges (abandon formalism), hackers (open source), and voters (educate yourself, and don't give up hope), but ultimately concludes with a call for more lawyerly deliberation: if only our leaders could act more like lawyers, telling stories that persuade "not by hiding the truth or exciting the emotion, but by using reason," and our fellow citizens could act like juries, resisting the fleeting passions of the mob and making decisions based on the facts alone, then perhaps we could overcome the architectural challenges of both cyberspace and physical space.

Story Boards and Internet Constitutions

Notwithstanding its solipsistic aspects, advice like that discussed in the last section is hard to find fault with. Professor Lessig is unquestionably a person of good conscience, and has a long, distinguished, and very well-documented record of putting this advice into practice in a wide range of good causes, including many that are wholly unrelated to code or cyberspace. However, one could argue (perhaps with equal solipsism) that many of the behaviors and virtues that he commends are now regularly on display in the mailing lists, message boards, and other deliberative records of the Regional Internet Registries, the IETF, and the IAB—in particular in discussions on the form that IPv4 and IPv6 address-allocation policies should take, in the design of future routing systems that balance scalability with the freedom to choose between competing providers, and in the reconciliation of traditional policies and their beneficiaries with the changing realities of Internet resource stewardship. Closer scrutiny of these records reveals that successful consensus policies are almost invariably borne of good, well-reasoned stories, the vast majority of which are offered by individuals who are affiliated neither with government agencies nor with any of the largest and most powerful ISPs. Many of the storytellers are old hands, but new voices regularly emerge and command attention based on nothing more than the strength of their reasoning. Participating in these discussions, one can *occasionally* experience the same feeling that inspires Lessig in the courtroom, where "some, for the first time in their lives, see power constrained by reason. Not by votes, not by wealth, not by who someone knows—but by an argument that persuades."

That this "architectural" work has gone largely unrecognized to date in law schools, university humanities and social science departments, and even in some civil society-oriented Internet governance fora is not entirely unexpected, because the context and terminology of those discussions is invariably technical, even if many participants recognize that the underlying principles are essentially "constitutional" in nature. No doubt a more complete conversation between code writers and constitutionalists is inevitable over time, and with luck more cross-fertilization will lead to better protocols, better policies, and better architecture.

However, this rapprochement is unlikely to be initiated by technologists seeking to take up the study and application of legal principles. Lessig, whose own intellectual project builds substantially on the antiformalist, “legal realist” school of thought, should understand this reality better than most. In the crudest of forms, legal realism holds that “the Law is whatever lawyers happen to say it is.” Stated as neither a boast nor a claim of entitlement but rather as a practical observation of the challenges that lawyers face in applying ambiguous old laws to incommensurable new circumstances, this maxim nevertheless clearly conveys a sense of both the great responsibility and the great power that lawyers command. Perhaps it is time that Mr. Lessig and his counterparts consider the possibility that a similar school of thought may inform (consciously or unconsciously) the perspectives of network builders and code writers. Being of no less good conscience, perhaps code writers and other “cyberspace realists” are merely waiting for the moment when the Law and lawyers come calling with a good story, under the banner of reason rather than power. So long as the story now unfolding continues to make sense and satisfy the ever-expanding audience, we needn’t fear either.

Code may not be *that* particular story, but it’s an excellent read, and an important contribution to a dialogue that must be engaged.

—Tom Vest
tvest@eyeconomics.com

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. Contact us at ipj@cisco.com for more information.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Global Policy Proposal for Remaining IPv4 Address Space

Global Internet Number Resource Policies are defined by the *Address Supporting Organization* (ASO) MoU^[1]—between the *Internet Corporation for Assigned Names and Numbers* (ICANN) and the *Number Resource Organization* (NRO)—as “Internet number resource policies that have the agreement of all RIRs according to their policy development processes and ICANN, and require specific actions or outcomes on the part of the *Internet Assigned Numbers Authority* (IANA) or any other external ICANN-related body in order to be implemented.” Attachment A of this MoU describes the *Development Process of Global Internet Number Resource Policies*, including the adoption by every *Regional Internet Registry* (RIR) of a global policy to be forwarded to the ICANN Board by the ASO, as well as its ratification by the ICANN Board. In this context, the ICANN Board adopted its own Procedures^[2] for the Review of Internet Number Resource Policies Forwarded by the ASO for Ratification.

Among other features, these Procedures state that the Board will decide, as and when appropriate, that ICANN staff should follow the development of a particular global policy, undertaking an “early awareness” tracking of proposals in the addressing community. To this end, staff should issue background reports periodically, forwarded to the Board, to all ICANN Supporting Organizations and Advisory Committees and posted at the ICANN Web site.

At its meeting on 20 November 2007, the Board resolved to request tracking of the development of a global policy proposal for allocation of remaining IPv4 address space, under discussion in the Regional Internet Registries. The status overview presented below is compiled in response to this request and will be further updated as developments proceed, for information to ICANN entities and the wider community. This is the fifth issue of the tracking of this policy.

Originally, two slightly different global policy proposals were introduced for allocation of the remaining IPv4 address space:

- A version (1) “Global Policy for the Allocation of the Remaining IPv4 Address Space,” first presented at LACNIC X in May 2007
- A version (2) “End Policy for IANA IPv4 allocations to RIRs,” first presented at APNIC 24 in September 2007

Both featured the same approach, distribution of an equal number N of /8 IPv4 address blocks to each RIR when the IANA free pool would reach the threshold value of $5 \times N$, but differed in the proposed value of N , notably 2 or 1, respectively. The proposals were discussed in parallel in the RIRs and regarded essentially as one proposal, with a view to converging on a value for N . In February 2008, agreement was reached for a unified proposal (3).

The current proposal is thus:

- Version (3) “Global Policy for the Allocation of the Remaining IPv4 Address Space,” first presented at APNIC 25 in February 2008.

The proposal was introduced at the subsequent meetings of all other RIRs. It has now been adopted in ARIN, AfriNIC, LACNIC and RIPE, and is in final call in APNIC. If adopted by all the RIRs, the proposal will subsequently be handled by the NRO Executive Council and the ASO Advisory Council according to their procedures before being submitted to the ICANN Board for ratification. A table^[3] can be found on the ICANN Website that indicates the status within each RIR for the current proposal. Hyperlinks are included for easy access.

It should be noted that other policy proposals have been put forward and are being discussed regarding IPv4 address space exhaustion, although only those mentioned above have been scoped as global policy proposals in the sense of the ASO MoU, that is, focusing on address allocation from IANA to the RIRs, and recognized by the ASO AC as global policy proposals in that meaning.

[1] <http://aso.icann.org/docs/aso-mou2004.html>

[2] <http://icann.org/en/general/review-procedures-pgp.html>

[3] <http://www.icann.org/en/announcements/proposal-ipv4-report-29nov07.htm>

Upcoming Events

The *Internet Engineering Task Force* (IETF) will meet in Minneapolis, Minnesota, November 16 – 21, 2008. In 2009, IETF meetings are scheduled for San Francisco, California (March 22 – 27), Stockholm, Sweden (July 26 – 31) and Hiroshima, Japan (November 8 – 13). For more information see <http://www.ietf.org/>

The *North American Network Operators’ Group* (NANOG) will meet in Los Angeles, California, October 12 – 14. Immediately following the NANOG meeting, the *American Registry for Internet Numbers* (ARIN) will meet in the same location, October 15 – 17. See <http://nanog.org> and <http://arin.net>

The *Internet Corporation for Assigned Names and Numbers* (ICANN) will meet in Cairo, Egypt, November 2 – 7, 2008. For more information see: <http://icann.org>

The *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) will be held in Manila, Philippines, February 18 – 27, 2009. See: <http://www.apricot2009.net/>

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2008 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA
--

The Internet Protocol Journal

December 2008

Volume 11, Number 4

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Wi-Fi, Bluetooth and WiMAX.....	2
The End of Eternity	18
Remembering Jon	29
Letters to the Editor.....	33
Book Reviews	36
Fragments	41
Call for Papers.....	43

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

FROM THE EDITOR

Response to our use of a new printing paper has been very positive, so we will continue to use the uncoated and recycled Exact® paper introduced with our September 2008 issue. We are still interested in hearing your feedback on the paper, as well as any other aspect of this journal. Send your comments to: ipj@cisco.com

The last decade has seen many developments in the area of *wireless* networking technologies. Wireless Internet access is now available in thousands of locations ranging from private homes to hotels, trains, airplanes, ships at sea, and even entire cities. Wireless systems, specifically Bluetooth, are also used for short-range device connectivity such as between a mobile phone and a headset, while WiMAX systems are being deployed for larger area coverage. In our first article, T. Sridhar gives an overview of Wi-Fi, Bluetooth, and WiMAX.

As stated in our previous issue, the topic of IP Version 4 address exhaustion and migration to IP Version 6 is being debated in many Internet-related organizations, including the IETF, *Internet Corporation for Assigned Names and Numbers* (ICANN), and the *Regional Internet Registries* (RIRs). In our last issue, Geoff Huston outlined the history of IPv4 address depletion. This time we bring you the first in a two-part series of articles entitled “The End of Eternity.” The article is by Niall Murphy and David Wilson. Part Two will follow in our March 2009 issue. As you will see from our “Letters to the Editor,” views on the right way to tackle the address exhaustion and protocol migration challenge abound, and I predict we will carry yet more articles on this topic in future issues.

Just over 10 years ago, Jonathan B. Postel, Internet pioneer and a key player in many core Internet activities, passed away. In this issue we bring you a remembrance article written by another Internet pioneer, Vint Cerf. In connection with this anniversary, special events were held in Minneapolis in conjunction with the 73rd meeting of the IETF. The *Jonathan B. Postel Service Award* for 2008 was awarded to EsLaRed of Venezuela by a committee of former award winners. You will find more information about the award in our “Fragments” section on page 42.

Remember to let us know if your mailing address changes and to visit our online companion, *The Internet Protocol Forum*, where you will find additional articles and other material: <http://ipjforum.org>

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

Wi-Fi, Bluetooth and WiMAX—Technology and Implementation

by T. Sridhar, Flextronics

Wireless networks can be classified broadly as *Wireless Personal-Area Networks* (WPAN), *Wireless LANs* (WLANs), and *Wireless Wide-Area Networks* (WWANs). WPANs operate in the range of a few feet, whereas WLANs operate in the range of a few hundred feet and WWANs beyond that. In fact, wireless WANs can operate in a wide range—a metropolitan area, cellular hierarchy, or even on intercity links through microwave relays.

This article examines wireless technologies for the WLAN, WPAN, and WWAN areas, with specific focus on the IEEE 802.11 WLAN (often known as Wi-Fi®), *Bluetooth* (BT) in the WPAN, and WiMAX for WWAN as representative technologies. It discusses key aspects of the technology—medium access and connectivity to the wired network—and concludes by listing some common (mis)perceptions about wireless technology.

WLANs

The *Institute of Electrical and Electronic Engineers* (IEEE) defined three major WLAN types in 802.11–802.11 b and g, which operate in the 2.4-GHz frequency band, and 802.11a, which operates in the 5-GHz band. The 2.4- and 5-GHz bands used here are in the license-free part of the electromagnetic spectrum, and portions are designated for use in *Industrial, Scientific, and Medical* (ISM) applications—so these portions are often called ISM bands. More recently, a high-speed 802.11 WLAN has been proposed—the 802.11n WLAN, which operates in both the 2.4- and 5-GHz bands.

The 2.4-GHz frequency band used for 802.11 is the band between 2.4 and 2.485 GHz for a total bandwidth of 85 MHz, with 3 separate nonoverlapping 20-MHz channels. In the 5-GHz band, there are a total of 12 channels in 3 separate subbands—5.15 to 5.25 GHz (100 MHz), 5.25 to 5.35 GHz (100 MHz), and 5.725 to 5.825 GHz (100 MHz).

The more common mode of operation in 802.11 is the *infrastructure* mode, where the stations communicate with other wireless stations and wired networks (Ethernet typically) through an *access point*. The other mode is the *ad-hoc* mode, where the stations can communicate directly with each other without the need for an access point; we will not discuss this mode in this article. The access point bridges traffic between wireless stations through a lookup of the destination address in the 802.11 frame (see Figure 1a).

The *Media Access Control* (MAC) header of 802.11 has four addresses. Depending upon the value of a *FromDS* (from access point), or a *ToDS* (to access point) bits in the header (see Figure 1b), the addresses have different connotations. The first two addresses are for the receiver and transmitter, respectively.

Figure 1a: WLAN Network with Ethernet Connectivity

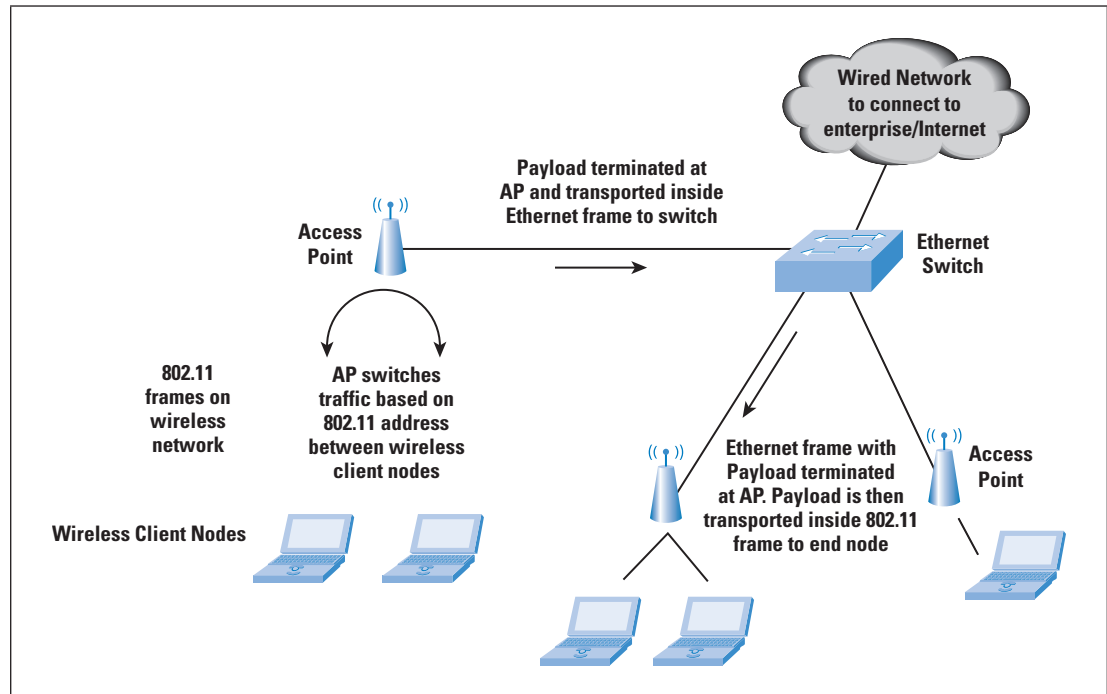
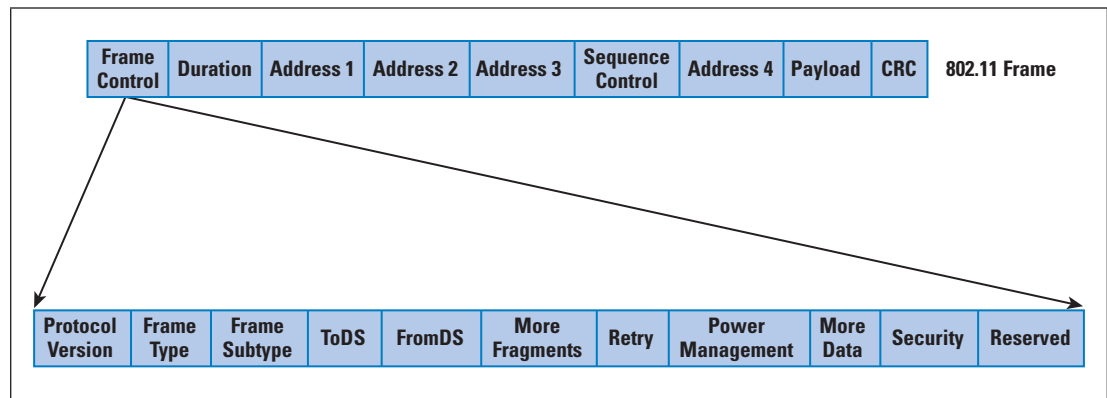


Figure 1b: 802.11 Frame Format



Address 4 is not used except when both FromDS and ToDS are set to 1—it is for a special mode of communication for access point-to-access point traffic, whence addresses 3 and 4 refer to the source- and destination-station MAC addresses, respectively, whereas addresses 1 and 2 refer to the access point addresses (that is, the transmitter and receiver on this inter-access point channel). When FromDS is set to 1, address 1 is the destination-station MAC address, address 2 is the access point address, and address 3 is the source-station MAC address. When ToDS is set to 1, address 1 is the access point MAC address, address 2 is the transmitting-station MAC address, and address 3 is the destination-station MAC address.

Although earlier versions of 802.11 LANs used *Frequency Hopping Spread Spectrum* (FHSS), 802.11b typically uses *Direct Sequence Spread Spectrum* (DSSS) for 1-, 2-, 5.5-, and 11-Mbps speeds. Both schemes involve transmission of a narrowband signal over a wider frequency range to mitigate the possibility of interference at any one frequency. The nodes and access points typically transmit at the highest data rate possible based on the current signal-to-noise ratio.

At the MAC level, 802.11 LANs involve the use of *Carrier Sense Multiple Access/Collision Avoidance* (CSMA/CA). Stations back off if they detect that another station is transmitting on that channel. The station then waits for a random period after the end of the transmission before it attempts to transmit on that channel. In addition, control frames such as *Request to Send* (RTS) and *Clear to Send* (CTS) are used to facilitate the actual data transfer. The CTS control frame has the duration for which the transmitting node is allowed to transmit. Other stations sense this frame and back off for at least the specified duration before sensing the radio link again.

When the access points are connected through a LAN, the entire system is known as a *Distribution System*. The access points perform an integration function—that is, bridging between wired and wireless LANs. In this scenario, (see Figure 1a) the wireless control and data frames are terminated at the access point or tunneled from the access point to a centralized controller over Ethernet. When terminated at the access point, the payload is transmitted from the access point to the network over Ethernet. This transmission is done in the following manner:

The source and destination addresses are set to the station and access point addresses, respectively. At the access point, the payload is stripped from the 802.11 data frame and sent as part of an Ethernet packet either as a broadcast packet or to a specific destination. If the packet sizes (when reassembled) are larger than the Ethernet frame size, they are discarded. In the reverse direction, the Ethernet frame can be directly encapsulated into an 802.11 frame for transmission from the access point to the end node. At the WLAN end node, the complete Ethernet frame shows up at the driver level as though it were a frame received on a pseudo Ethernet interface.

The most common 802.11b WLAN speed is 11 Mbps. However, based on the interframe spacing, preamble, header encapsulation, and acknowledgements for frames required, the actual throughput for user data would be about 50 percent of the actual speeds. This throughput of 50 percent of actual link speed is a common theme on 802.11g and 802.11a also.

Stations connect to the access point through a scanning process. Scanning can be passive or active. In the passive mode, the station searches for access points to find the best access point signal (which contains the *Service Set Identifier* [SSID], data rates, and so on).

The access point frame that the stations look for is a management frame known as the *beacon frame*. In the active mode, the station initiates the process by broadcasting a probe frame. All access points that receive the probe send back a probe response, helping the station build up the list of available access points. The sequence of a station “connecting” to an access point involves two steps. The first is *authentication*, where the station sends an authentication request frame to the access point. Depending upon the authentication through 802.1X or internal configuration, the access point can accept or reject the request with an authentication response. The second step is *association*, which is required to determine the data rates supported between the access point and the station. At the end of the association phase, the station is allowed to transmit and receive data frames.

Power Concerns in 802.11

Although it is not a part of the standard, the access points might adjust their transmitting power based on the environment they are in (they do have maximum limits based on regional restrictions). If they do not perform this adjustment, all the stations might connect to the access point with the highest transmitting power, even if the access point is far away. The other concern is, of course, the interference between access points. The power adjustment is usually done through configuration and, in some cases, through a monitoring function on the network. In the latter case, the monitoring function reports the information to a central controller.

A new initiative within the IEEE (802.11k) has been started to improve traffic distribution within the network. Specifically, it addresses the problem of access point overloading so that stations can connect to underused access points for a more efficient use of network resources.

With respect to power management on the client side, a station can indicate that it is going into a “sleep” or low-power state to the access point through a status bit in a frame header (refer to Figure 1b). The access point then buffers packets for the station instead of forwarding them to the station as soon as they are received. The sleeping station periodically wakes up to receive beacons from the access point. The beacons include information about whether frames are being buffered for the station. The station then sends a request to the access point to send the buffered frames. After receiving the frames, the station can go back to sleep.

802.11a/g Technology—Orthogonal Frequency-Division Multiplexing

Sometimes called *discrete multitone* (DMT) in the *Digital Subscriber Line* (DSL) world, *Orthogonal Frequency-Division Multiplexing* (OFDM) is used as the underlying technology in 802.11g and 802.11a. OFDM is a form of *Frequency-Division Multiplexing* (FDM); normally, FDM uses multiple frequency channels to carry the information of different users. OFDM uses multicarrier communications, but only between one pair of users—that is, a single transmitter and a single receiver.

Multicarrier communications splits a signal into multiple signals and modulates each of the signals over its own frequency carrier, and then combines multiple frequency carriers through FDM. OFDM uses an approach whereby the carriers are totally independent of (orthogonal to) each other. Note that the total bandwidth consumed with OFDM is the same as with single carrier systems even though multiple carriers are used—because the original signal is split into multiple signals. OFDM is more effective at handling narrowband interference and problems related to multipath fading, simplifying the building of receiver systems.

We can illustrate this process with a simple example—one often used in discussions about OFDM. For a “normal” transmission at 1 Mbps, each bit can take 1 microsecond to send. Consider bit 1 and bit 2 sent with a gap of 1 microsecond. If two copies of bit 1 are received at the destination, one of them is the reflected or delayed copy. If the delay is around 1 microsecond, this delayed copy of bit 1 can interfere with bit 2 as it is received at the destination because they arrive at approximately the same time. Now consider an OFDM transmission rate of 100 kbps, that is, the bits are sent “slower” but over multiple frequencies. A multipath delay of around 1 microsecond will not affect bit 2, because bit 2 is now arriving much slower (around 10 microseconds). The delay in bit arrival (1 microsecond in our example) is not a function of the transmission—rather it is due to the various paths taken by the signal.

Orthogonal Frequency-Division Multiple Access (OFDMA) superimposes the multiple-access mechanism on OFDM channels, so that multiple users can be supported through subsets of the subcarriers assigned to different users. Note that 802.16-2004 (“Fixed” WiMAX) uses OFDM, whereas 802.16e-2005 (“Mobile” WiMAX) uses OFDMA.

MIMO and 802.11n

Multiple Input Multiple Output (MIMO) antennas are the basis for the 802.11n wireless LAN standard, currently in draft form but on the way to final standardization. Signals often reflect off objects and are received at different times and strengths at the receiver, resulting in a phenomenon called *multipath distortion*. (Note: 802.11n in this article implies the draft 802.11n standard at the time of writing.) MIMO actually takes advantage of this distortion by sending a single data stream split into multiple parts to be transmitted from multiple antennas (typically 3 in 802.11n) and letting the reflected signals be processed at the receiver (through multiple antennas). The transmission of multiple data streams over different spatial channels, sometimes known as *Space Division Multiplexing* (SDM), also allows a larger amount of data to be sent over the air. Through advances in the *Digital Signal Processing* (DSP)-based processing, the receiver can process the signals, cross-correlate them, and reconstitute them accurately despite interference. Also, because of the multiple signals received over multiple paths, link reliability is increased.

The 802.11n standard uses three antennas and also supports two radios (for the 2.4- and 5-GHz bands where 802.11n can operate). It can also use 40-MHz channels through *channel bonding*—that is, two adjacent 20-MHz channels are combined into a single 40-MHz channel, possibly resulting in a data rate of up to 150 Mbps of effective throughput.

One concern with 802.11n that is starting to gain attention is the power requirement of 802.11n access points. With radios in both bands and the use of MIMO, 802.11n access points tend to consume more power than the 802.11 a/b/g access points, leading to problems when the access point is powered by *Power over Ethernet* (PoE) power-sourcing equipment. The 802.3af standard permits a maximum of 12.95W per Ethernet port, which is often less than the power that most 802.11n APs need. The IEEE 802.3at working group is working toward a higher-power PoE standard. This initiative, commonly called *PoE Plus*, will peak at 25W per Ethernet port (on Category 5 Ethernet cable).

Ethernet Backhaul

The access point has two primary functions—connecting wireless clients to each other as well as connecting wireless and wired clients. In the latter, the access point can act as an Ethernet bridge by passing Layer 2 frames between the wired and wireless networks, or as a router, terminating WLAN and Ethernet Layer 2 frames and performing IP-level forwarding. The Layer 3 routing model is less popular and we will not consider it here.

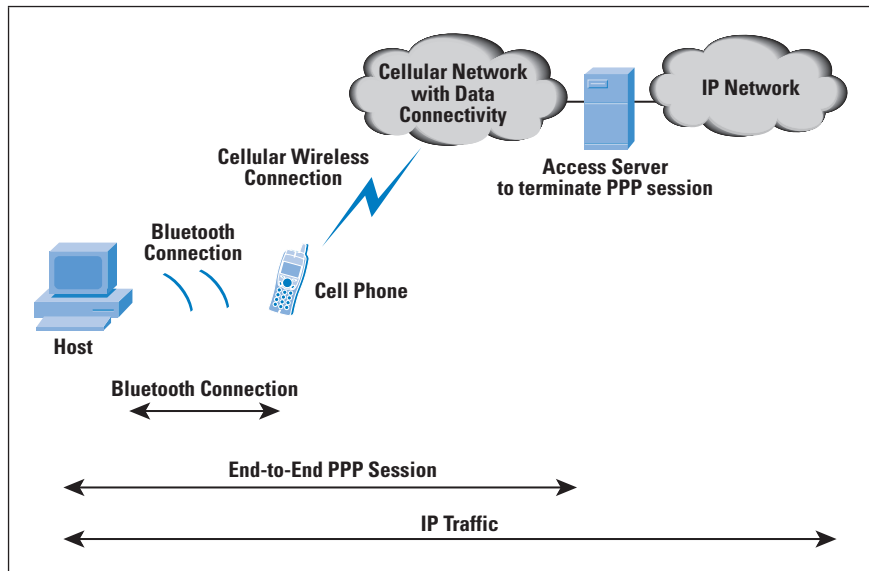
The access point typically terminates WLAN management and control frames. However, there is another model of a *thin access point* wherein these frames can be backhauled to a WLAN switch for processing. The access point connection to the wired network is typically an Ethernet link to a dedicated Ethernet switch port at 100-Mbps or Gigabit Ethernet speeds. With the advent of 802.11g and 802.11a WLANs, 10-Mbps links are not sufficient because these WLANs can operate at close to 27-Mbps throughput over the wireless network.

When considering 802.11n, we find that 100-Mbps backhaul links to the switch are insufficient for the 802.11n throughput of 150, or even 300 Mbps with channel bonding. Gigabit Ethernet links are often considered for connectivity between the 802.11n access point and the Ethernet switch. The next speed for Ethernet connectivity is 10 Gbps, which is well-established in the enterprise for data center and core Ethernet network applications. Work is ongoing in the IEEE for 40- and 100-Gbps Ethernet, so that should cover advances in wireless speeds for efficient backhaul to the wired network.

Bluetooth

Bluetooth started as a “wire-replacement” protocol for operation at short distances. A typical example is the connection of a phone to a PC, which, in turn, uses the phone as a modem (see Figure 2). The technology operates in the unlicensed 2.4-GHz ISM band. The standard uses FHSS technology. There are 79 hops in BT displaced by 1 MHz, starting at 2.402 GHz and ending at 2.480 GHz.

Figure 2: Typical Use of a Bluetooth enabled phone as a data modem for a PC



Bluetooth belongs to a category of *Short-Range Wireless* (SRW) technologies originally intended to replace the cables connecting portable and fixed electronic devices. It is typically used in mobile phones, cordless handsets, and hands-free headsets (though it is not limited to these applications). The specifications detail operation in three different power classes—for distances of 100 meters (long range), 10 meters (ordinary range), and 10 cm (short range).

Bluetooth operates in the unlicensed ISM band at 2.4 GHz (similar to 802.11 b/g wireless), but it is most efficient at short distances and in noisy frequency environments. It uses FHSS technology—that is, it avoids interference from other signals by hopping to a new frequency after transmitting and receiving a packet. Specifically, 79 hops are displaced by 1 MHz, starting at 2.402 GHz and finishing at 2.480 GHz.

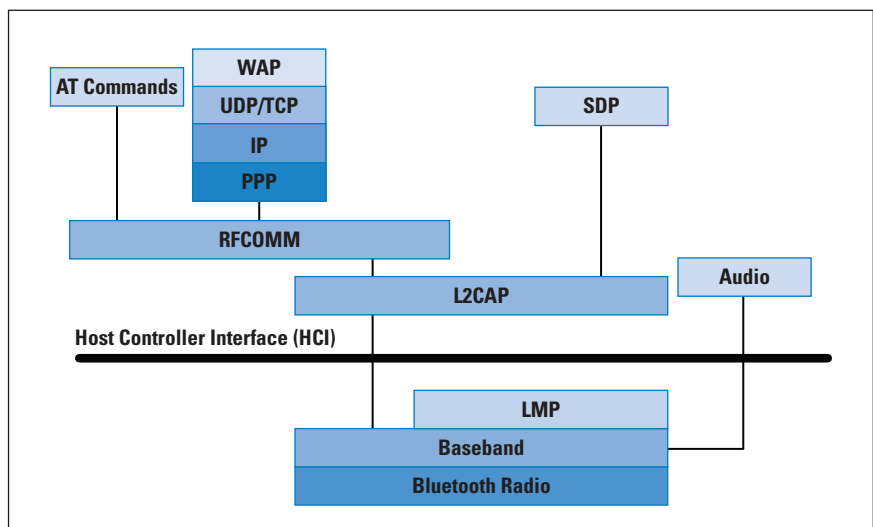
Bluetooth can operate in both point-to-point and logical point-to-multipoint modes. Devices using the same BT channel are part of a *piconet* that includes one master and one or more slaves. The master BT address determines the frequency hopping sequence of the slaves. The channel is also divided into time slots, each 625 microseconds in duration. The master starts its transmission in even-numbered time slots, whereas the slave starts its transmission in odd-numbered slots.

BT specifies two types of links, a *Synchronous Connection-Oriented* (SCO) link and an *Asynchronous Connectionless Link* (ACL). The SCO link is a symmetric point-to-point link between a master and a single slave in the piconet, whereas the ACL link is a point-to-multi-point link between the master and all the slaves participating in the piconet. Only a single ACL link can exist in the piconet, as compared to several individual SCO links.

Bluetooth Stack

Other than the radio and baseband components (the physical layer of Bluetooth that manages physical channels and links), the Bluetooth stack (see Figure 3) includes a *Link Manager Protocol* (LMP) used for link management between the endpoints, a *Logical Link Control and Adaptation Protocol* (L2CAP) for the data link, a *Radio Frequency Communication* (RFCOMM) protocol to provide emulation of serial ports over L2CAP, and a *Service Discovery Protocol* (SDP) for the dynamic discovery of services—because the set of services changes dynamically based on the RF proximity of the devices. In addition, the *Host Controller Interface* (HCI) provides a uniform command interface to the baseband controller and the link manager to have access to the hardware registers.

Figure 3: Key Elements of the Bluetooth Stack



LMP is required for authentication, encryption, switching of roles between master and slave, power control, and so on. L2CAP provides both connection-oriented and connectionless data services functions, including protocol multiplexing, segmentation and reassembly, and piconet-based group abstraction. As part of the multiplexing function, L2CAP uses the concept of channels, with a channel ID representing a logical channel endpoint on a BT device. L2CAP offers services to the higher layers for connection setup, disconnect, data reading and writing, pinging the endpoint, and so on.

RFCOMM, which provides emulation of serial ports on the BT link, can support up to 60 simultaneous connections between two BT devices. The most common emulation is of the RS-232 interface, which includes emulation of the various signals of this interface such as *Request To Send* (RTS), *Clear To Send* (CTS), *Data Terminal Ready* (DTR), and so on. RFCOMM is used with two types of BT devices—endpoints such as printers and computers and intermediate devices such as modems. In Figure 3, the IP stack over *Point-to-Point Protocol* (PPP) over RFCOMM emulates the mode of operation over a dialup or dedicated serial link. Because the various BT devices in a piconet may offer or require a different set of services, the *Service Discovery Protocol* (SDP) is used to determine the nature of the services available on the other nodes. SDP uses a request-response packet scheme for its operation.

Bluetooth Profiles

BT includes multiple profiles that correlate to the type of services that are available from BT nodes. For example, the BT headset profile is used between an audio source and a headset, both connecting wirelessly through BT—it involves a subset of the well-known AT commands used with modems. The audio source (typically a cell phone or cordless phone) implements the BT audio gateway profile for communicating with the device implementing the headset profile. Other profiles include a basic printing profile (often used for printing between a PC and a BT-enabled printer), dialup networking profile, fax profile, cordless telephony profile, *Human Interface Device* (HID) profile, and so on. The last profile is used for BT-enabled keyboards and mice—it is based on the HID protocol defined for USB.

The Bluetooth dialup networking profile is interesting from an IP perspective; as shown in Figures 2 and 3, it involves the IP stack running over RFCOMM to provide the appearance of a serial port running PPP, which is very similar to dialup networking over a basic telephone service line.

Bluetooth Frame Format and Speeds

The frame format in BT consists of a 72-bit field for the access code (including a 4-bit preamble, 64-bit synchronization field, and 4 bits of trailer), followed by a 54-bit header field that includes information about the frame type, flow control, acknowledgement indication, sequence number, and header error check. Following the header field is the actual payload, which can be up to 2745 bits. In all, the frame length can be a maximum of 2871 bits. Whereas synchronous BT traffic has periodic reserved slots, asynchronous traffic can be carried on the other slots.

BT ranges can vary from a low-power range of 1 meter (1 mW) for Class 3 devices, 10 meters (2.5 mW) for Class 2 devices, to 100 meters (100 mW) for Class 1 devices. BT Version 1.2 offers a data rate of 1 Mbps, and BT Version 2.0 with *Enhanced Data Rate* (EDR) supports a data rate of 3 Mbps. BT Version 1.1 was ratified as the IEEE Standard 802.15.1 in 2002.

Bluetooth versus Wi-Fi

A few years ago, some marketing literature tried to emphasize BT and Wi-Fi as competing technologies. Though both operate in the ISM spectrum, they were invented for different reasons. Whereas Wi-Fi was often seen as a “wireless Ethernet,” BT was initially seen purely as a cable- or wire-replacement technology. Uses such as dialup networking and wireless headsets fit right into this usage model. Recently, the discussion has focused more on coexistence instead of competition because they serve primarily different purposes. There are still some concerns related to their coexistence because they operate over the same 2.4-GHz ISM band.

To recapitulate, the Bluetooth physical layer uses FHSS with a 1-MHz-wide channel at 1600 hops/second (that is, 625 microseconds in every frequency channel). Bluetooth uses 79 different channels. Standard 802.11b/g uses DSSS with 20-MHz-wide channels—it can use any of the 11 20-MHz-wide channels across the allocated 83.5 MHz of the 2.4-GHz frequency band. Interference can occur either when the Wi-Fi receiver senses a BT signal at the same time that a Wi-Fi signal is being sent to it (this happens when the BT signal is within the 22-MHz-wide Wi-Fi channel) or when the BT receiver senses a Wi-Fi signal.

BT 1.2 has made some enhancements to enable coexistence, including *Adaptive Frequency Hopping* (AFH) and optimizations such as Extended SCO channels for voice transmission within BT. With AFH, a BT device can indicate to the other devices in its piconet about the noisy channels to avoid. Wi-Fi optimization includes techniques such as dynamic channel selection to skip those channels that BT transmitters are using. Access points skip these channels by determining which channels to operate over based on the signal strength of the interferers in the band. Adaptive fragmentation is another technique that is often used to aid optimization. Here, the level of fragmentation of the data packets is increased or reduced in the presence of interference. For example, in a noisy environment, the size of the fragment can be reduced to reduce the probability of interference.

Another way to implement coexistence is through intelligent transmit power control. If the two communicating (802.11 or Wi-Fi) devices are close to each other, they can reduce the transmit power, thus lowering the probability of interference with other transmitters.

WiBree to Low-Energy Bluetooth

WiBree is a technology first proposed by Nokia to enable low power communication over the 2.4-GHz band for button cell (or equivalent) battery-powered devices. A consequence of the low power requirement is the need for the wireless function to perform a very small set of operations when active and go back to the sleep or to standby mode when inactive.

The WiBree technology has been adapted by the *Bluetooth Special Interest Group* (SIG) as part of the lower-power BT initiative—also known as *Low Energy* (LE) BT technology. The LE standard is expected to be finalized sometime in 2009. When this standardization is completed, three types of BT devices will be available: traditional BT, LE BT, and a mixed or dual-mode BT. A mixed-mode device can operate in low power mode when communicating with other LE devices (for example, sensors) and traditional BT mode when communicating with BT devices, implying the presence of both a BT stack and an LE stack on the same device.

WiMAX

WiMAX stands for *Worldwide Interoperability for Microwave Access* and is defined under the IEEE 802.16 working group. Two standards exist for WiMAX—802.16d-2004 for fixed access, and 802.16e-2005 for mobile stations^[9]. The WiMAX forum certifies systems for compatibility under these two standards and also defines network architecture for implementing WiMAX-based networks.

WiMAX can be classified as a last-mile access technology similar to DSL, with a typical range of 3 to 10 kilometers and speeds of up to 5 Mbps per user with non-line of sight coverage. WiMAX access networks can operate over licensed or unlicensed spectra in various regions or countries—though licensed spectrum implementations are more common. WiMAX operation is defined over frequencies between 2 and 66 GHz, parts of which may be unlicensed spectrum deployments in some countries. The lower frequencies can operate over longer ranges and penetrate obstacles, so initial network roll-outs are in this part of the spectrum—with 2.3-, 2.5-, and 3.5-GHz frequency bands being common. Channel sizes vary from 3.5, 5, 7, and 10 MHz for 802.16d-2004 and 5, 8.75, and 10 MHz for 802.16e-2005. WiMAX networks are often used to backhaul data from Wi-Fi access points. In fact, they are often envisaged as replacements for the current implementation of metro Wi-Fi networks that use 802.11b/g for client access and 802.11a for backhaul to connect to the other parts of the network.

Technology

The 802.16d-2004 standard uses OFDM similar to 802.16a and 802.16g, whereas 802.16e-2005 uses a technology called *Scalable Orthogonal Frequency Division Multiplexed Access* (S-OFDMA). This technology is more suited to mobile systems because it uses subcarriers that enable the mobile nodes to concentrate the power on the subcarriers with the best propagation characteristics (because a mobile environment has more dynamic variables). Likewise, the 802.16e radio and signal processing is more complex.

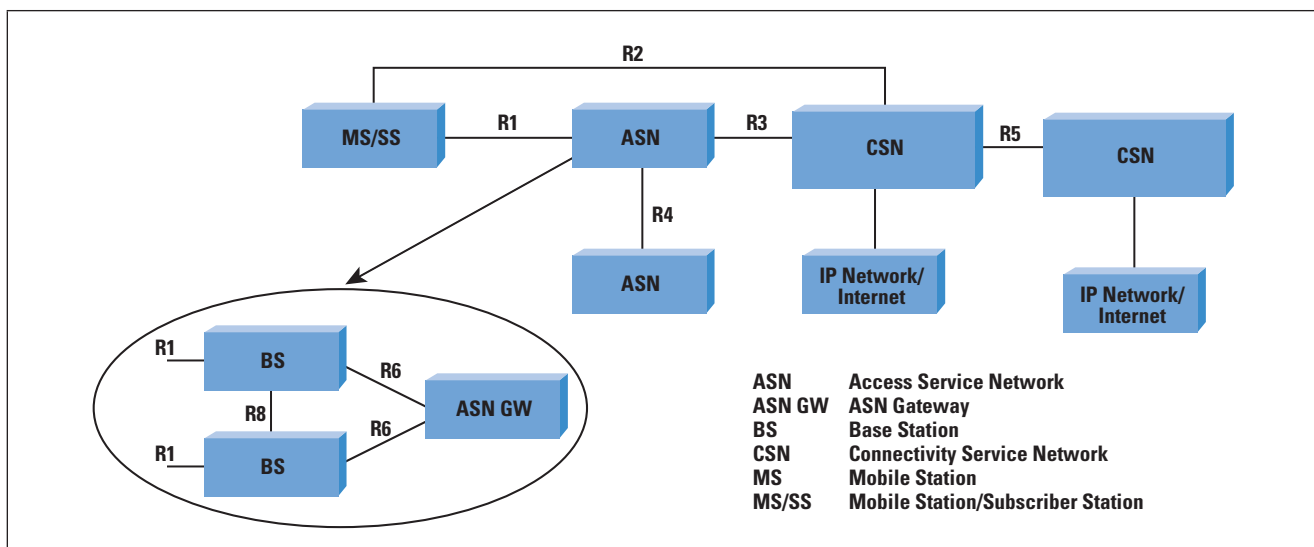
Unlike 802.11, which supports only *Time-Division Duplexing* (TDD)—where transmit and receive functions occur on the same channel but at different times), 802.16 offers TDD, *Frequency-Division Duplexing* (FDD) (transmit and receive on different frequencies, which could also be at different times). Another innovation in WiMAX is similar to the scheme in *Code Division Multiple Access* (CDMA)—subscriber stations are able to adjust their power based on the distance from the base station, unlike the case of client stations in an 802.11 network.

WiMAX base stations use a scheduling algorithm for medium access by the subscriber stations. This access is through an access slot that can be enlarged or contracted (to more or fewer slots) that is assigned to the subscriber stations. *Quality-of-Service* (QoS) parameters can be controlled through balance of the time-slot assignments among the base stations. The base-station scheduling types can be unsolicited grant service, real-time polling service, non-real time polling service, and best effort. Depending upon the time of traffic and service requested, one of these scheduling types can be used.

WiMAX Network Architecture

The WiMAX network architecture is specified through functional entities (see Figure 4), so you can combine more than one functional entity to reside on a network element. The *Mobile Station* (MS) connects the *Access Service Network* (ASN) through the R1 interface—which is based on 802.16d/e. The ASN is composed of one or more *base stations* (BSs) with one or more ASN gateways to connect to other ASNs and to the *Connectivity Service Network* (CSN). The CSN provides IP connectivity for WiMAX subscribers and performs functions such as *Authentication, Authorization, and Accounting* (AAA)^[10,11], ASN-CSN tunneling, inter-CSN tunneling for roaming stations, and so on. A critical tenet of the WiMAX Forum network architecture is that the CSN must be independent of the protocols related to the radio protocols of 802.16.

Figure 4: WiMAX Forum Network Architecture Functional Blocks and Interface Points



The R3 interface (reference point) is used for the control-plane protocols and bearer traffic between the ASN and CSN for authentication, policy enforcement, and mobility management. The base station connects to an ASN gateway to provide the MS with external network access. The R6 interface between the BS and ASN-GW could be open or closed based on the profile—in fact, you could have a co-located base station and *ASN gateway* (ASN-GW), depending upon the network implementation. The ASN gateway uses the R3 interface to communicate with the AAA services in the visited CSN (that is, the CSN “corresponding” to the ASN). The servers in the visited CSN can communicate with the home CSN (that is, the CSN corresponding to the “home” network of the MS). In the simplest case multiple ASNs (WiMAX networks) connect through ASN gateways to the public Internet (that is, there is only one *Network Service Provider* (NSP) and the visited and home CSNs are the same). Note that you could implement a WiMAX network with just one ASN and one CSN—in that case, the R3 interface would be completely internal and not exposed.

Three profiles are identified to map ASN functions into ASN-GW and BS functions. These profiles are considered an implementation guideline for how you would build the various devices implementing these functions. Profile A is a strict separation of the BS and ASN-GW functions, where the ASN-GW controls and manages radio resources that are located on the BS and also provides the handover and data-path functions. The R6 interface is exposed in this profile.

Profile B is a more integrated function, where the BS has more functions than in profile A; in fact, the BS might even integrate most of the ASN functions. The R6 interface is a closed interface in this profile. The third profile is profile C, which is similar to profile A except that the base stations incorporate more functions, including radio resource management and control as well as hand-offs.

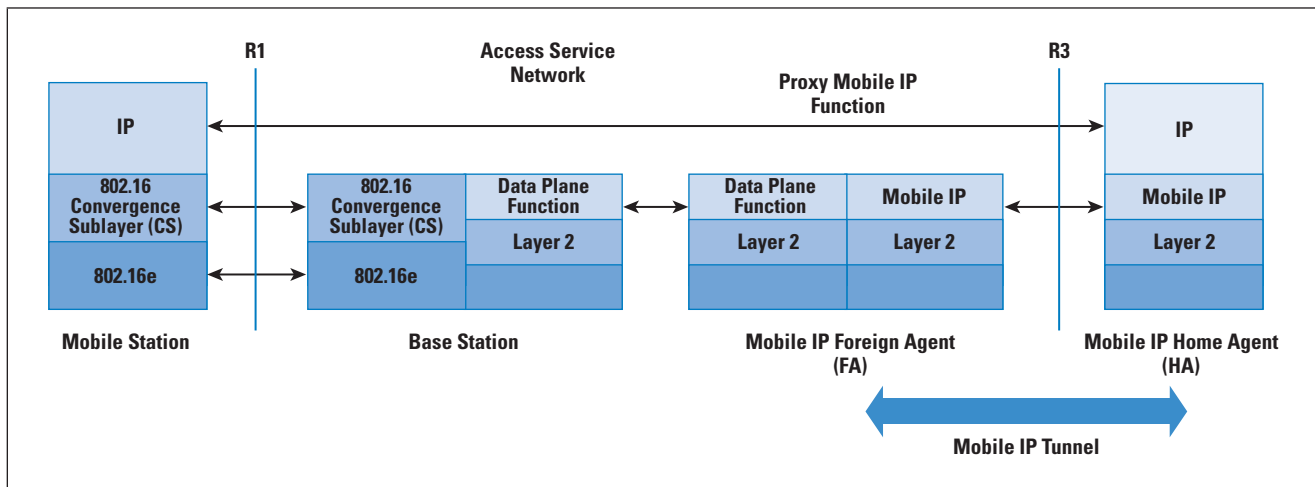
IP Connectivity and Data Transfer

The MS can be a fixed IP gateway (think of an 802.11 access point that provides connectivity to users in a coffee shop and connects to the IP network of the service provider through WiMAX) or a mobile end node (for example, a laptop with WiMAX connectivity). The IP address used by the gateway on the connection to the WiMAX network is known as the *Point of Attachment* (PoA) address. A third type of access is nomadic access, where the IP gateway can be moved from one location to another but connects to the network only after it has been relocated.

When the station is mobile, the WiMAX Forum specifies that the *Mobile IP* (MIP) architecture and protocols should be used. There are two types of Mobile IP possible: *Client Mobile IP* (CMIP) and *Proxy Mobile IP* (PMIP). The former involves changes to the MS protocol stack, but the latter does not.

The architecture can support both models. In the P-MIP scenario (see Figure 5), the ASN implements the Foreign Agent (see William Stallings' article in IPJ on Mobile IP^[8]), and terminates Mobile IP tunnels for the various mobile stations in the same ASN.

Figure 5: Data Transport and Proxy Mobile IP in WiMAX



In the figure, the MS has an address at the point of attachment that is used to forward packets from the MIP Foreign Agent inside the ASN. Because the ASN acts as a proxy of the attached MS, this implementation is known as a *Proxy MIP* implementation—also, there is no need for the MS to be aware of the MIP function being performed by the network.

Perspective on WiMAX versus Cellular Services

The WiMAX Forum has specified that the *Network Working Group* (NWG) architecture should be capable of supporting voice, multimedia services, and priority services such as emergency voice calls. It also supports interfacing with interworking and media gateways. Also, the service permits more than one voice session per subscriber, as well as simultaneous voice and data sessions. Support of IP Broadcast and Multicast services over WiMAX networks is also included. The architecture is also expected to support differentiated QoS levels at a per-MS or -user level (coarse grained) and at a per-service flow (fine-grained) level. It shall also support admission control and bandwidth management.

Initially, WiMAX was touted by some as a replacement for cellular services. An important consideration was using *Voice over IP* (VoIP) for voice calls—that is, where voice was another service over the data network. This model was in contrast to the existing cellular service where data was an adjunct to the basic service of TDM-based voice. More recently, WiMAX is being positioned as a data-connectivity option for remote locations, especially where it would be difficult to lay new copper or optical cable. Not surprisingly, these options are being pursued aggressively in developing countries.

Common Misperceptions About Wi-Fi, BT, and WiMAX Technologies

We have considered the key aspects of the three technologies—Wi-Fi, BT, and WiMAX—and their position in IP networks. In this section, we will outline and clarify some common perceptions and misperceptions about these technologies.

1. *BT and Wi-Fi are competing technologies*—Actually, they address a different set of requirements despite operating in the same 2.4-GHz space. BT is a “wire replacement” usually for short distances. Wi-Fi is typically used for data, voice, and video traffic over distances up to 300 meters.
2. *WiMAX is Wi-Fi on steroids*—To clarify, this statement is an oversimplification used often in the trade press. WiMAX operates in licensed spectra and uses a different network architecture as compared to Wi-Fi, which is in the unlicensed spectrum and uses a simple access point to wired Ethernet architecture. One overlapping function is for backhauling Wi-Fi traffic, which can be done by Wi-Fi (typically 802.11a) or WiMAX.
3. *Unlike BT, Wi-Fi cannot be used for voice*—This perception is not true because you can send multimedia traffic over Wi-Fi networks implementing 802.11e QoS functions that rely on the access point and stations implementing priority-based traffic transmission and scheduling.
4. *Wireless networks are not secure*—Although there is some validity to this argument because it is easier to eavesdrop on wireless networks, implementation of security schemes such as *Wi-Fi Protected Access* (WPA/WPA2) will help alleviate this problem.
5. *Wireless and radio technologies consume more power*—This statement is often true if the devices transmit continuously or have to increase their power because of the distance between the transmitter and receiver. Noisy channels contribute to this power use also. However, with careful engineering of the wireless implementation and techniques such as power save (in Wi-Fi) and short duty cycle transmissions, the power requirement can be lowered.

Summary

In this article, we have provided a flavor for IEEE 802.11 WLAN, Bluetooth, and WiMAX technologies and their implementation—specifically, how the nodes on these networks connect to an IP network. These technologies often serve complementary functions for end-to-end connectivity.

For Further Reading

- [1] IEEE 802.11 Standard, <http://standards.ieee.org/get-ieee802/download/802.11-2007.pdf>
- [2] Edgar Danielyan, “IEEE 802.11,” *The Internet Protocol Journal*, Volume 5, No. 1, March 2002.
- [3] T. Sridhar, “Wireless LAN Switches—Functions and Deployment,” *The Internet Protocol Journal*, Volume 9, No. 3, September 2006.
- [4] IEEE 802.16-2004 IEEE Standard for Local and metropolitan area networks Part 16: Air Interface for Fixed Broadband Wireless Access Systems.
- [5] IEEE 802.163-2005 IEEE Standard for Local and metropolitan area networks Part 16: Air Interface for Fixed and Mobile Broadband Wireless Access Systems Amendment for Physical and Medium Access Control Layers for Combined Fixed and Mobile Operation in Licensed Bands, IEEE, <http://standards.ieee.org/getieee802/802.16.html>
- [6] Bluetooth Special Interest Group Publications, <http://www.bluetooth.com/Bluetooth/Technology/>
- [7] http://www.wimaxforum.org/technology/documents/WiMAX_Forum_Network_Architecture_Stage_2-3_Rel_1v1.2.zip
- [8] William Stallings “Mobile IP,” *The Internet Protocol Journal*, Volume 4, No. 2, June 2001.
- [9] Jarno Pinola and Kostas Pentikousis “Mobile WiMAX,” *The Internet Protocol Journal*, Volume 11, No. 2, June 2008.
- [10] Convery, S., “Network Authentication, Authorization, and Accounting – Part One: Concepts, Elements, and Approaches,” *The Internet Protocol Journal*, Volume 10, No. 1, March 2007.
- [11] Convery, S., “Network Authentication, Authorization, and Accounting – Part Two: Protocols, Applications, and the Future of AAA,” *The Internet Protocol Journal*, Volume 10, No. 2, June 2007.

T. SRIDHAR is Vice President of Technology at Flextronics in San Jose, California. He received his BE in Electronics and Communications Engineering from the College of Engineering, Guindy, Anna University, Madras, India, and his Master of Science in Electrical and Computer Engineering from the University of Texas at Austin. He can be reached at T.Sridhar@flextronics.com

The End of Eternity

Part One: IPv4 Address Exhaustion and Consequences

by Niall Murphy, Google, and David Wilson, HEAnet

“Eternity is a very long time, especially towards the end,” said Woody Allen^[22,23], and he was mostly right. The eternity that the 32 bits of IPv4 address space promised is now almost at an end, and we are faced with the task of deciding what to do after the “end of eternity.”

The size of the problem of IPv4 exhaustion is, unfortunately, also proportional to its longevity^[1,2,3]. Although the next-generation (IPng) effort^[4] kick-started the development of IPv6 partially in response to concern about the IPv4 consumption rate, the industry as a whole largely ignored the problem after *Classless Inter-Domain Routing* (CIDR) and the *Regional Internet Registries* (RIR) system contained the depletion problem to a manageable horizon. More recently, after Geoff Huston’s^[5] work showing that the expected depletion time was sooner than many organizations had expected, the concern has received considerable attention in address-allocation policy circles.

In this article, we examine IPv4 exhaustion in more detail. We talk about what exactly exhaustion will mean and what we can do about it, and then present a vision for the postexhausted world. Those familiar with our RIPE-55 talk^[6] will find much that is familiar, but the arguments have been expanded for a more general audience. The authors, as in that talk, are speaking only for themselves, and not their organizations.

What Does Exhaustion Mean?

Trivially, the point of IPv4 exhaustion is the point at which the guaranteed-free-and-unused pool runs out and the current allocation mechanism comes to an end. Although the depletion of the free pool defines the technical point of exhaustion, it is not the depletion itself that is of primary importance. After all, if it were, we could simply declare a moratorium on allocations with immediate effect, to preserve the resource for some notional future requirements. Rather, it is the effect on the practices and procedures, within the RIRs and within the *Local Internet Registries* (LIRs), administrative and technical, that will practically define exhaustion. These practices, which have grown to fit around the current behavior of the addressing system, the free pool, and so on, will require urgent reform after exhaustion, as indeed will the RIR system in general.

Currently organizations use and require new addresses for essentially every IP-related additional deployment (for example, adding customers to a publicly numbered DSL service, adding extra *Secure Sockets Layer* (SSL)-enabled websites to a Web hosting service, and adding extra publicly reachable servers to almost any service).

It has been emphasized that this problem affects only the *growth* of organizations performing IP deployments^[7]. Although it is important to acknowledge the partial correctness of this statement, much about the postexhaustion state could undermine the stability of well-established advertisements and routes unless the transition is well-handled. It seems intuitively correct that those who received allocations before exhaustion will be unaffected by exhaustion turmoil^[8], but we regard this premise as optimistic, as you will see later.

Along those lines, one less well-examined consequence of exhaustion is the erosion of the consensus model of Internet governance. There is potential for wide divisions to open up at the local and regional level unless this consensus is carefully conserved. No clear successor to the current model as yet exists; the RIRs appeared to be heading toward a spectrum of positions on, for example, the allocation of the last portions of the IPv4 free pool^[9, 10] until quite recently^[24].

The erosion of this model of governance as a consequence of exhaustion has been neither widely examined nor expected in the Internet community. Partially, this situation arose because of the useful and well-executed role that the RIRs have historically filled in providing sensible and stable conditions for decision making; some proportion of the membership of the RIRs might well feel that IPv4 exhaustion is a problem like any other, which the RIRs themselves are in the perfect position to resolve. However, although the atmosphere of mutual cooperation fostered by the RIRs has produced many useful service-related outputs (for example, the *Test Traffic Measurement* service of *Réseaux IP Européens* [RIPE]^[25]), one of the major nonobvious benefits they have brought is to provide a centralized focus for discussion with governments and regulatory agencies. Not only is it more efficient and therefore less time-wasting to centralize through one representative organization, it has also created expectations that similar matters can be dealt with in the same coordinated way—a very valuable expectation, which has helped to increase the credibility of industry self-regulation. This credibility allowed, for example, the *Number Resource Organization* (NRO) to help forestall a proposal to allocate IPv6 according to geographical boundaries^[27, 28].

Indeed, without credible industry self-regulation, it is not at all clear that this community could have grown as fast as it did. Although it seems clear today that the RIRs are the correct place for this kind of activity to go on (witness RIPE’s “enhanced cooperation” task force^[26]), if they had not been around, government would either have had to deal with an organization with less of a pedigree or one with more inherent bias, or multiple organizations with competing biases, all of which could compel them to distrust the results of their liaisons. Unfortunately, in this respect the RIRs have been a victim of their own success.

Just as the consensus model in domains broke down when top- and second-level domains became monetized, so it is likely that the inherent win or loss for any given holder in any policy changes will undermine attempts to build consensus for address policy in a monetized IPv4 world. Absent this consensus, many of the RIR services that we rely upon will be undermined—not least the veracity of the WHOIS database and subsequent reliability of our routing filters, but also the RIR and *Internet Corporation for Assigned Names and Numbers* (ICANN) representations toward governments.

What Are the Problems with Exhaustion?

The biggest problem is the simplest one: existing organizations whose business model or operations are *solely* predicated on an ongoing flow of IPv4 addresses will fail. This premise would seem an extreme, even theoretical, characterization, but the size of this category in the real world is larger than you might think. Numerous organizations are also in trouble, perhaps less predicated upon IPv4 than the others, but that—for example—might have financial or operational difficulty in making the postexhaustion transition happen internally. They would also be placed at risk. Finally, there are those organizations that might rely on others to perform their transition correctly in order for them to continue effective operations: less directly at risk, but still probably affected.

Those who deal with the operation of the Internet on a daily basis are well-aware of the workarounds available that could save organizations from the doomsday scenario. It is unfortunate, then, that many of us have looked to the simplest cases in our immediate experience in order to form our opinions of the scale of the problem. It is indeed true that, in the short term, the client-side problem has largely been solved—provided that your customers or developers never have expectations in line with an end-to-end Internet. (It would seem that address-space pressure is likely to erode whatever end-to-end expectations still remain in today's Internet.)

However, the server-side problem (for example, SSL Website hosting, *IP Security* [IPsec] VPN endpoints, ...) remains unsolved. Workarounds exist^[11], but whether they will be ready and deployed in time remains an open question. There are, therefore, organizations operating at this moment that depend upon the continued availability of IPv4 addresses. Adequate workarounds have yet to be developed—never mind proven—for these businesses.

The situation becomes more complicated when we consider the candidate solutions. For example, such organizations as described previously cannot solve their problem by deploying IPv6 alone prior to the end of the transition, because they require universal reachability. Without universal reachability, support costs will rise, the quality of the user experience will decrease, and the credibility of Internet governance will be threatened. The only available evidence shows our position on the IPv6 transition curve being at the very beginning^[12].

Therefore it is difficult to emphasize this enough—new entrants providing Internet services *cannot expect to compete equally with existing operations*—because they have a very high barrier to entry formed not by the natural action and development of competitors, but by the resource scarcity of new addresses. Without new addresses, they cannot have an IPv4 *Default-Free Zone* (DFZ) routing-table entry; without a DFZ entry, they cannot be multihomed; without multihoming, they cannot offer sufficiently redundant Internet service; and without sufficiently redundant Internet service, they cannot meaningfully compete with existing operators.

A variety of poor-quality “fudges” are possible, of course: they could use the address space of their upstream operators (and run the risk of having that address space pulled or charged for), or they could outsource any address-requiring services to another organization (and be unable to control their service quality, as well as dependent upon their continuing operation), or they could host through some kind of public proxy network that redirects to their back-end servers through various hard-coded means (and create a fragile, difficult-to-operate network with higher running costs per unit customer than their competitors).

We will examine the other negative consequences of exhaustion in more detail later in the discussion; meanwhile, let us assume that the scenario described previously is undesirable enough for us to ask whether we can actually do anything to forestall it.

Can We Practically Defer Exhaustion?

What we would ideally like is some policy or algorithm that would give us more time—how much time is open to question—without producing its own set of ill effects. (We can certainly defer exhaustion by ceasing to allocate new IPv4 addresses tomorrow, but that solution is hardly practical.) Unfortunately, this problem is very difficult to resolve. Such direct precedents that appear clearly related to the current situation provide no useful guidance. Many resource-exhaustion problems have been faced before, but ultimately the solutions for those can be categorized into three kinds:

- *Make the resource renewable*: In this case, the resource is in danger of running out, but can be replenished by some means. Often this replenishment involves constraining production predicated on the resource to some smaller value, particularly when there is a natural rate of renewal—for example, fishing stocks. In the case of IPv4, it is fundamentally nonrenewable in that the resource is of a finite size. (As we discussed previously, current reclamation efforts^[13], although worthy of pursuit as a low-overhead task, cannot be a solution.)
- *Move to another resource*: This solution is already under way in the sense that we are engaged in the transition to IPv6. However, adoption of IPv6 will not happen fast enough to prevent the negative consequences of exhaustion.

- *Divide the resource more fairly:* This solution is useful primarily in the case where hoarding is taking place, causing resource problems for some significant proportion of a resource-using population. We are dividing the resource fairly as it is, and certainly since the emergence of the RIRs. For reasons discussed later, husbanding the resource more carefully is unlikely to actually be a solution.

We have faced other abstract exhaustion problems before as well: for example, phone-number depletion is somewhat similar to our current problem. However, phone-number depletion admits of a simpler solution—the creation of extra digits in the number space—because of the centralization of network knowledge in a comparatively small number of switches. For the Internet, where every deployed host would have to be informed about changes to the number space, such an approach is not operationally feasible. Furthermore, adding extra digits to the number code is not in fact simple, and telecommunications companies have experienced a wide range of problems with such approaches in the past, to say nothing of the loss of revenue and the failure of calls to connect because of customer confusion^[14, 15]. We see no historical situation that provides a clear precedent and a clear way forward.

SimLIR

Accordingly, to help answer the question posed in the preceding section, we wrote a tool, *SimLIR*, to explore exhaustion and post-exhaustion scenarios. Rather than being a tool influenced primarily by computations based on growth curves, a “top-down” approach, it is a modeling tool that examines how changes in behavior affect relative consumption rates. Roughly 6,000 lines of Python, the tool is due to be open-sourced at its Google Code page^[16] shortly after this article is available. The tool models the whole *Internet Assigned Numbers Authority* (IANA)—>RIR—>LIR hierarchy, and currently maps LIRs to countries; it uses the same publicly available data as Geoff’s work. We would appeal to the community to help improve the program, because more research is desperately needed in this area.

Running the tool under various scenarios has produced preliminary results indicating that we cannot meaningfully defer exhaustion, given our current growth rates. It can be used to compare the effect of policy adjustments on known historical and simulated behavior. For example, one simple policy adjustment that has been informally suggested is to decrease the initial allocation size for new LIRs. Modeling this allocation with the tool, we halve the size the LIRs receive at the time of initial membership. If we allow this scenario to run to completion, we have seen that it allows us to defer exhaustion by less than a week. Intuitively, we might expect this assumption to be realistic because startup activity, although important, is relatively small in terms of proportion of allocations. New LIRs numbered approximately 500 in 2006^[17], and any scheme that attempted to defer exhaustion based on such a small proportion of overall operations could not practically succeed.

The question then arises whether any other scheme based upon treating some partition of the request-space differently could have a significant positive effect. However, such a scheme necessarily assumes that some set of requests are oversized, and can in fact be shrunk with no ill effects. Even if they are oversized, identifying them without inducing either unworkable bureaucracy or a chilling effect on the operations of the organization would be a significant task, not lightly undertaken. Furthermore, it would be in the self-interest of the current RIR membership not to agree to such a change in policy. With any such scheme, there would be a non-zero chance of their own requests being deemed faulty in some respect, thus leading to significant risk to their own operations. All of this process would of course be happening in the approach to exhaustion, where it would be more critical than ever to receive enough numbering resources! We can assume, therefore, that no such scheme would ever make it past the policy-making apparatus of bottom-up-influenced RIRs. Ironically, the easiest changes to enact are changes governing allocations to startup organizations; the affected organizations are not in the room at the time of policy formation, because they are not members yet. But such changes are highly unlikely to have a positive effect.

Finally, partitioning schemes are similar to other schemes proposed to rework the *End Game* for IPv4 allocation^[18, 19, 20] or retain a certain proportion of the free pool for as-yet-unknown future needs, in that we put RIRs in the awkward situation of having to decide that some requests are more legitimate than others, at a time when these requests are likely to be particularly urgent. RIRs should not be in the business of deciding who gets to have new customers, and partitioning the request space invites the possibility of preferential treatment. We can be sure that any preferential treatment at this crucial time, accidental or otherwise, would attract lawsuits. Judicial involvement in the allocation process close to the time of exhaustion would benefit almost nobody.

It is important to note that these risks are mainly specific to partitioning the request space from the RIR to the LIR; in other words, imposing criteria at the time of request. Partitioning the remaining pool per RIR, that is, imposing criteria at the time of division, such as proposed by the $n = 1$ policy^[24], does not suffer from “favoritism.” Indeed, even if there were blatantly iniquitous division at the IANA-to-RIR level, although various checks and balances exist to ensure there is not, it would be unlikely to affect those with resources sufficient to possess an office in the region in question, or to open one up; it is patently clear that the requests will follow where the space is, and it is highly unlikely that any single RIR with a large amount of space left after others have been exhausted would be in any kind of position to pass a discriminatory policy.

We make these points to highlight that any scheme based upon LIR partitioning presents immense difficulties of principle. Even if these difficulties are worked out, they seem unlikely to meaningfully defer exhaustion: the current run rate for IPv4 address space will exhaust the space within a 5-year timeframe anyway, even if all practically possible measures are taken.

The Consequences of Scarcity

Suppose for the moment that at the time of exhaustion, Internet-connected organizations have to fend for themselves, with no particularly well-defined industry strategy in place. We would then expect to see a broad movement within the industry to conserve precious public IPv4 address space. One obvious way for an organization to obtain more usable IPv4 space is to move previously publicly-numbered resources behind *Network Address Translation* (NAT) gateways. Other, less-legitimate sources of new addresses will probably also be explored, and these actions, combined with the generally uncoordinated changes, may well trigger the following negative consequences:

- *Inability to measure clients, and difficulty of supporting them:* As we see more layers of NAT within networks, it becomes gradually more difficult to establish who is actually connecting to you, and what problems they are having. Cookies are a partial solution for only one important protocol. Measurement becoming harder means that support costs will rise.
- *Address-space hijacking:* As organizations become more desperate for space, it is entirely feasible that they will begin to cast around for space not explicitly unavailable in order to meet their business needs. How widespread this practice would be remains an open question, but effective barriers to this behavior are not currently available. We would expect a general deterioration in the quality of routing.
- *WHOIS database quality down:* Coupled with layers of NAT hiding more and more networks from direct sight, transfers of address space (legitimate or otherwise) will cause the WHOIS database to become gradually less and less accurate, leading to...
- *Distributed denial-of-service (DDoS) tracking trouble:* Problems tracking DDoS attacks and abuse origins of all kinds make law enforcement and network operators equally unhappy.
- *Connection quality down:* Connection quality, in terms of connections that complete successfully and have tolerable latency, will go down as a function of client growth behind gateways.
- *RIR billing model under pressure:* The RIRs will need to find a new way to pay their costs or go out of business—gradually, but inevitably. Of course the RIRs, like every other organization, must serve a need, but they currently provide a large number of ancillary services not directly related to IP allocation, and those services would also be under threat.

- *Consensus undermined:* This consequence is possibly the most dangerous of them all. If a chaotic state of affairs is allowed to continue for too long, our very ability to make decisions as a community will be undermined as organizations abandon the RIR model that has failed them. We will have squandered, in a way, the foundation of trust that allows such ethical codes as we have developed in Internet operations to persist. That foundation will not be easily recovered.

(Note that all of these are effects that are likely to emerge to varying degrees with the onset of scarcity, however it takes place; in other words, if the RIRs engage in a program of scarcity management by partitioning requests, it is highly likely that the scenario described previously will happen no matter what is left in the free pool.)

In any large shock such as we describe, there will be operational turmoil. Organizations will attempt to employ the technologies they need to dig themselves out of trouble, or bend the rules to the same end. There will be financial turmoil as the ability of each business to scale in the new regime is tested. Turmoil for existing businesses and new entrants will no doubt attract increased attention from governmental and quasigovernmental agencies of all kinds. Turnover in the routing table will increase as uncoordinated deaggregation of prefixes takes place. Unwelcome as all these consequences are, we will probably be far too preoccupied with our own individual problems to take care of the broader picture.

Postexhaustion Vision

Although we hope it is clear, given the previous discussion—that IPv4 addresses will still be required after exhaustion—our highest aspiration cannot be an Internet confined in perpetuity to IPv4 alone. If we are to continue in a manner resembling our current operations, we require continued address plenty, even by today's rather restricted standards. The End Game, therefore, is an IPv6 Internet, or at least enough of one to keep off address scarcity for a workable subset of the industry.

So, the problem can then be characterized as the transition toward this state of affairs—the gap between the end of the old allocation model and the emergence of an adequate replacement. Any solution will have to either make the gap shorter, by bringing users to the IPv6 Internet sooner, or make it less painful, by helping IPv4-dependent organizations survive. (Note that a solution that makes the gap less painful may well cause it to lengthen.)

With the problem stated this way, we can evaluate possible solutions in this context. A hurried, stimulated transition of popular services to IPv6 will quite likely shorten the gap, although a mass transition is also likely to be an unstable one and so rather painful.

A voluntary release of unused addresses may help reduce the pain, but is unlikely to service the run rate adequately, given its voluntary nature, and in any event will prolong dependence on IPv4, thus lengthening the gap. Tweaking policies to make remaining IPv4 addresses arbitrarily difficult to get merely introduces the effects of scarcity still sooner, helping neither goal.

That said, our initial examination of the problems of exhaustion indicate that there will be a group of people who will require IPv4 addresses after the exhaustion point, and it is also clear that there are those who have addresses, such as the lucky recipients of class A addresses in the early days, but no particular incentive to give them up. We do not actually want to recycle these prefixes indefinitely, however; that just sustains the current model. Optimally, we should provide whatever opportunity we can to those who require IPv4 addresses, to get them (and us) toward the End Game of an adequate global IPv6 deployment.

We do not require an unlimited IPv4 supply to accomplish this goal. We do, however, require liquidity: the ability to transfer, with incentives to transfer. Although it is very difficult for a centralized system (such as an RIR) to reclaim adequate space, the effort/reward ratio is much more favorable for an individual organization that knows its own network. So we must provide some stimulus for them to increase liquidity, while imposing some realistic restriction on demand. It must of course be scrupulously fair.

Stated in this way, a market-based trading exchange is not just one way of attempting to solve the problem—such an exchange, properly regulated, is arguably the most neutral and fairest way to manage the problem of scarcity.

In the next article we will explore how such a market system should work, discuss what new problems it is likely to create, and consider the potential effect on the routing table.

References

- [1] <ftp://ftp.ietf.org/ietf-online-proceedings/94dec/area.and.wg.reports/ipng/ale/ale-minutes-94dec.txt>
- [2] <http://tools.ietf.org/html/rfc2008>
- [3] Hain, Tony, “A Pragmatic Report on IPv4 Address Space Consumption,” *The Internet Protocol Journal*, Volume 8, No. 3, September 2005
- [4] <http://playground.sun.com/ipv6/doc/history.html>
- [5] <http://ipv4.potaroo.net>
- [6] <http://www.ripe.net/ripe/meetings/ripe-55/presentations/murphy-simlir.pdf>
- [7] http://www.isoc.org/educpillar/resources/ipv6_faq.shtml
- [8] <http://www.ietf.org/internet-drafts/draft-narten-ipv6-statement-00.txt>
- [9] <http://www.apnic.net/meetings/24/program/sigs/policy/presentations/el-nakhal-prop-051.pdf>
- [10] <http://www.ripe.net/ripe/policies/proposals/2007-06.html>
- [11] http://www.switch.ch/pki/meetings/2007-01/name-based_ssl_virtualhosts.pdf
- [12] For example, http://h.root-servers.org/128.63.2.53_2.html versus http://h.root-servers.org/h2_5.html
- [13] <http://www.ripe.net/ripe/meetings/ripe-55/presentations/vegoda-reclaiming-our.pdf>
- [14] A “smooth and convenient” dialing plan for India.
<http://www.mycoordinates.org/indias-phone-june-06>
- [15] http://en.wikipedia.org/wiki/UK_telephone_code_misconceptions
- [16] <http://code.google.com/p/simlir/>
- [17] <http://www.ripe.net/docs/ripe-407.html#membership>
- [18] <http://www.ripe.net/ripe/policies/proposals/2007-03.html>
- [19] <http://www.ripe.net/ripe/policies/proposals/2007-06.html>

- [20] <http://www.ripe.net/ripe/policies/proposals/2007-07.html>
- [21] <http://kuznets.fas.harvard.edu/~aroeth/alroth.html>
- [22] Woody Allen, "Side Effects," 1980.
- [23] Woody Allen through (most famously) Stephen Hawking, <http://www.cnn.com/2006/WORLD/asiapcf/07/04/talkasia.hawking.script/index.html>
- [24] <http://icann.org/en/announcements/proposal-ipv4-report-29nov07.htm>
- [25] <http://www.ripe.net/ttm/>
- [26] <http://www.ripe.net/ripe/tf/enhanced-cooperation/index.html>
- [27] <http://www.nro.net/documents/nro18.html>
- [28] <http://www.ripe.net/maillists/ncc-archives/im-support/2004/index.html>
- [29] Huston, G., "The Changing Foundation of the Internet: Confronting IPv4 Address Exhaustion," *The Internet Protocol Journal*, Volume 11, No. 3, September 2008.

NIALL MURPHY holds a B.Sc. in Computer Science and Mathematics from University College Dublin. While in university, he founded the UCD Internet Society, which provided Internet access to approximately 5,000 students. He went on to work for (and found) various organizations: the **.IE** domain registry, Club Internet (now Magnet Entertainment), Ireland On-Line, Enigma Consulting, Bitbuzz, and Amazon.com. He is currently in Site Reliability Engineering at Google. He is the coauthor of numerous articles, some RFCs, the O'Reilly book *IPv6 Network Administration*, and is a published poet and keen amateur landscape photographer. E-mail: **niallm@avernus.net**

DAVE WILSON holds a B.Sc. in Computer Science from University College Dublin, not coincidentally from around the same time as Niall. He has worked at HEAnet, the Irish National Research & Education Network, for more than 10 years, maintaining an involvement with RIPE and with the pan-European research network Géant. Dave is a member of the ICANN Address Supporting Organization Address Council; he helped to found the Irish IPv6 task force, which has the support of the national government there. E-mail: **dave.wilson@heanet.ie**

Remembering Jon: Looking Beyond the Decade

by Vint Cerf, Google

A decade has passed since Jon Postel left us.^[0] It seems timely to look back beyond that decade and to look forward beyond a decade hence. It seems ironic that a man who took special joy in natural surroundings, who hiked the Muir Trail and spent precious time in the high Sierras, was also deeply involved in that most artificial of enterprises, the Internet. As the *Internet Assigned Numbers Authority* (IANA)^[1] and the *Request for Comments* (RFC) editor, Jon could hardly have chosen more polar interests. Perhaps the business of the artificial world was precisely what stimulated his interest in the natural one.

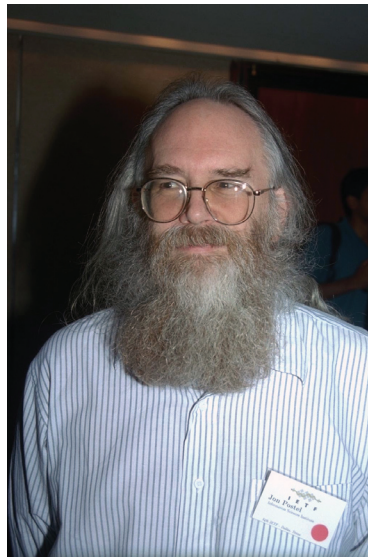


Photo: Peter Löthberg

As a graduate student at UCLA in the late 1960s, Jon was deeply involved in the ARPANET project, becoming the first custodian of the RFC note series inaugurated by Stephen D. Crocker. He also undertook to serve as the “Numbers Czar,” tracking domain names, Internet addresses, and all the parameters, numeric and otherwise, that were critical to the successful functioning of the burgeoning ARPANET and, later, Internet protocols. His career took him to the east and west coasts of the United States but ultimately led him to the University of Southern California’s *Information Sciences Institute* (ISI), where he joined his colleagues, Danny Cohen, Joyce K. Reynolds, Daniel Lynch, Paul Mockapetris, and Robert Braden, among many others, who were themselves to play important roles in the evolution of the Internet.

It was at ISI that Jon served longest and as the end of the 20th century approached, began to fashion an institutional home for the work he had so passionately and effectively carried out in support of the Internet. In consultation with many colleagues, but particularly with Joseph Sims of the Jones Day law firm and Ira Magaziner, then at the Clinton administration White House, Jon worked to design an institution to assume the IANA responsibilities. Although the path to its creation was rocky, the *Internet Corporation for Assigned Names and Numbers* (ICANN)^[2] was officially created in early October 1998, just two weeks before Jon’s untimely death on October 16.

In 1998 an estimated 30 million computers and 70 million users were on the Internet. In the ensuing decade, the user population has grown to almost 1.5 billion and the number of servers on the Internet now exceeds 500 million (not counting episodically connected laptops, *personal digital assistants* [PDAs], and other such devices). As this decade comes to a close, the *Domain Name System* (DNS) is undergoing a major change to accommodate the use of non-Latin character sets in recognition that the world's languages are not exclusively expressible in one script^[7]. A tidal wave of newly Internet-enabled devices as well as the increasing penetration of Internet access in the world's population is consuming what remains of the current IPv4 address space, accelerating the need to adopt the much larger IPv6 address space in parallel with the older one. More than three billion mobile devices are in use, roughly 15 percent of which are already Internet-enabled.

Jon would take considerable satisfaction knowing that the institution he worked hard to create has survived and contributed materially to the stability of the Internet. Not only has ICANN managed to meet the serious demands of Internet growth and importance in all aspects of society, but it has become a worked example of a new kind of international body that embraces and perhaps even defines a multi-stakeholder model of policy making. Governments, civil society, the private sector, and the technical community are accommodated in the ICANN policy development process. By no means a perfect and frictionless process, it nonetheless has managed to take decisions and adapt to the changing demands and new business developments rooted in the spread of the Internet around the globe.

Always a strong believer in the open and bottom-up style of the Internet, Jon would also be pleased to see that the management of the Internet address space has become regionalized and that five *Regional Internet Registries* (RIRs)^[3] now cooperate on global policy, serving and adapting to regional needs as they evolve. He would be equally relieved to find that the loose collaboration of DNS root zone operators has withstood the test of time and the demands of a much larger Internet, showing that their commitment has served the Internet community well. Jon put this strong belief into practice as he founded and served as ex-officio trustee of the *American Registry for Internet Numbers* (ARIN)^[4].

As the first individual member of the Internet Society he helped to found in 1992, Jon would certainly be pleased that it has become a primary contributor to the support of the Internet protocol standards process, as intended. The Internet Architecture Board and Internet Engineering and Research Task Forces, as well as the RFC editing functions, all receive substantial support from the Internet Society.

He might be surprised and pleased to discover that much of this support is derived from the Internet Society's creation of the *Public Interest Registry* (PIR)^[5, 6] to operate the **.org** top-level domain registry. The Internet Society's scope has increased significantly as a consequence of this stable support, and it contributes to global education and training about the Internet as well as to the broad policy developments needed for effective use of this new communication infrastructure.

As a computer scientist and naturalist, Jon would also be fascinated and excited by the development of an interplanetary extension of the Internet to support manned and robotic exploration of the Solar System. In October 2008, the Jet Propulsion Laboratory began testing of an interplanetary protocol using the Deep Impact spacecraft now in eccentric orbit around the sun. This project began almost exactly 10 years ago and is reaching a major milestone as the first decade of the 21st century comes to an end.

It is probable that Jon would not agree with all the various choices and decisions that have been made regarding the Internet in the last 10 years, and it is worth remembering his philosophical view: "Be conservative in what you send and liberal in what you receive."

Of course he meant this idea in the context of detailed protocols, but it also serves as a reminder that in a multi-stakeholder world, accommodation and understanding can go a long way toward reaching consensus or, failing that, at least toleration of choices that might not be at the top of everyone's list.

No one, not even someone of Jon's vision, can predict where the Internet will be decades hence. It is certain, however, that it will evolve and that this evolution will come, in large measure, from its users. Virtually all the most interesting new applications of the Internet have come not from the providers of various Internet-based services, but from ordinary users with extraordinary ideas and the skills to experiment. That they are able to experiment is a consequence of the largely open and nondiscriminatory access to the Internet that has prevailed over the past decade. Maintaining this spirit of open access is the key to further development, and it seems a reasonable speculation that if Jon were still with us, he would be in the forefront of the Internet community in vocal and articulate support of that view.

A 10-year toast seems in order. Here's to Jonathan B. Postel, a man who went about his work diligently and humbly, who served all who wished to partake of the Internet and to contribute to it, and who did so asking nothing in return but the satisfaction of a job well done and a world open to new ideas.

References

- [0] Vint Cerf, “I Remember IANA,” *The Internet Protocol Journal*, Volume 1, No. 3, December 1998. Also published as RFC 2468, October 1998.
- [1] <http://www.iana.org>
- [2] Vint Cerf, “Looking Toward the Future,” *The Internet Protocol Journal*, Volume 10, No. 4, December 2007.
- [3] Daniel Karrenberg, Gerard Ross, Paul Wilson, and Leslie Nobile, “Development of the Regional Internet Registry System,” *The Internet Protocol Journal*, Volume 4, No. 4, December 2001.
- [4] <http://www.arin.net>
- [5] <http://www.pir.org>
- [6] <http://www.isoc.org>
- [7] Huston, G., “Internationalizing the Domain Name System,” *The Internet Protocol Journal*, Volume 11, No. 1, March 2008.

VINTON G. CERF is vice president and chief Internet evangelist for Google. Cerf served as a senior vice president of MCI from 1994 through 2005. Widely known as one of the “Fathers of the Internet,” Cerf is the co-designer of the TCP/IP protocols and the architecture of the Internet. He received the U.S. National Medal of Technology in 1997 and the 2004 ACM Alan M. Turing award. In November 2005, he was awarded the Presidential Medal of Freedom. Cerf served as chairman of the board of the *Internet Corporation for Assigned Names and Numbers* (ICANN) from 2000 through 2007 and was founding president of the Internet Society. He is a Fellow of the IEEE, ACM, the American Association for the Advancement of Science, the American Academy of Arts and Sciences, the International Engineering Consortium, the Computer History Museum, and the National Academy of Engineering. He is an honorary Freeman of the City of London. Cerf holds a Bachelor of Science degree in Mathematics from Stanford University and Master of Science and Ph.D. degrees in Computer Science from UCLA. E-mail: vint@google.com

Letters to the Editor

IPv4 Address Exhaustion

I read with interest your article in *The Internet Protocol Journal* (Volume 11, No. 3, September 2008) regarding the IPv4 address exhaustion problem. It occurs to me that two approaches for encouraging the public and *Internet Service Provider* (ISP) community to migrate to IPv6 are being dismissed somewhat, but used creatively together might offer some hope for pushing us in that direction: government regulation and changing the fact that there isn't a public interest in IPv6.

What if government regulation forced a new or currently existing common service to use IPv6? One obvious possibility is video content. Since the broadcast industry is already regulated by the FCC, further regulation providing for governance of this type of application isn't too much of a stretch. Consumer demand is likely to increase in this area as broadband continues to be widely deployed, and if the public were required to run in dual-stack mode to access it, the likelihood of adoption would be much greater. It would also incent the ISPs to provide connectivity to the IPv6 address space, possibly even with a revenue-generating model behind it.

I reluctantly bring up the pornography industry as another type of content that could be relegated to the IPv6 address space. It is my understanding that this type of traffic as a percentage of the total is quite large. Based on this assumption, it would have the same effect of forcing the large portions of the public and ISPs to provide connectivity to the IPv6 address space. Again, I mention this industry reluctantly, but from a political perspective regulation of this industry and its content is likely to be an easier proposal for the public to support since you could use the "value" of disconnected portions of the Internet to best advantage.

I realize that the global nature of the Internet makes regulation and the subsequent enforcement extremely difficult. But, I also assume that even if our enforcement were controlled only at the perimeter of the U.S. traffic it would have a strong effect on the behavior of the public and ISPs.

Best regards,

—John Newell, INX Inc.
jcnnewell@gmail.com

The author responds:

Thanks for your response. It is true to say that various efforts have been undertaken across many years to find a "killer-app" for IPv6, if I may be permitted to use that overabused and by now very tired term. To date these efforts have not been successful. That's not because of any lack of trying.

There have been some really quite innovative ideas for IPv6 over the years, and so far most of them have been retrofitted into IPv4 one way or another. From one perspective this retrofit is entirely logical, given that good ideas tend to thrive in locations where audiences are receptive, and today's IPv4 Internet is still a very fertile place for good ideas to flourish.

The other part of the problem is that service providers tend to create innovative services with existing markets in minds, so these days the novel applications and services that appear to gain the attention of significant parts of the user base tend to operate in the IPv4 network, and by necessity such applications and services account for *Network Address Translation* (NAT) devices and various forms of filters and firewalls.

These observations indicate that a certain reinforcing cycle exists that cements the existing role of the IPv4 Internet, and tends to work against the widespread deployment of innovative services that are feasible only in the IPv6 environment.

So if the adoption of IPv6 is a carrot or stick affair, our efforts to find some tempting carrots have, so far, not been overly successful. We've been unable to identify particular goods or services for which there is a compelling case of consumer demand coupled with a set of technology constraints that imply that the service is feasible only across a deployed IPv6 infrastructure with IPv6 endpoints. So if the field we are working in is bereft of carrots, are there any available sticks that we can use instead? In this case there is the same old stick that originally motivated IPv6 in the first place: We are running out of IPv4 addresses. If we believe that there is more to do in the Internet, more people to connect, more devices to add, more conversations to have, more services to deploy, more ideas to realize, and more objectives to achieve, then IPv4 cannot in and of itself sustain that vision for the Internet. The threat here is that the growth of the IPv4 Internet may well cease when the supply of further IPv4 addresses is exhausted.

Is this threat of network stagnation going to be enough to propel us into an IPv6 Internet? Will it be an adequate motivator to encourage the necessary investment in network infrastructure and in the provision of goods and services that first operate in a transitional dual-stack environment, and ultimately in an IPv6 world? I hope that the answers are "yes," as do many others I'm sure.

But I'm also worried that it may not be enough and that we may spin off into an entirely different trajectory that ultimately dismantles most of the attributes of today's Internet. I worry that instead of an open network that fosters innovation and creativity we might end up with "vertical integration" and "transparent convergence" and a network that actively resists new services and applications.

So for me, and I hope many others, IPv6 needs no new “killer-app.” IPv6 does not need television or pornography to succeed. IPv6 is an imperative for the Internet simply because the alternatives to IPv6 appear to offer us a leap backward in technology and a leap backward in the elastic ways we’ve been able to use networks—and in the process we are going to destroy the Internet as we know it!

Regards,

—Geoff Huston, APNIC
gih@apnic.net

Dear Ole,

In his latest IPJ article (Volume 11, No. 3), Geoff Huston highlights the significance of NAT as a mechanism enabling service providers to externalize the costs and risks arising from IPv4 address scarcity. While acknowledging the increased burden and uncertainty borne by end users and NAT-traversing applications, Geoff speculates that the success of this mechanism is likely to inspire the deployment of yet another level of (“carrier grade”) address translation, to further prolong if not absolutely preclude the incorporation of IPv6 by incumbent service providers. While entirely plausible, such a move would create the same kind of “double blind” conditions for Internet service delivery that prevailed in financial markets when debt securitization was coupled with the externalization of asset depreciation risks in the form of *Credit Default Swaps*. In such cases, the second layer of indirection tends to make it all too easy to maintain self-serving assumptions (and/or plausible deniability) about the true nature and purpose of the first layer, and thus to fuel the perpetuation of unsustainable industry practices unto the point of industry collapse. Given the now inescapable lessons of the recent financial sector collapse, it would be nice if we didn’t have to learn this particular one again the hard way.

—Tom Vest
tvest@eyeconomics.com

On Paper

I just received the September issue (Volume 11, No. 3) of IPJ and wanted to make a quick comment about the paper change. Upon reading the section on the change I quickly dug up the previous copy of IPJ and compared the two. I personally like the new paper much better. The main reason I like it is because it is much easier on the eyes, I think mostly because it no longer has a glare from overhead lighting reflecting like the old paper type did. It’s a welcomed change from my take.

—David Swafford,
Network Engineer for CareSource, Dayton, OH, US
david@davidswafford.com

Book Reviews

A Dictionary and a Handbook

Hundreds of telecom books are published each year, but it is unusual to find a really good one. There must have been a blue moon (I'll have to check my almanac) this month, for I found two new and quite remarkable books by the same author, Ray Horak. One is a dictionary and the other an encyclopedic work, both covering the full range of voice, data, fax, video, and multimedia technologies and applications that comprise contemporary telecommunications. Further, they do so in such a plain-English, commonsense manner that you don't need to be a serious telecom student or professional to benefit from them—any layperson with a serious need to know will find them to be of great value. Finally (and this is rare in a technical book), both are actually relatively easy and certainly interesting reads, with liberal doses of fascinating historical context. In fact, they are even strong on entertainment value, with humorous observations and quotations sprinkled throughout. Horak has written each book in a different style for a different purpose, so they are best acquired together—as a set.

Webster's New World Telecom Dictionary

Webster's New World Telecom Dictionary, by Ray Horak, ISBN-10: 047177457X, ISBN-13 978-0471774570, Wiley Publishing Inc., 2007.

In order to communicate effectively in a contemporary telecom conversation, one must speak a special language rife with technical terminology, much of which is in the form of abbreviations, acronyms, contractions, initialisms and portmanteaux. To add to the confusion, many terms have multiple very precise—and occasionally imprecise—meanings, depending on the context. Writing a telecom dictionary must be a formidable task, one which only either the very brave or very foolhardy would even attempt. I'm not sure into which category Ray Horak falls, but his *Webster's New World Telecom Dictionary* is an excellent piece of work.

Organization

Dictionaries are in alphabetical order, of course, with chapters thrown in for symbols and numbers. Because the introduction of symbols requires special treatment, within each of the 28 chapters Horak organizes the approximately 4,600 definitions in ASCII order, perhaps as an accommodation for the binarians among us. The book includes an appendix of standards organizations and special interest groups, which can be useful if you need more information on a subject or need to know exactly to whom to complain about a *standard* or *specification*, both of which terms are defined clearly in the dictionary, of course.

Comparisons: Comprehensive and Correct

In my opinion, the best telecom dictionary ever written, aside from *Webster's*, is the *Communications Standard Dictionary*, by Martik H. Weik. That book unfortunately is out of print, with the final 3rd edition dated 1996. At 1095 pages, it is a bit overwritten and way too technical for most purposes, reading much like an IEEE dictionary. At this point, it certainly is out-of-date.

A handful of other telecom dictionaries and encyclopedias are currently in print, by far the most popular of which is *Newton's Telecom Dictionary*. Because *Newton's* dominates the market and has done so for many years, any telecom dictionary or encyclopedia is inevitably compared to that work. *Webster's New World Telecom Dictionary* is no exception, particularly because Ray Horak was the contributing editor to *Newton's* from the 12th through the 22nd editions.

Although *Webster's* defines only 4,600 terms in comparison to *Newton's* highly dubious claim of some 24,500 terms, *Webster's* definitions are much better researched, much more precise, and much more efficiently worded (that is, there is much less “fluff”). Even if *Webster's* almost certainly will gain in bulk as future editions expand the coverage of the telecom domain, it contains all of the essential telecom and IT terms, and defines them clearly and concisely. *Webster's* includes many humorous definitions but, unlike *Newton's*, they are all relevant and meaningful. For example, Horak lists three types of standards—*de jure*, *de facto*, and *du jour*. According to him, a *du jour* standard is defined as follows:

“From French, meaning *of the day*. The popular standard of the day. One day 10 years ago, ATM was really hot and a lot of people made a lot of money talking about ATM and selling products based on ATM. It seemed like only the next day that IP was really cool. (I made this one up.)”

Other humorous definitions include analogue, endianness, Hellenologophobia, hoot 'n' holler, OCD, PC, and WMBTOTCITB-WTNTALI. All of these, and more, serve to lighten the load, so to speak, but none of this humor detracts from what is a serious book on a serious subject. *Newton's*, on the other hand, is so full of personal observations and anecdotes, irrelevant humor (?), and inaccurate definitions as to make you wonder why bother to make the comparison at all. Horak states that he wrote *Webster's* partly to atone for his sins in contributing to *Newton's*, but mostly to put an authoritative reference book in his own hands, and those of others involved in litigation support. He apparently does a fair amount of work as an expert witness in intellectual property (the other IP) cases and on innumerable occasions has been asked to define and opine on terms such as link, circuit, channel, call, connection, switch, router, and PSTN. Now he can testify in court with one hand on the Good Book and the other on *Webster's*.

Recommended

Webster's New World Telecom Dictionary is an excellent piece of work. Ray Horak and his technical editor, Bill Flanagan, have collaborated to create a well-written, authoritative work that clearly sets a new standard for telecom dictionaries. I highly recommend it to anyone serious about telecom.

Telecommunications and Data Communications Handbook

Telecommunications and Data Communications Handbook, by Ray Horak, ISBN-10: 0470041412, ISBN-13: 978-0470041413, John Wiley & Sons, 2007.

Unless you have really big hands, you may wonder how it is that a tome of 791 pages that weighs more than 3 pounds could possibly be called a handbook. Well, the term “handbook” actually is fairly imprecise, but Ray Horak's *Telecommunications and Data Communications Handbook* certainly is not. Actually, it is about as compact as it can be, given its encyclopedic nature, and it is very precise, indeed. The book covers the entire telecom landscape, from wireline to wireless, from copper to radio and fiber, from electrical to optical, and from the customer premises to the cloud. It discusses voice, data, fax, video and multimedia technologies, systems, and applications in great detail, and in the LAN, MAN, and WAN domains. The handbook explores every relevant technology, standard, and application in the telecom and datacom space.

Horak is a well-known telecom consultant, author, writer, columnist, and lecturer. The *Telecommunications and Data Communications Handbook* is based on his best-selling *Communications Systems and Networks* (1997, 2000, 2002), but is considerably more technical and broader in scope. It is exceptionally well-written in Horak's plain-English, commonsense style, making it just as helpful to the neophyte and layperson as to the serious student or seasoned IT professional. Horak makes liberal use of well-constructed graphics to illustrate system and network architectures, topologies, and applications.

Organization

The Handbook begins with an excellent table of contents (20 pages) and ends with an excellent index (29 pages), both of which are crucial to a good book. After all, it doesn't make any difference how good the information is if you can't find it. The book is logically organized into 15 chapters and 2 appendixes.

Chapter 1 is devoted to fundamental concepts and definitions, thereby building a firm foundation of concepts and terminology upon which subsequent chapters build. Terms such as two-wire, four-wire, circuit, link, channel, switch, and router are clearly defined, compared, and contrasted. Chapter 2 explores the full range of transmission systems, including twisted pair (UTP, STP, and ScTP), coaxial, microwave, satellite, *Free Space Optics* (FSO), fiber-optics, *powerline carrier* (PLC), and hybrid systems.

Chapter 3 examines voice communications systems: KTS, PBX, Centrex, and ACD. Chapter 4 discusses messaging systems in detail, including facsimile (fax), voice processing, and e-mail and instant messaging, concluding with a detailed discussion of unified messaging and unified communications. Chapter 5 is dedicated to the *Public Switched Telephone Network* (PSTN) and addresses *Numbering Plan Administration* (NPA), regulatory domains, rates and tariffs, signaling and control systems, and network services. Chapter 6 returns to fundamentals, this time in the data communications domain, with detailed explanations of *Data Communications Equipment* (DCE) such as modems, codecs, CSUs, and DCUs, and then moves on to protocol basics, code sets, data formats, error control, compression techniques, network architectures, and security mechanisms.

Chapter 7 deals with conventional digital and data networks such as DDS, Switched 56, VPNs, T/E-carrier, X.25, and ISDN. Chapter 8 treats *Local-Area Networks* (LANs) and *Storage Area Networks* (SANs) exhaustively, including transmission media, topologies, broadband vs. baseband, equipment, operating systems, and standards. This chapter covers 802.3, 802.11, HiperLAN, Bluetooth, IEEE 1394, Fibre Channel, and iSCSI in considerable detail. Chapter 9 is devoted to broadband network infrastructure, including both access technologies (for example, xDSL, CATV, WLL, PON, and BPL) and transport technologies (for example, SONET/SDH and RPR). Chapter 10 offers an exhaustive study of broadband network services, including Frame Relay, ATM, Metropolitan Ethernet, B-ISDN, and AINs.

Chapter 11 discusses wireless, with an emphasis on mobility, covering both broad concepts and technical specifics of *Specialized Mobile Radio* (SMR), paging, cellular (1G, 2G, 2.5G, 3G, and beyond), packet data radio networks, and mobile satellite networks (GEOs, MEOs, and LEOs). Chapter 12 thoroughly treats video and multimedia networking, including a detailed discussion of video and multimedia standards (for example JPEG, MPEG, and H.320), *Session Initiation Protocol* (SIP), and IPTV. Chapter 13 exhaustively and insightfully explores the Internet and *World Wide Web* (WWW), including a thorough discussion of the IP protocol suite. Chapter 14 briefly examines convergence, and Chapter 15 examines telecom regulation, with a focus on the United States.

Appendix A is something of a decoder for abbreviations, acronyms, contractions, initialisms, and symbols. Appendix B gives a complete listing of relevant standards organizations and special interest groups, including full contact information, in case you need more information or want to offer comments on a particular subject.

Comparisons

It is hard to make a valid direct comparison to this book. *The Irwin Handbook of Telecommunications*, by James Harry Green, is good, but less complete, less technical, and drier, if such a combination is possible. The most recently published 5th edition also is apparently out of print. *The Voice & Data Communications Handbook*, by Regis “Bud” Bates, is written at a lower level; and, the *Essential Guide to Telecommunications*, by Annabel Dodd, at a much lower level. These latter two books are breezy reads and appeal more to a mass market than to a serious student or professional.

The *Telecommunications and Data Communications Handbook* compares more correctly to some of the more seminal works of Gilbert Held or James Martin, but covers a much wider range of subject matter and is a much easier and more pleasant read.

Recommended

The *Telecommunications and Data Communications Handbook* is written for the academic and professional community, but is just as relevant to anyone who needs to understand telecommunications system and network technologies and their meaningful applications. It is an exceptional work that should be on every IT professional’s bookshelf...when not in his or her hands.

—John R. Vacca,
jvacca@frognet.net

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. Contact us at ipj@cisco.com for more information.

Itojun Service Award Launched

A new award, providing recognition and support for those progressing IPv6 development on the Internet, was announced in November. The *Itojun Service Award* honors the memory of Dr. Jun-ichiro “Itojun” Hagino, who passed away in 2007, aged just 37^[1]. The award, established by the friends of Itojun and administered by the *Internet Society* (ISOC), recognizes and commemorates the extraordinary dedication exercised by Itojun over the course of IPv6 development. Itojun worked as a Senior Researcher at the *Internet Initiative Japan* (IIJ), was a member of the board of the *Widely Integrated Distributed Environment* (WIDE) Project, and from 1998 to 2006 served on the groundbreaking KAME project in Japan as the “IPv6 Samurai.” He was also a member of the *Internet Architecture Board* (IAB) from 2003 to 2005.

At the time of his passing, Russ Housley, *Internet Engineering Task Force* (IETF) Chair, and Olaf Kolkman, IAB Chair, issued a joint statement, praising Itojun’s service to IPv6 developments, saying that he had “inspired many and will be missed.”

The Itojun Service Award will run for 10 years, presented annually to an individual who has made outstanding contributions in service to the IPv6 community. The award includes a presentation crystal, a US\$3,000 honorarium, and a travel grant. The Award will honor an individual who has provided sustained and substantial technical contributions, service to the community, and leadership. With respect to leadership, the selection committee will place particular emphasis on candidates who have supported and enabled others in addition to their own specific actions.

The selection committee members for the Itojun Service Award are: Jun Murai, Hiroshi Esaki, Ole Jacobsen, Bob Hinden, Randy Bush, Bill Manning, Tatuya Jinmei, Kazu Yamamoto, and Kenjiro Cho.

Memorial donations to the Itojun Service Award Fund are welcomed and the Internet Society has established an account for donations. Details of the fund, as well as more information about Jun-ichiro “Itojun” Hagino and the Itojun Service Award are available on the ISOC Web site: <http://www.isoc.org/awards/itojun/>

The WIDE Project has also established a Japanese bank account to collect donations in Japanese Yen, the details of which are available here: <http://www.wide.ad.jp/itojun-award>

[1] Hinden, Bob, “Remembering Itojun: The IPv6 Samurai,” *The Internet Protocol Journal*, Volume 10, No. 4, December 2007.

EsLaRed Receives 10th Annual Postel Service Award

ISOC awarded the *Jonathan B. Postel Service Award* for 2008 to *La Fundación Escuela Latinoamericana de Redes* (EsLaRed) of Venezuela for its significant contributions to promote information technologies in Latin America and the Caribbean.

It is now ten years since the passing of Internet pioneer Jonathan B. Postel, the inspiration for this prestigious award. To mark this event in a special way, ISOC formed a *10th Anniversary Award Committee* including all the past award recipients, which has formally recognised EsLaRed for “its sustained efforts to bring scientific, technical, and social progress in Latin America and the Caribbean through education, research, and development activities on technology transfer.”

ISOC presented the award, including a US\$20,000 honorarium and a crystal engraved globe, in November during the 73th meeting of the IETF in Minneapolis, USA.

Accepting the award for EsLaRed was its President, Professor Ermanno Pietrosemoli. “We’re very excited to be honored in this way,” said Professor Pietrosemoli. “In the developing world, having access to the Internet, which gives us access to things like scientific journals and medical information, is not easy and it is not taken for granted. It is wonderful for us to be able to help people improve their conditions and to see first hand how the Internet can change people’s lives,” he said.

“On behalf of the ISOC community, it is my great pleasure to congratulate Professor Pietrosemoli and his dedicated colleagues at EsLaRed for their achievements over the years,” said Lynn St. Amour, President and CEO of ISOC. “EsLaRed’s commitment to the Internet has been at the forefront of regional development and their leadership has been an instrumental element in forming today’s dynamic Latin American and Caribbean Internet community,” said Ms St. Amour. For more information about this year’s recipient see:

<http://www.isoc.org/awards/postel/eslared.shtml>

The Postel Service Award was established by ISOC to honor individuals or organisations that, like Jon Postel, have made outstanding contributions in service to the data communications community. The award is focused on sustained and substantial technical contributions, service to the community, and leadership. Previous recipients of the Postel Award include Jon himself (posthumously and accepted by his mother), Scott Bradner, Daniel Karrenberg, Stephen Wolff, Peter Kirstein, Phill Gross, Jun Murai, Bob Braden and Joyce K. Reynolds (jointly), and Nii Quaynor. The award consists of an engraved crystal globe and a US\$20,000 honorarium. For more information see: <http://www.isoc.org/awards/postel/>

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ will contain standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2008 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA
--

The Internet Protocol Journal

March 2009

Volume 12, Number 1

A Quarterly Technical Publication for
Internet and Intranet Professionals

In This Issue

From the Editor	1
The End of Eternity	2
Resource Certification	13
Host Identity Protocol	27
Fragments	33
Call for Papers	35

FROM THE EDITOR

IP Version 4 address exhaustion and migration to IP Version 6 continues to be the focus of many Internet-related organizations and events. The *Regional Internet Registries* (RIRs), still debating what will happen as the IPv4 address pool runs out, are developing policies for how to manage address-block transfers between address holders. One potential result of the address shortage is that a *market* (official or otherwise) will develop for the buying and selling of IPv4 addresses. In our last issue, we brought you the first in a two-part series of articles entitled “The End of Eternity,” by Niall Murphy and David Wilson. Part Two, included in this issue, discusses what a market-based IP trading exchange might look like.

IP address allocation, transfers, and even the potential trading market for addresses is ultimately dependent on a reliable and trusted registry for this information. The RIRs have been working on a way to ensure that information about *IP Number Resources* (that is, IPv4 addresses, IPv6 addresses, and *Autonomous System* [AS] numbers) are securely stored and distributed so that users of such information can be assured that it is authentic. The underlying technology is a *Resource Certificate Public Key Infrastructure* (RPKI), and it is described in our second article by Geoff Huston.

The Internet technical community is discussing the so-called *identifier/locator split* as a major change to the Internet architecture. The IETF is developing several proposals, including the *Locator Identifier Separation Protocol* (LISP) discussed in our March 2008 issue. In this issue we look at another proposal, the *Host Identity Protocol* (HIP). The article is by Andrei Gurtov, Miika Komu, and Robert Moskowitz.

You will notice that our back cover has a new look. This layout is not the result of any creative design urges, but rather a change in U.S. Postal Service regulations regarding the placement of the subscriber address label. I guess the Internet isn’t the only place where addressing is a major topic.

As always, your comments, suggestions, and contributions are welcome, including Letters to the Editor, Book Reviews, and of course full-length articles. Our Call for Papers is included on page 35. Contact us by e-mail at ipj@cisco.com

—Ole J. Jacobsen, Editor and Publisher

ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

The End of Eternity

Part Two: Address Space Trading and the Routing Table

by Niall Murphy, Google, and David Wilson, HEAnet

In our last article^[0], we wrote about the onset of scarcity and the problems that are likely to ensue as a result. We characterized the problem we face as the *gap*, the length of time between the end of IPv4 plenty and the beginning of a universally reachable IPv6 Internet. Noting that any solution should either make the gap shorter, by bringing forward full IPv6 deployment, or make it less painful, by reducing the pressure of IPv4 scarcity, we propose that the fairest, most neutral way to encourage networks out of IPv4 while providing help for those who need it is to introduce a market-based IP address trading exchange. Let us explore now how such a system could work.

Possible Market Structures: Advantages and Drawbacks

An exchange could be set up and operated in many ways. Our preference, however, is for such a service to be run by the existing, trusted, and stable *Regional Internet Registries* (RIRs). Not only are they experienced in maintaining the values that the community as a whole wants to see maintained—fairness and neutrality, transparency, etc.—the RIRs are also in an excellent position to establish the *quality* of prefixes traded in an exchange, having excellent service contracts and history with members. Furthermore, the RIRs are unlikely to be made available for onward sale or transfer to other organizations with “different values,” and would maintain their traditionally community-focused policy-making apparatus. They would also be in a position to act quickly to coordinate and assume responsibility if given sufficient authority by the membership.

It does not have to be an RIR, of course: we *could* set up another industry body, but it would take valuable time and require a new governance model. We could also outsource the whole thing to any professionally run auction-handling site, but for such a fundamental change in how we do things, it seems wise to keep it under direct control. Finally, the psychology of continuity is important; if organizations are used to dealing with the RIRs, it provides an important perception of stability to keep them as the interface to getting new addresses.

As with our previous article, we emphasize again that the RIRs have provided excellent service in focusing the consensus of the community in a form that can be passed back to governments and other stakeholders, both external and internal.

The shield provided by the RIRs, protecting the members from the outside and protecting the members from themselves, has worked well for three reasons:

- First, RIR consensus is widely seen to broadly reflect the wishes of their communities as a whole because of the extremely low barrier to representation—in essence anyone who cares can attempt to influence policy, and no formal attempt is made to weigh one set of opinions over another. As a result, RIR policy is a lowest common denominator that is in general free from many of the more partisan stances usually found in the telecommunications arena, leading to greater credibility outside the RIR system, and greater credibility within, because the oppression of a minority by the majority within the context of policy formation is very difficult.
- Secondly, possessing that credibility has led to repeated success for the RIRs in the arena of disseminating and explaining policies outward, and they have therefore reinforced the confidence their members have in them.
- Finally, the RIRs are also comparatively financially easy to run; in the *Réseaux IP Européens* (RIPE) region, fees are by no means excessive given the ratio of customers to addresses; they are observed and validated by RIPE Network Coordination Centre (NCC) members, and any competing industry body would have to duplicate not only all the previously mentioned activities, but also the large working surplus that allows the RIRs to ensure stability through more turbulent times. Or to put it another way, “it’s open, it works, and it’s cheap.” We would recommend that any significant extension to the RIR authority, such as running an exchange as proposed, should endeavour to preserve as many of these properties as possible.

So if RIRs are to be the point of contact and policy making, how might such an exchange operate? We have a few guidelines from a relatively new field of economics, called *Market Design Theory*^[21], that might help to inform our choices. Firstly, we must have *thickness*: we must have enough traders (both buyers and sellers) entering the market, such that the populace at large can be assured that if they need to perform a transaction, the exchange is the place to do it, rather than private trades. (Private trades, although they enable liquidity, have the disadvantages that the WHOIS database is not maintained, that policy cannot be centralized, that prefix de-aggregation can occur arbitrarily, and so on.) We should avoid *congestion*: so many participants that it becomes difficult to trade. Finally, we must have *safety*: the assurance that if a transaction is engaged in, it will complete, and buyers will receive what they want.

Although other properties exist, those are the main ones required for the exchange to operate successfully. On thickness, we think it is clear that attracting buyers in a time of scarcity will not be a problem. The problem will be attracting sellers from such constituencies as have them available (old *Internet Assigned Numbers Authority* [IANA]-allocation holders, dot-com failures, and so on). An open question is whether the exchange can do more to attract sellers than the monetary reward for selling would do on its own; more meaningful incentives for them are difficult to determine. Overall, congestion does not seem likely to be a concern, given that the RIR model most usefully supports only membership-based participation initially. (Furthermore, our guess is that the “product” will be quite homogenous, so performing trades will presumably be mostly a matter of determining price.)

Let us return to the question of prefix *quality*. The single most important measure of quality of a prefix, the attribute without which the prefix is useless, is *uniqueness*. One must be assured that the prefix one holds is acknowledged as being held by oneself, and that *Internet Service Providers* (ISPs) will accept its announcement from *no other parties*.

From a plentiful pool, where prefixes have no cost other than the service charge of the registry, ensuring uniqueness is perhaps not a simple task, but it is a relatively uncontroversial one. When scarce, prefixes become valuable and will be given a cash value, either officially or by other means. ISPs will then have a business reason to break with consensus on routing filters, as we discuss later in more detail; but regardless, prefixes allocated from the IANA free pool generally have an impeccable heritage and do not vary greatly in usability. There are, of course, the natural delays in having new /8s incorporated into routing filters across the world. Those delays do have real effects, but the recipient of these prefixes usually has good reason to believe that a) these problems will be corrected over time, and b) everyone else in the same /8 will have the same problem.

In the new paradigm, each prefix must be carefully examined by the recipient to test that it is uniquely held by the proffering organization, and the recipient will presumably have a further interest in its routability and membership in blacklists. The quality problem arises in both private and public trades; if the RIRs implemented a quality test, that would be yet another advantage of centralization to the benefit of everyone.

Closely associated with prefix quality is the question of *safety*. Again the RIRs are in an excellent position to provide the necessary support for good-faith transactions, certification of prefixes being the primary mechanism, although various other possibilities (such as membership controls) might also exist.

More pertinently, pricing of the goods traded in such an exchange is an important question. Various natural calculations might support the calculation of address costs, including but not limited to average revenue per address, operational costs averaged over all addresses held, and so on. Our primary contention here is that the RIRs should not engage in price setting directly. Doing so would at the very least invite regulation. There may be a case for placing caps on trades as an antispeculation measure, but that requires further analysis.

What exactly the “goods” are in this case also needs consideration. Our preference is that what is traded is the right to use a prefix, rather than a prefix itself. Quite apart from the inherent oddness in selling a 32-bit integer (with 5-bit netmask), we should avoid the land registry model, where all the previous history of a prefix must be checked before sale. We need the RIR to intermediate itself and provide quality evaluation services rather than leaving it up to the end buyer. We should also not be selling rights to use prefixes of fixed sizes. The exchange needs to offer a spread of lengths in order to meet the needs of all potential customers.

You Say You Want a Revolution

To be sure, a change in the perceptual or legal status of IP addresses is a revolution in how we do things. The ramifications of IP addresses becoming property, or even acquiring intermediate states with property-like title rights, are manifold and they involve sweeping changes. Suddenly things that had no value have a clear public worth. Will organizations then be compelled to list addresses on their books as an asset? Could they then be taxed on them? What would such a tax rate be? Could organizations not actually using the asset (say, the RIRs) avoid this charge? Would transfers entail a taxable operation? These questions are significant and difficult. The right thing for the community is almost undoubtedly that IP addresses do not become simple property, but rather have (at a minimum) transfer and sale rights associated with them. In this way we could enable liquidity without complications, and avoid introducing extra complications at a difficult time. But it is unclear whether regulatory authorities will see it this way without the correct guidance.

The change in legal status of IP addresses is not the only violent change that could be unleashed by exhaustion. Consider, for example, the potential for litigation led by both new entrants unable to acquire an allocation to fulfill their business plan and incumbents seeking to either cause confusion (as an anticompetitive measure against just about anyone) or to try to disrupt any fragile consensus about how the last allocations play out. Leaving aside the question of whether simple prudence would recommend or deprecate such a move, there is a very clear risk of attempted litigation affecting the outcome of the end game.

However, one of the major benefits of a market is that it allows the RIRs to maintain a hands-off approach while still making it at least theoretically possible for an organization to get an independent allocation. The community can be doing all that it realistically can to continue the flow of IPv4, in terms of creating conditions fostering its dissemination, while being seen to be doing such, rather than simply running out of ideas and giving up. It could, of course, be seen—not unfairly—that participating in the transition to a market mechanism might amount to the effective transference of title to those who happened to be in the room at the time of exhaustion, an effective “insider privatization.”

Yet, if a market does not emerge, it is hard to see how any new entrants can have a business plan not directly dependent on incumbents. Although there are plenty of incumbents who would value having more address space to continue their business over the cash value of their addresses, so rendering entrance to the market impossible, there are plenty of other organizations that have only ever used a portion of their first allocation and would theoretically be well motivated to disburse these addresses accordingly.

To avoid exceptional attention from regulatory authorities, and to prevent the exchange from failing, we should design the exchange to deter in a systematic way the misbehavior of markets: speculation, hoarding, cartels, price fixing, and regional disadvantage should all be made as difficult as possible within the context of running a limited-membership market.

If we define *speculation* as short-term dealing with no expectation of use, we may be able to limit this kind of behavior naturally as a consequence of the membership-based participation inherent in the RIR model, and as a function of the periodic nature of routing filter generation. Increasing the price with short-term speculation disincentivizes the end purchaser with a use expectation from actually buying the prefix, because there will be a time delay before it can be used; therefore the purchaser with no use expectation will find it more difficult to find a buyer if the price rises to unreasonably high levels.

Hoarding, defined as long-term speculation with no use expectation, is bad for the exchange in that thickness is reduced, but also bad for the hoarder because the long-term value of the asset should decrease, in line with the increase in deployment of IPv6.

The formation of *cartels* would actually be quite a practical difficulty, especially under the closer attention likely to be paid to the exchange by competition authorities. Notwithstanding the coordination difficulties, we are inclined to say again that enough buyers should help to control this problem sufficiently to make the exchange work.

Regional disadvantage is, however we look at this situation, a problem. If scarcity is likely to lead to some monetary value being placed on address space, we face a vista where regional disadvantage can only be reduced, not eliminated. The inequality is, ultimately, one of the most compelling reasons to minimize the length of the transition period, and it would benefit us all to do so. Some measures go part way toward alleviating the problem. For instance, regional cooperation can help—in a market, if buyers cooperate and bulk buy, the threshold for organizations that would otherwise be facing a prohibitive barrier to entry would be reduced.

If we do not have a globally accessible exchange, it does not necessarily mean that the organizations will simply fail, entrenching the regional inequality, but they may respond by trying to fulfill their customer requirements by means of private, uncoordinated trading, with all the problems that entails.

We note that it is probably best to structure the actual trades as *auctions*, rather than facilitated marketplace transactions. When quality is asserted, one prefix is much like another—at least compared to prefixes of a similar size—and treating them as a commodity in this way facilitates the enforcement of policies on a centralized basis.

Drawbacks of a Market

Many cautionary tales about the operation of markets exist. Irrational exuberance, long-lasting depressions, fraudulent or exploitative behavior of all kinds—all of these effects, either enabled or supported by market mechanisms, are well known. Do we have any reason to believe either that these consequences will be not serious in our particular domain or that we have any new way of preventing them from happening?

In truth, we have no particular reason to believe that they won't happen, but there is a structural reason to believe that they might not matter to the exclusion of all else: the worse the situation becomes in the IPv4 marketplace, the more incentive there is to move to IPv6. To that extent, the market might be considered as providing a somewhat self-regulating reason for transition. Of course, we can put various mechanisms in place to help mitigate unstable behavior, as we suggested previously, but ultimately this is a fundamentally new way of doing things that we are ill equipped to understand the full consequences of.

Perhaps the largest drawback, outside of the practical difficulties in getting IPv4 addresses to organizations, is the philosophical impediments that come inherent with switching to a market-based model for allocation. Although a market cannot be said to rule out the consensus model that has turned out well for the Internet community, it also cannot be said to fully support it. This change may be a cultural one we find difficult to reverse, and it might undermine any future attempt by the community to try to differentiate itself on governance model.

Even though we have proposed the market model in good faith, as an attempt to meet the needs of new entrants and existing organizations—and as a boost to the faster deployment of IPv6—if it proves to be a failure in meeting those needs, there may be no more credible strategies left if governments insist on action. That in itself might represent even larger, more unpredictable change for the industry.

Effects on the Routing Table

Another inescapably important question is what will happen to the *Default Free Zone* (DFZ) routing table. A world in which address blocks transfer without the aggregating procedures of the RIRs is naturally a cause for concern, and when needs-based allocation comes to an end, a change in the rate of growth does seem inevitable. We can, however, make some observations that might reassure us, to some extent, that the rate of growth will not be calamitous.

First, as we go from a time of address plenty to address scarcity, one can assume that the ongoing fulfilled demand for address space will be no greater than it is now. Hence, the future growth in the number of prefixes in the routing table—regardless of prefix length—would seem to have an upper limit consistent with the number of allocations by RIRs to *Local Internet Registries* (LIRs) at the moment. This limit is still a multiple of the current curve, because we lose the benefit of the aggregation function performed by LIRs, but it suggests that we will at least not face an order-of-magnitude step change as a result of a disorderly competition.

Then there is the question of the routability of smaller prefixes. There is, at the moment, a *de facto* longest prefix size of around /24 that has close to universal reachability on the general Internet. One might assume that this prefix size will grow inexorably during and after exhaustion, as existing space is broken up into smaller and smaller blocks. Implicit in that assumption is the notion that such block sizes will be adequate for users and worthwhile for ISPs to route; we should probably not rely on networks “making do” with smaller and smaller chunks of address space.

Simultaneously, inexorably growing prefix lengths in the DFZ can only come about because of operator action. In particular, although there is a rough consensus in DFZ operators at the moment that /24 is routable and /25 is not, this policy is not a consensus-approved policy of the RIRs or the IETF. Each operator makes its own decision, based on its own customer needs, its own network, and the expectation of routability with other networks.

Reachability, therefore, depends on ISPs cooperating, and universal reachability depends on ISPs cooperating universally. An ISP may well choose to carry smaller prefixes on behalf of its customers, but unless this practice becomes widespread, no expectation can be made of universal reachability, and the practice will remain a minority one conducted by cooperating ISPs, as occasionally happens from time to time today, and this situation will little affect the size of the routing table for those involved.

Is there a competitive advantage to the largest of the ISPs in investing in very large routers that can carry many millions of prefixes, more than the smaller ISPs can support? If there were, it could perhaps lead to a concentration of power in the tier-one providers (who, as inevitable parts of any lengthy path across the Internet, have the greatest influence on the *de facto* longest routable prefix.) This situation could perhaps be true if routers are price-limited by the supportable number of prefixes, but this characteristic is typically a secondary one at worst. Routers are grouped by the bandwidth they can support, and priced accordingly; a 100-Mbps router that can support a million prefixes will certainly be more expensive than a 100-Mbps router that can support only ten thousand, but there is an order of magnitude step from either router to a router that can support 10 Gbps.

Inaction Leads to Harm

In fact the argument that the effect on the routing table will be unsustainable is opposed to the argument that there may not be adequate liquidity to sustain the market. It is true that we could find ourselves in the latter position, and so the effect of this system on reducing the problem (characterized as “the gap”) will be smaller than we might like—but, as a best-effort scenario, not negligible, particularly in regard to showing good stewardship of the resource to potential outside influence. Compared to any other proposal, and particularly compared to voluntary release or a locking down of the address space, we think that this way is the best way to assure that we make available what liquidity there is.

It is difficult to see any model—even an idealized one—that could possibly service the run rate while maintaining aggregatability. The sparse allocation model used by the RIRs is dependent upon the continued availability of large, clean blocks of space, that is, /8s from IANA. With this address plenty comes freedom in our choice of policies, and with that freedom comes relatively quick consensus.

Post-exhaustion, the space will not be plentiful, and regardless of whether a monetary cost is attached, it will no longer be free. At this point, the legitimacy of the consensus of the RIR fora becomes critical. It is a fiercely defended bottom-up process. As the legitimacy of policies in the *Domain Name System* (DNS) world comes from consensus to abide by a single `root.cache`, so the legitimacy of policies in routing comes from general agreement on route filters and the authenticity of data in the RIR WHOIS databases.

We have also learned from the DNS world what happens to operational consensus when the resource becomes in some way valuable. Although the current RIR meetings are able to come to decisions that roughly reflect the consensus of the operational Internet, the necessarily tougher decisions forced upon us will challenge those who participate directly in policy making to reach conclusions that will satisfy operators who are not present. In principle it should not be necessary to account for those who do not represent themselves, but when the legitimacy of our policies is derived from their operational choices, the burden rests on us to ensure that our processes are truly representative.

If we are unsuccessful in doing so, or indeed if we choose to maintain the status quo, we cannot assume that the policies implemented on the operational Internet will themselves remain static. It is already the case that ISPs will work together, as is their entitlement, to agree to route prefixes for the benefit of their mutual customers. It is not unusual for one ISP to accept the announcement from a customer of a subnet of another ISP's address space. This decision is one for those ISPs to make about their own operational environments.

If we choose not to endorse a particular short-term solution to depletion, it falls upon ISPs themselves to find a way to continue their business operations, and resolve their customers' problems. If they cannot get address space from themselves, it will be their *duty* to their customers to get routable address space from somewhere—by negotiating, if necessary, with their peers and upstream providers to change the definition of “routable address space.” Ultimately we may assume that if we do not provide a solution to the industry, the industry will invent one—or several competing ones.

Because we assert that the solution that best solves this problem is an address space trading exchange, we may well end up getting one—but one (or more) that is private, and out of sight of our existing policy-making structure. Worse still, competing exchanges would not have access to the RIRs data, and so would not be in a position to assure the quality of a prefix—a situation that could threaten all transactions.

Without exaggerating, it is likely that what we do in response to this crisis will determine the architecture of the Internet for a long while to come. Although we are reminded of Woody Allen's quote wherein he “... hope[s] mankind has the wisdom to choose correctly... between utter hopelessness and total extinction^[22, 23],” there are, as we have outlined, measures we can take to survive the coming storm. They are not beautiful solutions. They are not how we have traditionally done things, or even how we would like to do things. Adopting them will almost certainly result in someone being worse off than if we had simply done nothing. But they represent, to our minds, the best, most realistic chance to avoid widespread difficulties and the loss of many of the principles we in the networking community hold dear, to ourselves and in our institutions. Let us begin this process now.

Acknowledgements

The authors would like to gratefully acknowledge help and support from Léan Ní Chuilleanáin, Emma Apted, and David Malone for diligent editing.

References

- [0] Murphy, Niall and Wilson, David, “The End of Eternity Part One: IPv4 Address Exhaustion and Consequences,” *The Internet Protocol Journal*, Volume 11, No. 4, December 2008.
- [1] <ftp://ftp.ietf.org/ietf-online-proceedings/94dec/area.and.wg.reports/ipng/ale/ale-minutes-94dec.txt>
- [2] <http://tools.ietf.org/html/rfc2008>
- [3] Hain, Tony, “A Pragmatic Report on IPv4 Address Space Consumption,” *The Internet Protocol Journal*, Volume 8, No. 3, September 2005.
- [4] <http://playground.sun.com/ipv6/doc/history.html>
- [5] <http://ipv4.potaroo.net>
- [6] <http://www.ripe.net/ripe/meetings/ripe-55/presentations/murphy-simlir.pdf>
- [7] http://www.isoc.org/educpillar/resources/ipv6_faq.shtml
- [8] <http://www.ietf.org/internet-drafts/draft-narten-ipv6-statement-00.txt>
- [9] <http://www.apnic.net/meetings/24/program/sigs/policy/presentations/el-nakhal-prop-051.pdf>
- [10] <http://www.ripe.net/ripe/policies/proposals/2007-06.html>
- [11] http://www.switch.ch/pki/meetings/2007-01/namebased_ssl_virtualhosts.pdf
- [12] For example,
http://h.root-servers.org/128.63.2.53_2.html versus
http://h.root-servers.org/h2_5.html
- [13] <http://www.ripe.net/ripe/meetings/ripe-55/presentations/vegoda-reclaiming-our.pdf>
- [14] A “smooth and convenient” dialing plan for India.
<http://www.mycoordinates.org/indias-phone-june-06>
- [15] http://en.wikipedia.org/wiki/UK_telephone_code_misconceptions

- [16] <http://code.google.com/p/simlir/>
- [17] <http://www.ripe.net/docs/ripe-407.html#membership>
- [18] <http://www.ripe.net/ripe/policies/proposals/2007-03.html>
- [19] <http://www.ripe.net/ripe/policies/proposals/2007-06.html>
- [20] <http://www.ripe.net/ripe/policies/proposals/2007-07.html>
- [21] <http://kuznets.fas.harvard.edu/~aroeth/alroth.html>
- [22] Woody Allen, "Side Effects," 1980.
- [23] Woody Allen through (most famously) Stephen Hawking,
<http://www.cnn.com/2006/WORLD/asiapcf/07/04/talkasia.hawking.script/index.html>
- [24] <http://icann.org/en/announcements/proposal-ipv4-report-29nov07.htm>
- [25] <http://www.ripe.net/ttm/>
- [26] <http://www.ripe.net/ripe/tf/enhanced-cooperation/index.html>
- [27] <http://www.nro.net/documents/nro18.html>
- [28] <http://www.ripe.net/maillists/ncc-archives/im-support/2004/index.html>

NIALL MURPHY holds a B.Sc. in Computer Science and Mathematics from University College Dublin. While in university, he founded the UCD Internet Society, which provided Internet access to approximately 5000 students. He went on to work for (and found) various organizations: the .IE domain registry, Club Internet (now Magnet Entertainment), Ireland On-Line, Enigma Consulting, Bitbuzz, and Amazon.com. He is currently in Site Reliability Engineering at Google. He is the coauthor of numerous articles, some RFCs, the O'Reilly book *IPv6 Network Administration*, and is a published poet and keen amateur landscape photographer. E-mail: niallm@avernus.net

DAVE WILSON holds a B.Sc. in Computer Science from University College Dublin, not coincidentally from around the same time as Niall. He has worked at HEAnet, the Irish National Research & Education Network, for more than 10 years, maintaining an involvement with RIPE and with the pan-European research network Géant. Dave is a member of the ICANN Address Supporting Organization Address Council; he helped to found the Irish IPv6 task force, which has the support of the national government there. E-mail: dave.wilson@heanet.ie

Resource Certification

by Geoff Huston, APNIC

Opinions vary as to what aspect of the Internet infrastructure represents the greatest common vulnerability to the security and safety of Internet users, but it is generally regarded that attacks that are directed at the network infrastructure are the most insidious, and in that case the choice is probably between the *Domain Name System* (DNS) and the interdomain routing system.

The question of how to improve the robustness of these functions has been a longstanding topic of study. For the DNS it appears that there is convergence on *Domain Name System Security Extensions* (DNSSEC) as the technical solution to securing DNS resolution operations, and the focus of attention in this space has shifted from technical behavior to topics relating to operational deployment. It has been a difficult time for DNSSEC and to say that there is an end in sight may well be premature at this stage, but there are definite signs of progress in this space. The same cannot be said of progress with securing routing, and particularly in securing interdomain routing. Here much remains to be done in order to achieve reasonable consensus on what technical measures to adopt, let alone the second step of study of how such measures could be deployed across the Internet.

The IETF's approach to addressing the topic of securing interdomain routing has followed a conventional IETF path. The first step has been to consider the nature of various vulnerabilities that exist within today's interdomain routing system and then develop a set of requirements that should be addressed in any solution space, without necessarily defining what such a solution may be. When the enumeration of requirements achieves a suitable level of consensus from the community, it is then possible to commence work on standardizing solutions. In the case of securing interdomain routing, the first steps were undertaken in *Birds of a Feather* (BOF) sessions and in the subsequently formed *Routing Protocol Security Requirements* (RPSEC) Working Group. This work is almost complete, and apart from some definitive statement relating to a requirement for securing the *Autonomous System* (AS) Path attribute in *Border Gateway Protocol* (BGP), the set of requirements for securing interdomain routing is now in an almost final state^[1]. The task of the *Securing Inter-Domain Routing* (SIDR) Working Group is to standardize technologies that can meet these requirements.

So where does "Resource Certification" fit in?

Public Key Cryptography

One commonly used security technology is *Public Key Cryptography*, a technique that is easily explained. The approach uses a pair of keys, A and B. Anything enciphered with key A can be deciphered only with key B, and conversely, and knowledge of the value of one key does not lead to discovery of the value of the other key. Key A is kept as a closely guarded secret, whereas key B is openly published. If I want to send you a message that only you can decipher and read, I should encrypt it using your public key. If I want to send you a message that only I could have sent (nonrepudiation), then I will generate a digital signature of the message using my private key. That way any attempts to alter the message will also be detectable.

This latter approach, of using keys to generate digital signatures of messages, lies at the heart of DNSSEC, because DNSSEC adds public keys and digital signatures to the DNS. A DNS query can generate a response that lists both the DNS answer and the digital signature of that answer. The DNS can also be queried to retrieve the public key used to sign all the components of that zone, so that the digital signature can be verified and the query agent can be assured that the response is a genuine one. But how can the key itself be verified? In DNSSEC the hierarchical nature of the DNS itself is exploited by having each zone “parent” sign the keys of its delegated “children.” So the zone key can be verified by retrieving the parent’s signature across that zone key, and so on to the root of the DNS. As long as the query agent knows beforehand the value of the public key used to sign the root zone of the DNS, and as long as DNSSEC is used universally, all DNS responses can be verified in DNSSEC.

Although this approach works in the interlocked hierarchical structure of the DNS, when we turn our attention to securing the use of IP addresses and AS numbers in the context of interdomain routing, there is no comparable hierarchy to exploit. In such cases a common solution is to turn to *Digital Certificates*.

Digital Certificates are digitally signed public attestations by a certification authority that associate a subject’s public key value with some attribute of the subject. A typical application is in identity certification, where the certification authority is attesting that the holder of the private key whose matching public key is provided in the certificate has met the authority’s certification criteria to be identified by a particular name. Digital certificates are useful in that they can reduce the number of trust points in a security domain, so that each member of the domain does not have to validate identity and exchange public keys with every other member of the domain, but can undertake a single transaction with a certification authority that is trusted by all the members of the domain. As long as every member of the domain carries the public key of the certification authority and can access all issued digital certificates, then the members of the domain can verify each other’s attestations and digital signatures.

Of course digital certificates are used for far more than attestations of identity, and can encompass the authority to perform specific tasks, undertake particular roles, or grant permissions and right-of-use authorities. It is this latter use case that is relevant to resource certification.

Resource Certificates

A Resource Certificate is a conventional X.509 certificate that conforms to the *Public Key Infrastructure Working Group* (PKIX) profile (RFC 5280) with one critical component, namely a certificate extension that lists a collection of IP number resources (IPv4 addresses, IPv6 addresses, and AS numbers)^[17].

These certificates attest that the certificate issuer has granted to the entity represented by the certificate subject a unique “right-of-use” of the associated set of IP number resources listed in the certificate extension, by virtue of an associated resource allocation. The unique “right-of-use” concept mirrors the resource allocation framework, where the certificate provides a means of third-party validation of assertions related to resource allocations^[2].

By coupling the issuance of a certificate by a parent *Certification Authority* (CA) to the corresponding resource allocation, a test of the validity of a certificate, including the IP number resource extension, can also be interpreted as validation of that resource allocation. Signing operations that descend from that certificate can therefore be held to be testable, under the corresponding hierarchy of allocation. In other words, if you received your address block from a particular *Regional Internet Registry* (RIR), then only that RIR can issue a Resource Certificate for you that includes your public key and the allocated number resources. Anything you sign using your private key can be verified through the RIR’s issued certificate.

Unlike certificates that relate to attestations of identity, Resource Certificates are not necessarily long-lived. When an additional allocation action occurs, the associated Resource Certificate is reissued with an IP number resource extension that matches the new allocation state. In the case of a reduction in allocated resources, the previously issued certificates are explicitly revoked when the new certificate is issued. In other cases there is no explicit revocation of the older certificates.

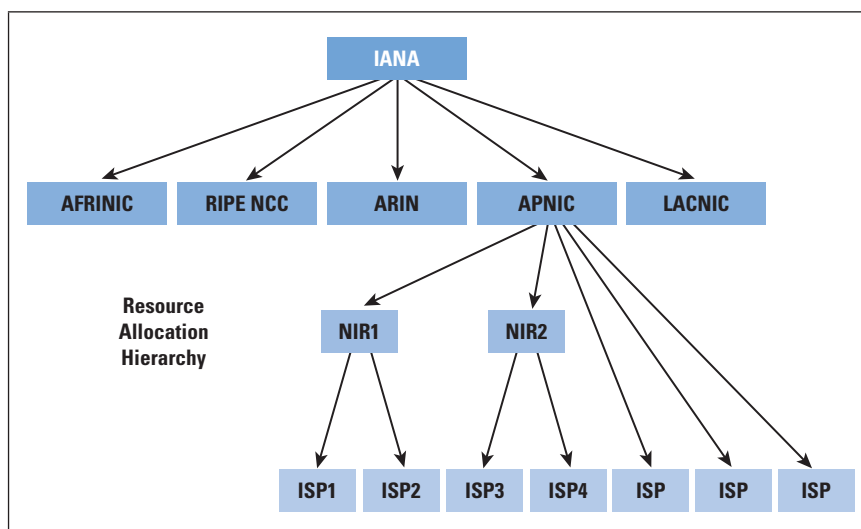
The intention here is that any instrument signed by the subject’s private key that relates to an assertion of resource control, whether it is a protocol message in a routing protocol or an administrative request to an *Internet Service Provider* (ISP) to route a prefix or as assertion of title over the “right-of-use” of a number resource, can be validated through the matching public key contained in the certificate and the IP number resource that is enumerated in this certificate. The Resource Certificate itself can be verified in the context of a Resource Certificate *Public Key Infrastructure* (PKI).

The Resource Certificate Public Key Infrastructure

The *Resource Certificate Public Key Infrastructure* (RPKI) describes the structure of the certification framework used by Resource Certificates. The intent of the RPKI is to construct a robust hierarchy of X.509 certificates that allows relying parties to validate assertions about IP addresses and AS numbers, and their use.

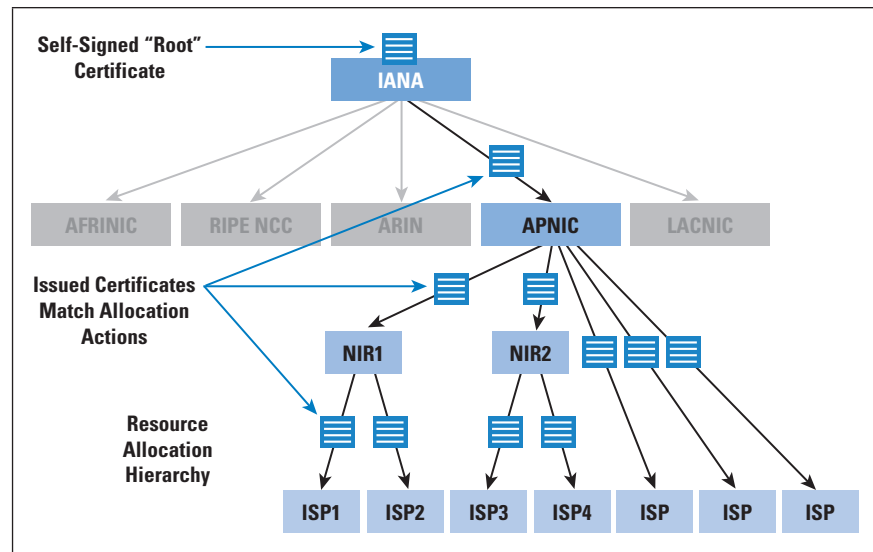
The structure of the RPKI as it relates to public use of IP number resources is designed to precisely mirror the structure of the distribution of addresses and ASs in the Internet, so a brief description of this distribution structure is appropriate. The *Internet Assigned Numbers Authority* (IANA) manages the central pool of number resources. The IANA publishes a registry of all current allocations. The IANA does not make direct allocations of number resources to end users or *Local Internet Registries* (LIRs), and instead allocates blocks of number resources to the RIRs. The RIRs perform the next level of distribution, allocating number resources to LIRs, *National Internet Registries* (NIRs), and end users. NIRs perform allocations to LIRs and end users, and LIRs allocate resources to end users (Figure 1).

Figure 1: Address Distribution Hierarchy for the Internet



The RPKI mirrors this allocation hierarchy. One interpretation of this model would send the IANA manager a root RPKI key, and using this key the IANA would issue a self-signed “root” certificate, and also issue subordinate certificates to each of the RIRs, describing in the resource extension to the certificate the complete set of number resources that have been allocated to that RIR at the time of issuance. The certificate would also hold the public key of the RIR and would be signed by the private key of the IANA. Each RIR would issue certificates that correspond to allocations made by that RIR, where the resource extension to those certificates lists all the allocated resources, and the certificate includes the public key of the recipient of the resource allocation, signed with the private key of the RIR. If the recipient of the resource allocation is an LIR or an NIR, then it too would also similarly issue resources certificates (Figure 2).

Figure 2: RPKI Resource Certificate Hierarchy



The common constraint within this certificate structure is that an issued certificate must contain a resource extension that contains a subset of the resources that are described in the resource extension of the issuing authority's certificate. This requirement corresponds to the allocation constraint that a registry cannot allocate resources that were not allocated to the registry in the first place. One implication of this constraint is that if any party holds resources allocated from two or more registries, then it will hold two or more Resource Certificates in order to describe the complete set of its resource holdings.

Validation of a certificate within this RPKI is similar to conventional certificate validation within any PKI, namely establishing a chain of valid certificates that are linked by issuer and subject from a nominated trust anchor CA to the certificate in question. The only additional constraints in the RPKI are that every certificate in this validation path must be a valid Resource Certificate, and the IP number of resources described in each certificate must be a subset of the resources described in the issuing authority's certificate.

Within this RPKI all Resource Certificates must have the IP addresses and AS resources present, and marked as critical extensions. The contents of these extensions correspond exactly to the current state of IP address and AS number allocations from the issuer to the subject.

Any holder of a resource who can make further allocations of resources to other parties must be able to issue Resource Certificates that correspond to these allocations. Similarly, any holder who wishes to use the RPKI to digitally sign an attestation needs to be able to issue an *End Entity* (EE) certificate to perform the digital signing operation.

For this reason all issued certificates that correspond to allocations are certificates with the CA capability enabled, and each CA certificate is capable of issuing subordinate CA certificates that correspond to further sub-allocations and subordinate EE certificates that correspond to a generation of digital signatures on attestations.

The RPKI makes conventional use of *Certificate Revocation Lists* (CRLs) to control the validity of issued certificates, and every CA certificate in the RPKI must issue a CRL according to the nominated CRL update cycle of the CA. A CA certificate may be revoked by an issuing authority for numerous reasons, including key rollover, the reduction in the resource set associated with the subject of the certificate, or termination of the resource allocation. To invalidate the authority or attestation that was signed by a given EE certificate, the CA issuing authority that issued the EE certificate simply revokes the EE certificate.

Resource Certificates are intended to be public documents, and all certificates and objects in the RPKI are published in openly accessible repositories. The set of all such repositories forms a complete information space, and it is fundamental to the model of securing the public Internet interdomain routing system that the entire RPKI information space is available. Other uses of the RPKI might permit use of subsets, such as the single chain from a given end-entity certificate to a trust anchor, but routing security is considered against all known publicly routable addresses and AS numbers, so all known resource certification outcomes must be available. In other words the intended use of the RPKI in routing contexts is not a case where each relying party may make specific requests for RPKI objects in order to validate a single object, but one where each relying party will perform a regular sweep across the entire set of RPKI objects in order to ensure that the relying party has a complete picture of the RPKI information space.

This aspect of the RPKI represents some interesting challenges, in that rather than having a single CA publish all the certificates produced in a security application at a single point, the RPKI permits the use of many publication points in a widely distributed fashion. Each CA can issue RPKI objects and publish them using a locally managed publication point. It is incumbent upon relying parties to synchronize a locally managed cache of the entire RPKI information space at regular and relatively frequent intervals.

For this reason the RPKI has introduced an additional mechanism in its publication framework, namely the use of a “manifest” to allow relying parties to determine whether they have been able to retrieve the entire set of RPKI published objects from each RPKI repository publication point, or if there has been some attempt to disrupt the relying party’s access to the entire RPKI information set.

It also implies that the RPKI publication point access protocols should support the efficient function of a synchronization comparison, so that a locally managed cache of the RPKI need only call for the uploading of those objects that have been altered since the previous synchronization operation.

Signed Attestations and Authorities

The underlying intent of digital certificates, and Resource Certificates in particular, is in terms of supporting a transitive trust relationship that allows a relying party to verify the authenticity of a signed artefact through verification of the signer's key using the PKI. So the obvious question is: what artefacts are useful to sign?

Much of the motivation for Resource Certificates has come from a desire to underpin efforts in securing aspects of interdomain routing. This effort goes well beyond securing the individual point-to-point connection used between BGP speakers, and refers to the matter of verifying the authenticity of the payload of the BGP protocol exchange. The specific question that may be posed is: how can a BGP speaker validate the authenticity of the route object being presented to it?

The approach being studied by the SIDR Working Group is to use structured attestations, where, like the digital certificate itself, the attestation is structured in an ASN.1 digital object, and this object is signed using a signing formation that is itself a piece of structured ASN.1, namely the *Cryptographic Message Syntax* (CMS)^[18].

The first of these attestations relates to the ability to verify the authenticity of the “origination” of an interdomain routing object. This verification refers to the address prefix and the originating AS, and the questions that this verification function is intended to answer include:

- Is this a valid address prefix and AS number? Have these resources been allocated through the IP number resource allocation process?
- Has the holder of the title of “right-of-use” for the address prefix authorized the AS holder to originate a routing advertisement for this prefix?

Here an address holder is authorizing a particular ISP to generate a route announcement for its particular address prefix. In this case the prefix holder would generate an EE Resource Certificate with the IP number resource extension spanning the set of addresses that match the address prefixes that are the intended subject of the routing authority, and place validity dates in the EE certificate that correspond to the intended validity dates of the routing authority.

The signed authority document would contain the AS number that is being authorized in this manner, a description of the range of prefixes that the prefix holder has authorized, and the EE certificate. The document would be signed by the EE certificate private key using a CMS signing structure. The resultant object is published in the RPKI distributed publication repository as a *Routing Origination Authorization* (ROA). A relying party can validate the ROA by checking to ensure that the digital signature in the ROA is correct, indicating that the authority document has not been tampered with in any way since it was signed, that the resources in the associated EE certificate encompass the prefixes specified in the document, and the EE certificate itself is valid in the context of the RPKI by verifying that there is an issuer-subject chain of valid certificates that link one of the relying party's nominated trust anchors to the EE certificate.

The ROA itself is valid as long as the signing EE certificate is valid. To withdraw the authority prior to the expiration of the EE certificate, the ROA publisher can simply revoke the EE certificate, leading to the concept of "one-off-use" EE certificates in the RPKI, where a key pair and a corresponding EE certificate are generated in order to sign a single attestation or authority. If the authority's lifetime is extended, the authority is reissued with a new EE certificate and a new digital signature, and, as noted, the authority can be prematurely terminated through revocation of the EE certificate, so at no stage is there a need to reuse the original signing private key. After the private key is used to sign this object, the key is destroyed, alleviating to some extent the key management load.

In any security system knowledge of what is authorized is helpful, but knowledge of what has not been authorized is perhaps even more helpful. For ROAs there is an analogous situation to DNSSEC, where DNSSEC is most effective from a client's perspective after the entire DNS space is DNSSEC signed. Where there are gaps in the DNSSEC signing chains the client is left in an uncertain state regarding the verification outcomes of the unlinked DNS sub-hierarchies. The same could apply to ROAs, in that in an environment where not every originated route object has a published ROA, the absence of a ROA does not necessarily indicate an unauthorized route origination. If one of the objectives of this study is to define a framework that can unambiguously identify the unauthorized use of IP number resources in routing (route "hijacks") even in a world where ROAs are used in a piecemeal fashion, then one possible refinement to the ROA model is the introduction of a comparable negative authority, the *Bogon Origination Attestation* (BOA).

In this case the prefix holder generates a signed attestation, or BOA, in a similar manner to the ROA, but does not provide any originating AS. Instead the BOA refers to "all originating ASs," and has the semantic interpretation that any use in the routing space of this address prefix described in the BOA, or any more specific address prefix, should be regarded as unauthorized and the route should be discarded.

Although this process makes the detection of route hijacks more direct in a world of piecemeal use of ROAs, there is now the added complication of having both “positive” and “negative” authorities. The proposed resolution of this dilemma is to use a relative priority rule that ROAs take precedence over BOAs, so that if a valid ROA and a valid BOA both exist that describe the origination component of a route, then the route can be regarded as authorized.

It should be noted, however, that at this stage these concepts are “work in progress,” and are part of the SIDR Working Group’s agenda of study, and the working group has not as yet reached any consensus regarding the decision to advance these proposals onward along the Internet Standards Process.

Also on the near-term horizon for SIDR is examining approaches to secure the AS path in BGP updates. The RPSEC Working Group has explored two approaches in this space. One involves an incremental multiple signature technique that allows a receiver of a BGP update to verify that the AS path described in the update is matched by a sequence of interlocking AS digital signatures using the RPKI. At the same time that an AS adds its own AS to the AS path prior to further *External Border Gateway Protocol* (eBGP) propagation of the route update, the AS would digitally sign over an analogous sequence of AS signatures. This approach allows a receiver to perform a match of the AS sequence in the AS path with the AS number sequence identified in the AS signature block. A match here would indicate that the BGP update has indeed been sequentially passed along the sequence identified by the AS path. This approach was originally proposed in the *Secure BGP* (sBGP) design^[21] and has attracted some comment related to the computation overhead associated with the application and validation of these AS path signature sequences. An alternative approach has been one that is described by RPSEC as being less rigorous, and refers to a “feasibility” check, which checks to ensure that each pair of ASs represented in the AS path has an associated verifiable assertion of inter-AS adjacency that is digitally signed by both ASs.

It should also be noted that this activity of addressing aspects of improving the robustness of interdomain routing has some previous context. In many parts of the Internet, some degree of routing integrity is managed through the use of *Internet Routing Registries* (IRRs) and the publication of routing policies through the use of *Routing Policy Specification Language* (RPSL) objects.

Although opinions vary as to the robustness of the security offered by the IRR approach, at the very least it can mitigate some weakness in the routing system through the use of a “second check” that can be used to filter the information that is being provided in a BGP feed.

The weaknesses in the IRR system tend to relate to the consistency, completeness, and authenticity of the IRR data, and in many cases the trust in the integrity of the data relies on the admission practices of the IRR itself, and individual data objects cannot be verified by clients of the IRR. One possible way to address this situation has been through the use of *Routing Policy System Security* (RPSS) measures, but the adoption of these measures has not been widespread, and the question still remains for the client that even if an IRR object was authenticated upon admission, it does not mean that when the object is subsequently used by an IRR client the information reflects the current situation, and the information could well be invalid or not reflect the current policies of the author of the IRR object.

One possible approach being considered by the SIDR Working Group is to implement the RPSS authentication models using object signing in the context of the RPKI. For example, the RPSS assumption that routes should be announced only with the consent of the holder of the origin AS number of the announcement and with the consent of the holder of the address space implies in RPSS that both parties should authorize the entry of a *route object* into the IRR. Translating this stipulation into an analogous model using the RPKI would require that a route object be signed with the digital signatures of both the AS holder and the address space holder, and a IRR client can verify this route object at the time of use by verifying both digital signatures. Either the address space holder or the AS holder can revoke authorization by revoking the EE certificate used to sign the route object, and the verification is independent of the particular IRR that has published the route object. It is also a possibility that the IRR itself can be folded into the RPKI distributed publication repository framework, because there is no particular requirement in such an environment for a disparate collection of IRRs with their own partial collections of routing policy information, although at this stage this discussion is heading into the realm of more advanced speculation about the potential for application of Resource Certificates and digital signatures to RPSL and the IRR framework.

Putting Resource Certificates into Context

Resource Certificates and the associated RPKI represent a major part of any effort to construct a secure interdomain routing framework. An RPKI, even partially populated with signed information, allows BGP speakers to make preferential selections to use routing information where the IP address block and the AS numbers being used are recognized as valid to use, and the parties using these IP addresses and AS numbers are properly authorized to so do. The RPKI can also be used to identify instances of unauthorized use of IP addresses and attempts to hijack routes.

However, the RPKI represents only one part of a larger framework of securing interdomain routing, and the next step is that of applying the RPKI to the local BGP processing framework. There is also the need to move beyond validation of route origination and look at the associated topic of validation of the AS path, and potentially to consider the most challenging task, of attempting to validate whether the initial forwarding decision associated with a route object actually represents the correct first hop along a usable forwarding path for packets to reach the network destination.

The concerns here include not only a consideration of what can be secured and validated, but matters of scalability and efficiency in terms of deployment cost. The various approaches to routing security studied so far offer a wide variety of outcomes in terms of the amount of routing information that is validated, the level of trust that can be placed in a validation outcome, and the overheads of generating and validating digital signatures on routing information. The next step appears to include the task of establishing an appropriate balance between the overheads of operating the security framework and the extent to which efforts to disrupt the routing system can be successfully deflected by such measures.

The RPKI has been designed as a robust, simple framework. As far as possible existing technologies and processes have been exploited, reflecting to some extent a level of conservatism of the routing community and the difficulty in securing widespread acceptance of novel technologies.

References and Further Reading

The following documents provide further detail about the IETF work on resource certification. The Internet Drafts listed here are still a “work in progress,” and although they are reflective of the areas of activity of the SIDR Working Group, they do not necessarily represent finished work.

Internet Drafts

Requirements:

- [1] B. Christian, T. Tauber, eds., “BGP Security Requirements,” work in progress, Internet Draft, **draft-ietf-rpsec-10.txt**, November 2008. *The report of the consensus outcomes of the RPSEC Working Group in enumerating the requirements for securing interdomain routing. The outstanding topic in this report remains in the area of AS path validation and the level of requirement associated with the two approaches described in the report.*

Architecture:

- [2] M. Lepinski, S. Kent, “An Infrastructure to Support Secure Internet Routing,” work in progress, Internet Draft, **draft-ietf-sidr-arch-04.txt**, November 2008. *An overview of the RPKI approach, describing the RPKI, the distributed repository structure, and common operations.*

Resource Certificates:

- [3] G. Huston, G. Michaelson, R. Loomans, “A Profile for X.509 PKIX Resource Certificates,” work in progress, Internet Draft, **draft-ietf-sidr-res-certs-15.txt**, November 2008. *The specification of the Resource Certificate.*

RPKI Repository Structure:

- [4] G. Huston, G. Michaelson, R. Loomans, “A Profile for Resource Certificate Repository Structure,” work in progress, Internet Draft, **draft-ietf-sidr-repos-struct-01.txt**, October 2008. *A description of the proposed distributed publication repository structure for the RPKI, including contents, access protocols, and object name conventions.*
- [5] R. Austein et al., “Manifests for the Resource Public Key Infrastructure,” work in progress, Internet Draft, **draft-ietf-sidr-rpki-manifests-04.txt**, October 2008. *A specification for repository manifests. Manifests are signed constructs that describe all the objects currently loaded into a repository publication point, and are used by relying parties as a means of ensuring that a local RPKI repository cache is correctly synchronized against the authoritative original publication point.*
- [6] G. Huston, R. Loomans, B. Ellacot, R. Austein, “A Protocol for Provisioning Resource Certificates,” work in progress, Internet Draft, **draft-ietf-sidr-rescerts-provisioning-03.txt**, August 2008. *A proposed protocol for use between a subject and a certificate issuer to ensure that certificate requests, the IP number resource allocation state, and the issued certificate status are correctly synchronized. This synchronization extends the conventional certificate request model into a transaction protocol that also includes the ability to perform certificate revocation requests and status queries from the subject.*

RPKI Signed Objects:

- [7] M. Lepinski, S. Kent, D. Kong, “A Profile for Route Origin Authorizations (ROAs),” work in progress, Internet Draft, **draft-ietf-sidr-roa-format-04.txt**, November 2008. *The specification of the syntax for signed ROAs.*
- [8] G. Huston, T. Manderson, G. Michaelson, “A Profile for Bogon Origin Attestations (BOAs),” work in progress, Internet Draft, **draft-ietf-sidr-bogons-02.txt**, October 2008. *The specification of the syntax for signed BOAs.*
- [9] G. Huston, G. Michaelson, “Validation of Route Origination in BGP Using the Resource Certificate PKI,” work in progress, Internet Draft, **draft-ietf-sidr-roa-validation-01.txt**, October 2008. *The specification of the semantics of ROAs and BOAs and the manner in which these objects may be interpreted in terms of the integration of these origination security credentials onto a BGP route-selection process.*

Certificate Policy and Practice Statements:

- [10] K. Seo, R. Watro, D. Kong, S. Kent, “Certificate Policy (CP) for the Resource PKI (RPKI),” work in progress, Internet Draft, **draft-ietf-sidr-cp-04.txt**, November 2008. *A description of the certificate policy that applies to all certificates issued within the RPKI framework.*
- [11] D. Kong, K. Seo, S. Kent, “Template for an Internet Registry’s Certification Practice Statement (CPS) for the Resource PKI (RPKI),” work in progress, Internet Draft, **draft-ietf-sidr-cps-irs-04.txt**, November 2008. *A template for the Practice Statement used by Internet Registries (IRs) to describe their operational practices in the issuance and management of Resource Certificates.*
- [12] D. Kong, K. Seo, S. Kent, “Template for an Internet Service Provider’s Certification Practice Statement (CPS) for the Resource PKI (RPKI),” work in progress, Internet Draft, **draft-ietf-sidr-cps-isp-03.txt**, November 2008. *A template for the Practice Statement used by ISPs to describe their operational practices in the issuance and management of Resource Certificates.*

Individual Submissions:

- [13] G. Huston, G. Michaelson, “A Profile for AS Adjacency Attestation Objects,” work in progress, Internet Draft, **draft-huston-sidr-aao-profile-00.txt**, September 2008. *The specification of the syntax for a pairwise inter-AS routing adjacency attestation.*
- [14] R. Kisteleki, J. Boumans, “Securing RPSL Objects with RPKI Signatures,” work in progress, Internet Draft, **draft-kisteleki-sidr-rpsl-sig-00.txt**, October 2008. *The specification of the addition of RPKI digital signatures to RPSL Objects in the context of an Internet Route Registry.*
- [15] T. Manderson, G. Michaelson, “RPKI Repository Retrieval Mechanism,” work in progress, Internet Draft, **draft-manderson-sidr-fetch-00**, October 2008. *A proposed mechanism to use the manifest as the basis for performing a synchronization operation between a local RPKI cache and a source point.*

RFCs:

- [16] D. Cooper et al., “Internet X.509 Public Key Infrastructure Certificate and Certificate Revocation List (CRL) Profile,” RFC 5280, May 2008.
- [17] C. Lynn, S. Kent and K. Seo, “X.509 Extensions for IP Addresses and AS Identifiers,” RFC 3779, June 2004.
- [18] R. Housley, “Cryptographic Message Syntax (CMS),” RFC 3852, July 2004.
- [19] C. Alaettinoglu, et al., “Routing Policy Specification Language (RPSL),” RFC 2622, June 1999.
- [20] C. Villamizar et al., “Routing Policy System Security,” RFC 2725, December 1999.

Other Documents:

- [21] Kent, S., “Securing BGP: S-BGP,” *The Internet Protocol Journal*, Volume 6, No. 3, September 2003.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. The author of numerous Internet-related books, he is currently the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He was a member of the Internet Architecture Board from 1999 until 2005, and served on the Board of the Internet Society from 1992 until 2001. E-mail: gih@apnic.net

Host Identity Protocol: Identifier/Locator Split for Host Mobility and Multihoming

by Andrei Gurtov and Miika Komu, Helsinki Institute for Information Technology,
and Robert Moskowitz, ICSAlab

A host and its location are identified using *Internet Protocol* (IP) addresses in the current Internet architecture. However, IP addresses can serve only as short-term identifiers because a considerable amount of hosts are *portable* devices and they change their IP addresses when moved from one network to another. Short-term identifiers disrupt long-term transport layer connections, such as Internet phone calls, and make locating the peer host more difficult. Therefore, mobility and multihoming are hard to implement securely in the present Internet. Upon changing an IP address, the host must prove to its peers that it is the same entity they communicated with before, requiring the use of cryptographic identities.

Another challenge the Internet faces is due to the fact that deployed protocols in the Internet are prone to *Denial-of-Service* (DoS) attacks. Substantial memory state can be created before the communicating peer is authenticated. Impersonation attacks are possible because IP addresses are relatively easy to forge. Because of difficulties in configuring *IP Security* (IPsec) for users, most Internet traffic is still transmitted in plaintext, making it easy for attackers to collect passwords or lists of visited websites, for example, in public *Wireless Local-Area Networks* (WLANs). As the IPv6 protocol is seeing gradual deployment, interoperating traditional IPv4 applications with new IPv6 applications remain a challenge.

The so-called *identifier/locator split* is recognized by the *Internet Engineering Task Force* (IETF) community as a next big change in the Internet architecture. Although the problem has been known for a long time^[17], it has only recently started to get sufficient attention. Developments in public key cryptography and increased computational resources of hosts enables the use of cryptographic mechanisms to securely handle identities. Several proposals are under consideration in the IETF, including the *Locator Identifier Separation Protocol* (LISP)^[16] for the network-based and the *Host Identity Protocol* (HIP) for the host-based approach. LISP focuses on improving scalability of the routing system, whereas HIP provides secure end-to-end mobility and multihoming. Therefore, the two proposals are complementary rather than competing.

HIP Architecture

The HIP architecture^[1,2] uses the identity/locator split advantage to address Internet architecture challenges in an integrated approach. HIP was proposed by Bob Moskowitz in 1999 and since then has been under active development in the IETF Working Group and *Internet Research Task Force* (IRTF) Research Group.

HIP enables host mobility and multihoming across different address families (IPv4 and IPv6), offers end-to-end encryption and protection against certain DoS attacks, allows moving away from IP address-based access control to permanent host identities, and restores end-to-end host identification in the presence of several addressing domains separated by *Network Address Translation* (NAT) devices.

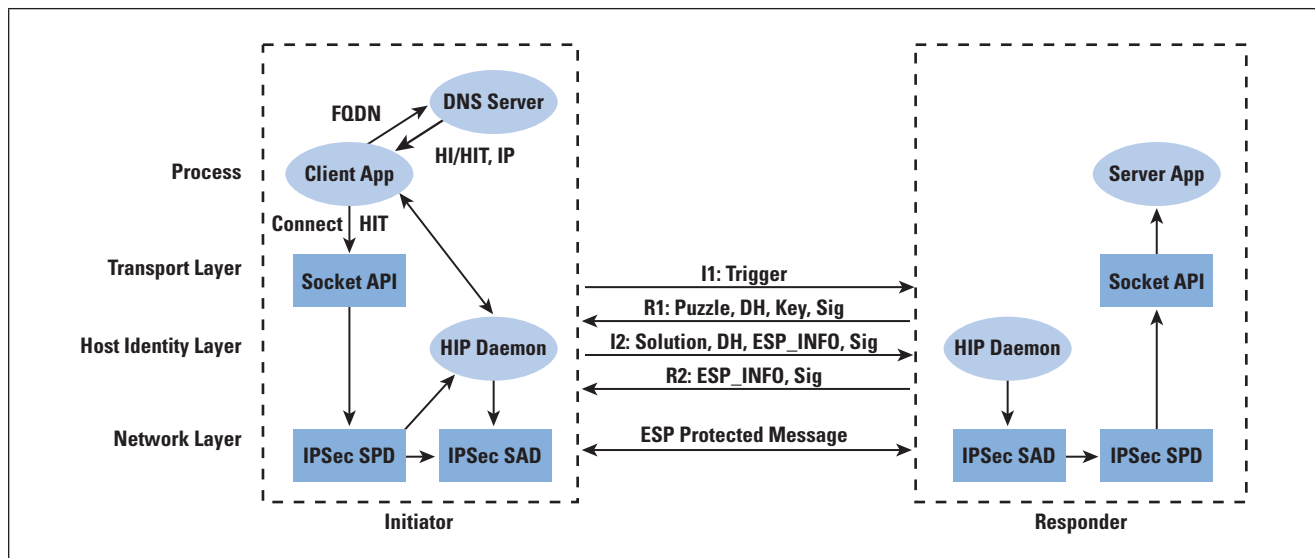
HIP separates the identity of a host from its location. The location of the host is bound to IP addresses and used for routing packets to the host in the same way as in the current Internet architecture. However, transport and application layers use *host identity*, consisting of the public key component of a private-public key pair. Each host is responsible for creating one or more public/private key pairs to provide identities for itself. Because the host identities are based on public key cryptography, they are computationally difficult to forge. Host identities are location-independent identifiers that allow a mobile host to preserve its transport layer connections upon movement. On the other hand, the host identity can be used for looking up the current location of a host because the host identity is a long-term identifier. A client host obtains the host identity of a server typically from the *Domain Name System* (DNS)^[7] or a *Distributed Hash Table* (DHT). However, the infrastructure may not support this DHT in certain scenarios, such as in peer-to-peer and temporary environments. In such cases, *opportunistic* HIP can be used for contacting a peer without prior information of the identity of the peer. Opportunistic HIP is based on a “leap-of-faith,” meaning that it is prone to man-in-the-middle attacks for the initial connection. It is similar to the *Secure Shell* (SSH) *Protocol*, where the public key of the server is added to the known host list after the first connection.

The problem of certifying the keys in *Public Key Infrastructure* (PKI) or otherwise creating trust relationships between hosts has explicitly been left out of the HIP architecture, because it is expected that each system using HIP may want to address it differently. For mere mobility and multihoming, the systems can work without any explicit trust management, in an opportunistic manner.

All other parties use the host identifier, that is, the *public key*, to identify and authenticate the host. Typically, a host identifier is a 128-bit-long bit string, the *Host Identity Tag* (HIT), as shown in Figure 1. A HIT is constructed by applying a cryptographic hash function over the public key. The introduction of new endpoint identifiers changes the role of IP addresses. When HIP is used, IP addresses become pure topological labels, naming locations in the Internet. One benefit of this identity/locator separation is that hosts in private address realms (behind NATs) can name each other in a unique way with HITs. A second benefit is that the hosts can change their IP address without breaking transport layer connections of applications and rely on HIP to manage host mobility; the relationship between location names and identifiers becomes dynamic.

To start communicating through HIP, two hosts must establish a HIP association. Known as the *HIP Base Exchange (BEX)*^[3], this process consists of four messages (I1, R1, I2, and R2) transferred between the initiator and the responder. After BEX is successfully completed, both hosts are confident that private keys corresponding to host identifiers (public keys) are indeed possessed by their peers. Another purpose of the HIP base exchange is to create a pair of *IPsec Encapsulated Security Payload (ESP) Security Associations (SAs)*, one for each direction. HIP uses *IPsec ESP Bound End-to-End Mode (BEET)*^[4,9] to provide data encryption and integrity protection for network applications.

Figure 1: HIP Architecture



Because neither transport layer connections nor security associations created after the HIP base exchange are bound to IP addresses, a mobile client can change its IP address (that is, upon moving, because of a *Dynamic Host Configuration Protocol* [DHCP] lease or IPv6 router advertisement) and continue to transmit ESP-protected packets to its peer. HIP supports such mobility events by implementing an end-to-end three-way UPDATE signaling mechanism^[8] between communicating nodes. HIP multihoming uses the same mechanisms as mobility for updating the peer with a current set of host IP addresses.

A rendezvous server^[6] provides a mechanism to locate a host, for example, when two communicating hosts move simultaneously. To employ a rendezvous mechanism, a host first must perform a registration procedure^[5], which is an extended version of the HIP base exchange.

The HIP control packets as well as ESP-encapsulated data packets have difficulties in going through NAT applications and firewalls. To traverse NAT, HIP uses *User Datagram Protocol* (UDP)-based encapsulation provided by the *Interactive Connectivity Establishment* (ICE) protocol.

It enables two hosts located behind NAT to communicate through a Rendezvous server. Bob Moskowitz suggests an alternative approach, where HIP always uses IPv6 for end-to-end communication and the *Teredo* protocol is employed to traverse NAT instances in IPv4 networks if native IPv6 connectivity is not available.

Most Internet applications can run unmodified over HIP^[10], although only HIP-aware (new) applications using the extended socket interface can take better advantage of the new features that HIP provides. As HIP secures application data traffic with IPsec that is located logically “deep” within the networking stack, the challenge is to provide proper and understandable security indicators to the user to convince the user that the connection, for example, to a banking website, is secured. Such indicators can be developed as extensions to applications (for example, a security plug-in to the *Firefox* browser) or within a hostwide HIP management utility that controls all applications.

HIP provides a network layer alternative to using *Secure Sockets Layer/Transport Layer Security* (SSL/TLS) for application security, which has its benefits and drawbacks. HIP is a generic solution that should work for any transport protocol, whereas until recently TLS supported only TCP. HIP enables host mobility and multihoming, which is not supported by TLS. TLS runs on top of TCP, leaving it vulnerable to various TCP attacks; for example, using spoofed *reset* (RST) packets or DoS attacks with SYNs. Applications must be designed explicitly to use TLS, whereas HIP can provide security as an add-on to existing traditional applications. On the other hand, TLS does not have a problem with traversing traditional middle-boxes such as NATs and firewalls that need special attention for HIP. Both protocols share the characteristic of endorsing host identity. TLS relies on certificates issued by one of the known Certification Authorities, whereas HIP can use *Domain Name System Security Extensions* (DNSSEC)^[18] or a PKI infrastructure.

There are currently three open-source interoperating HIP implementations. *OpenHIP* from Boeing runs on Linux, Windows, and Mac OS, whereas *HIP on Linux* (HIPL) runs on Linux and Symbian, and *HIP for Inter.net* from Ericsson runs on FreeBSD and Linux. Several testbeds are deployed based on HIP, including the Everett Boeing factory^[11], the P2PSIP pilot in Finland^[14], and Wi-Fi P2P Internet Sharing Architecture in Germany^[12]. Ericsson NomadicLab and TeliaSonera have demonstrated using HIP for transparent IPv4 and IPv6 handovers, mobile router, simultaneous multiaccess, and the use of proxy for traditional hosts^[13,15].

Acknowledgements

We are grateful to Pekka Nikander, Tom Henderson, and others in the IETF and the *Internet Research Task Force* (IRTF) community who were encouraging and contributing to the development of HIP. We thank Andrey Khurri for the figure on HIP architecture and Henry Sinnreich for encouraging us to write this article.

We also thank members of InfraHIP II project for comments helping to improve this article.

References

- [1] Moskowitz, R. and Nikander, P., “Host Identity Protocol Architecture,” RFC 4423, May 2006.
- [2] Gurtov, A., *Host Identity Protocol (HIP): Towards the Secure Mobile Internet*, ISBN 978-0-470-99790-1, Wiley and Sons, June 2008.
- [3] Moskowitz, R., Nikander, P., Jokela, P. and Henderson, T., “Host Identity Protocol,” RFC 5201, April 2008.
- [4] Jokela, P., Moskowitz, R. and Nikander, P., “Using the Encapsulating Security Payload (ESP) Transport Format with the Host Identity Protocol (HIP),” RFC 5202, April 2008.
- [5] Laganier, J., Koponen, T. and Eggert, L., “Host Identity Protocol (HIP) Registration Extension,” RFC 5203, April 2008.
- [6] Laganier, J. and Eggert, L., “Host Identity Protocol (HIP) Rendezvous Extension,” RFC 5204, April 2008.
- [7] Nikander, P. and Laganier, J., “Host Identity Protocol (HIP) Domain Name System (DNS) Extension,” RFC 5205, April 2008.
- [8] Nikander, P., Henderson, T., Vogt, C. and Arkko, J. “End-host Mobility and Multihoming with the Host Identity Protocol,” RFC 5206, April 2008.
- [9] Nikander, P. and Melen, J., “A Bound End-to-End Tunnel (BEET) Mode for ESP,” Internet Draft, Work in Progress, **draft-nikander-esp-beet-mode-09**
- [10] Henderson, T., Nikander, P. and Komu, M., “Using the Host Identity Protocol with Legacy Applications,” RFC 5338, September 2008.
- [11] Boeing, “Secure Mobile Architecture (SMA) for Automation Security,” http://www.isa.org/wsummit/presentations/Boeing-NGI_SMA_Automation_Security_Vancouver_ISA_presentationtemplates_7-23-07.ppt
- [12] Heer, T., Götz, S., Weingärtner, E. and Wehrle, K., “Secure Wi-Fi Sharing on Global Scales,” in Proceedings of the 15th International Conference on Telecommunication (ICT), St. Petersburg, Russian Federation, IEEE, 2008.
<https://www.ds-group.info/members/heer/publications-tobias-heer/pdfs/HeerEtAl2008.pdf>

- [13] Jokela, P., Ylitalo, J., and Salmela, P., “HIP Mobile Router Demo,” March 2007.
<http://www.ietf.org/proceedings/07mar/slides/HIPRG-3.pdf>
- [14] Koskela, J., Heikkila, J. and Gurtov, A., “A Secure P2PSIP System with SPAM Prevention,” Poster at ACM Mobicom, September 2008.
- [15] Korhonen, J., Mäkelä, A., and Rinta-aho, T., “HIP Based Network Access Protocol in Operator Network Deployments,” in First Ambient Networks Workshop on Mobility, Multiaccess, and Network Management (M2NM’07), Sydney, Australia, October 2007.
- [16] Meyer, D., “The Locator Identifier Separation Protocol (LISP),” *The Internet Protocol Journal*, Volume 11, No. 1, March 2008.
- [17] Saltzer J., “On The Naming and Binding of Network Destinations,” RFC 1498, September 1992.
- [18] Gieben, M., “DNSSEC: The Protocol, Deployment, and a Bit of Development,” *The Internet Protocol Journal*, Volume 7, No. 2, June 2004.
- [19] Sinnreich, H., “Letter to the Editor,” *The Internet Protocol Journal*, Volume 11, No. 3, page 37, September 2008.

ANDREI GURTOV received M.Sc and Ph.D. degrees in Computer Science from the University of Helsinki, Finland. He presently is Principal Scientist, leading the Networking Research group at the Helsinki Institute for Information Technology, focusing on distributed system security and next-generation Internet architecture. He co-chairs the IRTF research group on HIP and teaches as an adjunct professor at Helsinki University of Technology. He is a regular visitor of the ICSI Center for Internet Research (ICIR) at Berkeley. Andrei has co-authored more than 50 publications, including a book, research papers, patents, and RFCs. He can be reached through the webpage: <http://www.hiit.fi/~gurtov>

MIIKA KOMU received his M.Sc. from Helsinki University of Technology and continues his studies as a postgraduate student. He is working as a full-time researcher and software engineer at Helsinki Institute for Information Technology. He is an active IETF participant and co-author of RFC 5338. Miika is an open source advocate and martial arts fan. E-mail: miika.komu@hiit.fi

ROBERT MOSKOWITZ is senior technical director for ICSA Labs and is an active member in the IAB, IETF, and IEEE. At ICSA Labs, Moskowitz leads the IPsec product and system certification program. Prior to the ICSA, he led the adoption of the world’s largest IPsec network deployment servicing the automotive industry. As a former co-chair of the IPsec Working Group, Moskowitz provided a user set of multivendor, multipolicy, and multiuser requirements that galvanized many of the debates on the use of IPsec. A contributing editor for *Network Computing Magazine*, Moskowitz is currently helping define the new security component for the 802.11 standard. E-mail: rgm@htt-consult.com

Allocation Policy for the Remaining IPv4 Address Space Ratified by ICANN

On 6 March 2009, the *International Corporation for Assigned Names and Numbers* (ICANN) Board ratified the *Global Policy for the Allocation of the Remaining IPv4 Address Space*. The policy requires ICANN to reserve one /8 for each *Regional Internet Registry* (RIR) from the *Internet Assigned Numbers Authority* (IANA) free pool. This has been done. The remainder of the implementation will be done once the IANA free pool has been fully allocated to RIRs. There are currently 32 unallocated unicast IPv4 /8s. 27 are in the IANA free pool and five are reserved under the Global Policy for the Allocation of the Remaining IPv4 Address Space.

On 4 February 2009, the Chair of the *Address Supporting Organization Address Council* (ASO AC) forwarded the Proposed Global Policy for the Allocation of the Remaining IPv4 Address Space for ratification by the ICANN Board. On 5 March 2009, the ASO AC submitted advice in full support of the proposal to the ICANN Board. This proposed global policy had been submitted to the ASO AC by the Executive Council of the *Number Resource Organization* (NRO) on 3 December 2008, and adopted by the ASO AC on 8 January 2009. Each RIR community individually discussed the policy and approved its adoption via its own policy development process. The policy text is published on the ICANN web site at:

<http://www.icann.org/en/general/allocation-remaining-ipv4-space.htm>

ISOC's Trust and Identity Initiative

The Internet Society's *Trust and Identity Initiative* recognizes that in order to be trusted, the Internet must provide channels for secure, reliable, private, communication between entities, which can be clearly authenticated in a mutually understood manner. The mechanisms that provide this level of assurance must support both the end-to-end nature of Internet architecture and reasonable means for entities to manage and protect their own identity details.

A *trusted* Internet takes into account security, transaction protection, and identity assertion and management. Given the network dependence on unique numbers and the escalating amount of geolocation data being gathered, the privacy implications of the current Internet represent a significant and growing concern. Trust must be a primary design element at every layer of the architecture, and in some cases, existing elements may need to be redesigned or improved to meet emerging requirements.

In late 2007, the ISOC Board of Trustees held an intensive retreat to consider ISOC's role in identifying and pursuing trust and identity issues. The report arising from that meeting, "Trust and the Future of the Internet,"^[1] forms the basis of ISOC's current long term strategic initiative.

The Trust and Identity initiative focuses on the following major research programs:

- *Architecture and Trust:* This research program investigates the implementation of open-trust mechanisms throughout the full cycle of Internet research, standardization, development, and deployment.
- *Current Problems and Solutions and Trust:* This research program investigates the mitigation of the social, policy, and economic factors that may hinder development and deployment for trust-enabling technologies.
- *Identity and Trust:* This research program investigates the elevation of identity to a core issue in network research and standards development. ISOC is taking a lead role in reviewing the current Internet architecture and the model of Internet development and deployment. This includes active engagement with participants within the traditional ISOC sphere, as well as with the research, enterprise, and end-user communities. We offer the kind of support for research that enhances and facilitates trust and collaboration with the standards community and that advances the most interesting outcomes of that research.

ISOC is reaching out to the businesses and end users that rely on the Internet to exchange sensitive data. Their needs and concerns inform both our baseline research agendas and ongoing standards and development work. ISOC continues to support the advancement of current technical solutions and best practices through our existing programs.

[1] "Trust and the Future of the Internet,"

<http://www.isoc.org/isoc/mission/initiative/docs/trust-report-2008.pdf>

[2] "Trust and Identity Initiative" brochure,

<http://www.isoc.org/pubs/isoc/docs/trust.pdf>

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ contains standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2009 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

June 2009

Volume 12, Number 2

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
DNS Caching.....	2
IEEE 802.21	7
Book Review.....	28
Fragments	30
Call for Papers.....	31

After many years of using DSL as my only Internet access option from home, I recently upgraded to a broadband solution provided by a cable modem. As a result, I faced the task of renumbering (and partially rewiring) my home network. As you might have guessed, the addressing scheme provided by my new ISP offers *Network Address Translation* (NAT), as well as a small number (5) of fixed IPv4 addresses, the latter at an extra cost as you might expect. I probably should have tried to enable IPv6 just as an experiment, but this task will have to wait for another day. In the meantime, I was pleased to find a relatively user-friendly web interface to the cable modem that allows me to configure numerous parameters, including the range of the *Dynamic Host Configuration Protocol* (DHCP) pool so that certain devices (printers and wireless access points in particular) can have fixed IP addresses for ease of use and configuration. The entire exercise, which took a couple of hours on my very small network, reminded me of what network managers face every day, particularly as they consider the inevitable migration to IPv6. Let me take this opportunity to invite you to share your network management and operations experience, plans for IPv6 migration, and so on. You can send us Letters to the Editor or article proposals. The address, as always, is ipj@cisco.com

The *Domain Name System* (DNS) has been the target of attacks over its many years of existence. In recent years, new attacks have emerged that exploit some of the attributes of the DNS protocol and its implementation. One of the corrective measures is to improve the security of DNS caches. There are several ways to improve cache security, most of which involve changing the protocol. Another way, without changing the protocol, is to reduce the attack surface of your cache by shrinking the number of users of any given cache. Our first article, by Bill Manning, explores this view in more detail.

This journal has covered numerous current and emerging *wireless* technologies such as Bluetooth, Wi-Fi, WiMAX, and mobile cellular systems. In this issue, Esa Piri and Kostas Pentikousis describe *Media-Independent Handovers* (MIH), which allow mobile devices to use different wireless and wired network infrastructures transparently. The protocols associated with operation across such diverse access networks are being standardized by the IEEE 802.21 working group.

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

Intermediate DNS Caching as an Attack Vector

by Bill Manning

The *Domain Name System* (DNS) specification calls for the use of *caching*. Caching is expected to improve the overall responsiveness of the system by ensuring that answers to questions are known and stored locally and that the query load placed on the authoritative servers is minimized. Certain presumptions are associated with caches that may no longer hold. This article looks at some of these presumptions and explores some of the problems that emerge when they are violated. Based on our observations, we offer some recommendations on DNS cache best practices and show our results of testing these practices.

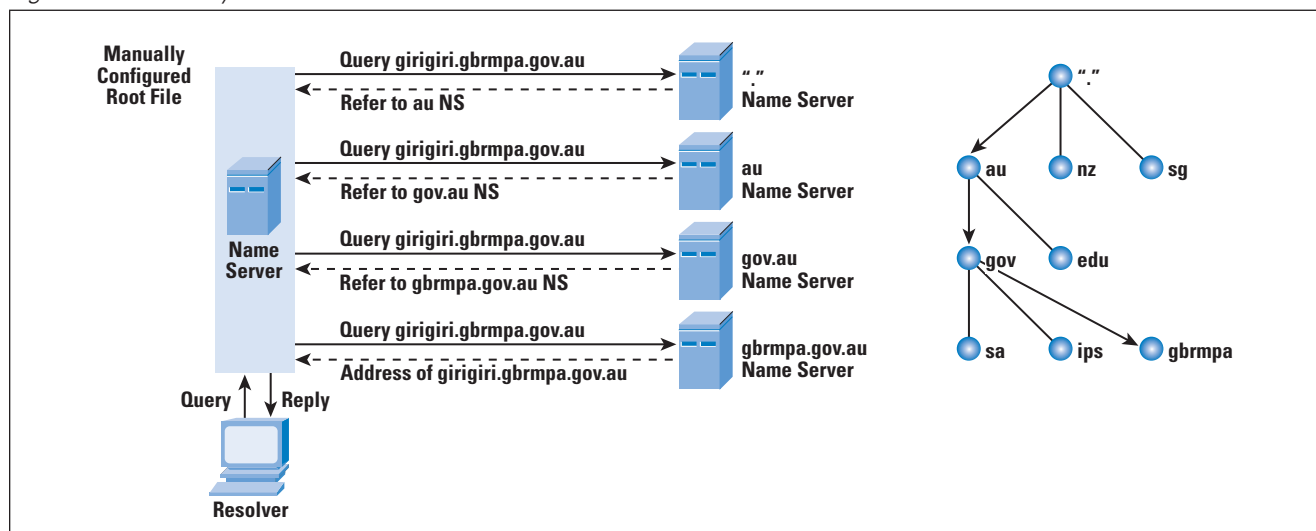
The Problem

A DNS resolver can no longer trust the data it gets—because the data generally comes from nonauthoritative nodes or caches operated by third parties, most of whom have no vested interest in providing accurate data. Removing or bypassing caching from the DNS and going directly to the authoritative servers is considered a fatal flaw because authoritative servers are presumed to have neither the bandwidth nor the processing power to accommodate the perceived demand from a cacheless service. This article looks at the bandwidth and processing capabilities of modern authoritative servers to ascertain the viability of these presumptions. We start by looking briefly at the DNS.

The DNS

The DNS namespace is made visible and useful by nodes publishing authoritative information about the namespace and *resolvers* that send queries about the namespace to these servers. As an optimization, other nodes may act as intermediates or proxies for the authoritative servers for one to many resolvers. These intermediate nodes are called *caching nameservers* or *iterative mode resolvers*. This flow is shown in Figure 1.

Figure 1: DNS Query Flow



Several assumptions about the use and placement of caches have been questioned recently. The simplest is one of placement. A cache works best when the *Round-Trip Time* (RTT) between the resolver and the cache is low. Historically, a cache was placed at traffic aggregation points such as an *Internet Service Provider* (ISP) operating a cache for its clients. With increased mobility of nodes, this presumption is no longer as firm. There are reported cases where resolvers continue to use caches 300 ms away, while an authoritative server is 15 ms away. So if the intent is to reduce network bandwidth, then a cache presuming its client resolvers are all “local” might be misconstrued.

Fixing a resolver to a specific cache does have the benefit of being tied to a known business relationship; for example, using your ISP’s caching service. In contrast, mobile nodes often get an IP address from a provider’s *Dynamic Host Configuration Protocol* (DHCP) servers, which also hand out more “local” caching servers to be used by the mobile node.

This scenario would be fine—as long as the DNS namespace was in fact a coherent, single space. Unfortunately it is not. So-called *Walled-Garden* networks that have their own versions of DNS namespace have been and remain common. In the Internet, there are more and more alternate root hierarchies that diverge from what most think of as “the” root namespace in either subtle or wildly divergent ways. To date, there is no deployed way for a resolver to determine the origin of the data stored in a cache. A resolver then has no way other than verification of the data to know that the locally assigned cache is in fact using the namespace desired. This situation represents one important reason for going back to a well-known cache, even if it is topologically remote. But this assumption may no longer be valid.

ISPs and even some caching service providers are starting to manipulate caches as a means to monetize their operations.^[1] Numerous techniques are in use, from the nominally benign method of using wildcards to more insidious capture and rewrite of NXDOMAIN replies, to outright intentional cache pollution.

In this climate, a resolver should choose its cache carefully. We argue that it is reasonable, in many of today’s environments, to place the cache within 1 ms of the resolver; for example, run a cache on the local node. This argument is an extension of the assertion^[2] that claims that caches are effective for client populations that are about 10 or fewer.

This technique has the added advantage of reducing the “attack surface” by reducing the effect of cache poisoning or rewriting replies to a small handful of nodes. The perceived disadvantage is the increased load on network bandwidth and query load on authoritative servers as the number of caches increases.

The Experiment

Our experiment has two parts: first we looked at authoritative server processing capabilities and then at the bandwidth effects of a larger number of caches.

Authoritative service is generally run on systems with modern software, supporting threading or precomputed responses. Independent testing shows that these stock software solutions can, on current hardware, support query rates in the hundreds of thousands of queries per second.^[3]

A brief survey of authoritative server operators indicates that normal query rates range from 12,000 to 64,000 queries per second.^[4,5,6]

On the surface, this result would indicate that there is enough overhead to be able to process more queries, regardless of how they are originated. Regarding bandwidth, a survey of *Top-Level Domain* (TLD) operators has shown that 92 percent of the delegations have two or more authoritative servers for that data on networks with a minimum uplink bandwidth of 100 Mbps. Selected path characterization from clients to target authoritative servers seems to support our presumption that bandwidth is not of concern.

The DNS was designed to function as a roughly symmetrical transfer of information: a request or query is sent and the reply reflects the query and supplies the answer and additional data. Historically, the request and reply were within the same order of magnitude. Into the future, this model may no longer be valid. With *Domain Name System Security Extensions* (DNSSEC), *IP Version 6* (IPv6), and *Naming Authority Pointer* (NAPTR) records being possible candidates in the *Resource Record set* (RRset), the traffic profile more resembles an HTTP request/response, with a significant amount of data being returned from a simple question.^[7]

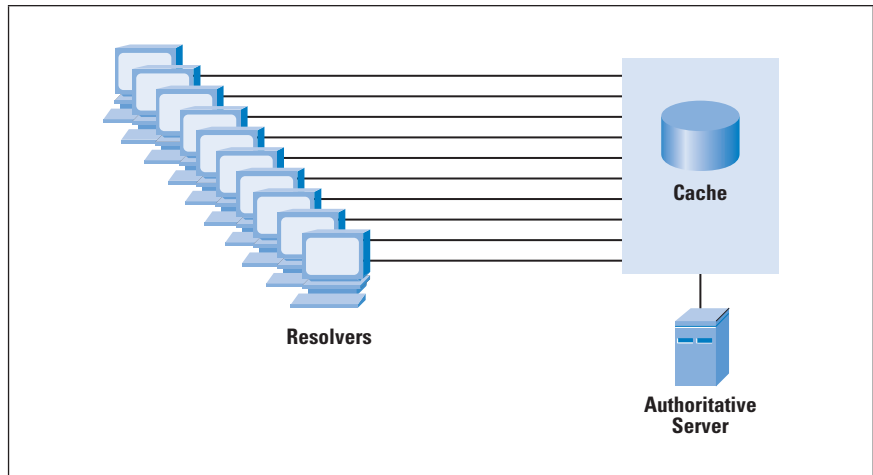
With this information, we can project a worse case in today's environment where a query/reply is about 260 bytes to a worst case in a future environment where a query/reply is about 9 KB, clearly indicating that the amount of bandwidth to authoritative servers needs to grow as new DNS capabilities are deployed, but for the nonce, most have a bandwidth overhead sufficient to absorb a modest change in the number of queries presented.

Modification of the Number of Caching Servers

We began with a cache that serviced 140 stub resolvers on the *University of Southern California's Information Sciences Institute* (USC/ISI) campus in a "normal" dense cache mode (Figure 2).

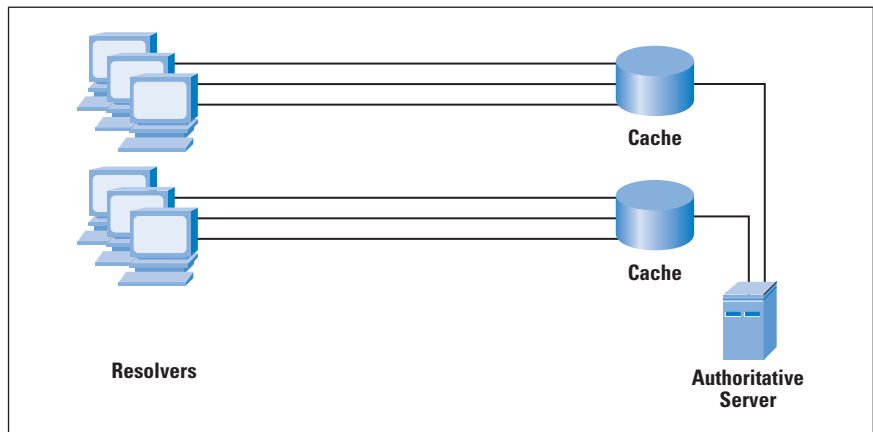
Traffic traces show a distribution of priming queries to 534 authoritative servers in the first 15 minutes of clearing the cache.

Figure 2: Dense Cache



We then added 9 new caches and redistributed the 140 stub resolvers among the 10 caches into a sparse cache mode (Figure 3) and restarted all the caches. In the first 15 minutes, the number of priming queries from each of the caches averaged 61, with a total of 622 unique priming queries for all caches. The number of “duplicate” queries between caches averaged 45. Although the number of queries to the authoritative servers was slightly higher, the results seem to indicate that there is a small but significant difference in each of the caches^[8].

Figure 3: Sparse Cache



Conclusions

Reducing the size of the user population for each cache reduces the attack surface for the DNS overall because we have effectively compartmentalized the threat to a small number of nodes. Generally, restarting a cache for a small number of nodes is considered acceptable, whereas restarting a cache for 10,000 or 100,000 nodes would significantly affect operations.

Moving the cache closer to the resolver increases overall response time and may support better mobility of the node. If validation is also placed with the cache, it is possible to increase the confidence of validation because that information may not have to use DNS protocols to send validation data over untrusted, open networks.

The concept of supporting larger numbers of full DNS servers on more nodes raises concerns, but most systems these days have enough processing power and bandwidth to support this application. Administrative and management processes can be fully automated. Overall, this design complements other, protocol-based attempts to increase DNS integrity.

References

- [1] “Preliminary Report on DNS Response Modification,” 20 June 2008, <http://www.icann.org/en/committees/security/sac032.pdf>
- [2] “DNS Performance and the Effectiveness of Caching,” Jaeyeon Jung, Emil Sit, Hari Balakrishnan, and Robert Morris, *IEEE/ACM Transactions on Networking*, Volume 10, No. 5, pp. 589–603, October 2002.
- [3] <https://www.dns-oarc.net/files/workshop-2006/Dickinson-Performance.pdf>
- [4] “An analysis of Wide-Area Name Server Traffic: A Study of the Internet Domain Name System,” Peter B. Danzig, Katia Obraczka, and Anant Kumar, *ACM SIGCOMM Computer Communications Review*, Volume 22, No. 4, pp. 281–292, 1992.
- [5] “An Analysis of the Queries from Caching Servers to Root Servers, Tsuyoshi Toyono, NTT Laboratories, 2007 OARC Workshop, <https://www.dns-oarc.net/files/dnsops-2007/Toyono-Caching-analysis.pdf>
- [6] RootServer supplied statistics:
<http://h.root-servers.org/>
<http://k.root-servers.org/index.html#stats>
<http://m.root-servers.org/>
- [7] http://snad.ncsl.nist.gov/dnssec/mem_usage.html
<http://snad.ncsl.nist.gov/dnssec/bandwidth.html>
- [8] “Sharp Transition Towards Shared Vocabularies in Multi-agent Systems,” Andrea Baronchelli, Maddalena Felici, Vittorio Loreto, Emanuele Caglioti, and Luc Steels, *Journal of Statistical Mechanics: Theory and Experiment*, 2006, P06014.
<http://www.iop.org/EJ/abstract/1742-5468/2006/06/P06014>

BILL MANNING has been in the network field since 1979, currently with the Keio University, Shonan Fujisawa Campus, and USC/ISI. He has been an IETF Working Group chair and RFC author, and he currently serves on numerous ICANN committees. He is part of the team that runs one of the Internet Root nameservers. E-mail: bmanning@sfc.wide.ad.jp

IEEE 802.21: Media-Independent Handover Services

by Esa Piri and Kostas Pentikousis, VTT Technical Research Centre of Finland

Popular mobile devices now ship with several integrated wired and wireless network interfaces. *Personal Digital Assistants* (PDAs) and smartphones, for example, are increasingly supporting communications through both cellular technologies and *Wireless LANs* (WLANs); laptops typically come with built-in Ethernet, Wi-Fi, and Bluetooth^[1]. As multiaccess devices proliferate, we move closer to a network environment that is often referred to as “*beyond 3G*” (B3G). Key success factors for cellular *third-generation* (3G) communications include better cell capacities, increased data rates, transparent mobility within large geographical areas, and global reachability. For B3G, the next frontier lies beyond transparent mobile connections within the same access technology because users will expect to be globally reachable anytime, anywhere, and remain “*always best-connected*” (ABC)^[2]. In order to select the best possible connectivity option (anytime, anywhere), mobile devices and access networks will have to work together in order to enable users to take full advantage of all available options.

The IEEE 802.21 working group (see www.ieee802.org/21) recently finalized the first standard for dealing with handovers in heterogeneous networks, also called *Media-Independent Handovers* (MIH)^[3]. The standard is expected to allow mobile users (and operators) to take full advantage of overlapping and diverse access networks. It provides a framework for efficiently discovering networks in range and executing intelligent heterogeneous handovers, based on their respective capabilities and current link conditions. This article aims to serve as a primer for those interested in the IEEE 802.21 standard. After introducing the IEEE 802.21 reference model, we present the MIH services and provide illustrative use cases that highlight the benefits of employing the Media-Independent Handover Services standard in heterogeneous networks.

Mobile and Wireless

The widespread success of 3G technologies^[4, 5] is evidenced by the rapid increase in the amount of data traffic over cellular networks in recent years. In Sweden, for example, the total amount of mobile data traffic leapt tenfold from just over 203 TB in 2006 to 2191 TB in 2007^[6]. This trend is expected to continue unabated with the deployment of *High-Speed Packet Access* (HSPA) and *Long-Term Evolution* (LTE) in the coming years. Of course, the amount of traffic over cellular networks is only a proportion of the traffic that originates from or terminates at WLANs worldwide. Campuswide deployments of WLANs are becoming the norm in developed countries, and we even find citywide WLANs, as in the case of the city of Oulu, Finland (see www.panoulu.net).

Finally, many anticipate that mobile WiMAX^[7] deployments will significantly affect telecommunications markets. In short, we are moving toward a far more heterogeneous network access environment than the one users and operators face today, with multiple overlapping mobile and wireless networks with diverse characteristics.

Multiaccess Devices in Heterogeneous Networks

As communication environments become more complex because of the diversity of network access technologies that support, for example, different access rates and *Quality of Service* (QoS) levels, users expect more from their wireless operator. Mobile devices, once featuring tiny screens, extremely limited processing and storage capacities, and narrowband connectivity^[8], now pack capabilities that just a few years ago were typical of high-end laptops. This scenario has allowed users to increasingly depend on mobile devices for e-mail and *Instant Messaging* (IM), but also for making *Voice over IP* (VoIP) calls, listening to streaming Internet radio, and watching online videos.

With respect to user mobility patterns, campuswide Wi-Fi users typically spend most of their connection time attached to a small set of access points located within a small radius^[9, 10]. This situation is not surprising, because Wi-Fi was originally designed and subsequently deployed mainly as an extension to wired infrastructures. In the future, however, we anticipate that multiaccess devices will employ different network interfaces to attach to different access networks, establishing multiple parallel connections over 3G/*Universal Mobile Telecommunications Service* (UMTS) and Wi-Fi, for example. With global reachability and ABC mechanisms in place, mobile devices will be able to selectively connect to different access networks depending on certain criteria. Keep in mind that from a conventional, IP-centered point of view, changing the *Point of Attachment* (PoA) calls for mobility management actions^[11, 12, 13], although in practice there may be no physical mobility whatsoever.

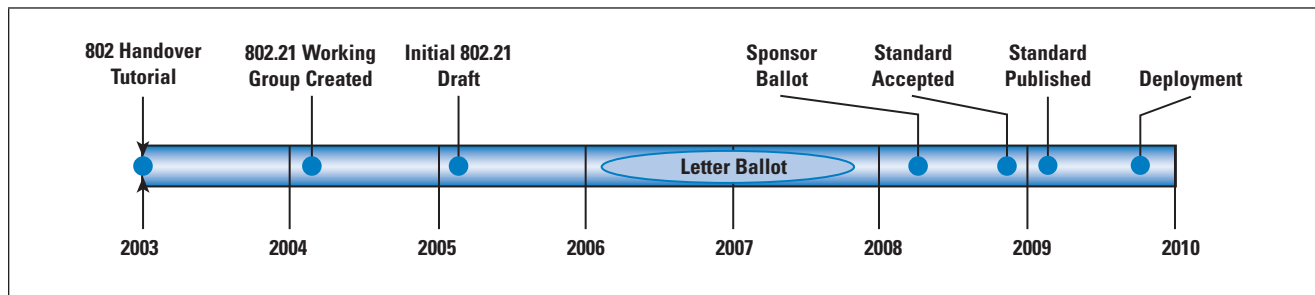
Given the diversity of networked applications running on mobile devices, knowledgeable network resource planning and operation is needed, in turn calling for a framework that allows users and their applications to state their network access preferences. This framework should also allow operators to steer terminal access patterns aiming at maximizing resource usage and increasing user satisfaction. For instance, podcasts can be downloaded only when connected to an uncongested WLAN, but web, map/navigation, and e-mail clients can use the cellular network or WLAN access on demand. Currently, this process can only be done manually: users need to be watchful for available access networks and choose which one to attach to based on very rudimentary information such as signal quality. If mobile nodes could collect timely and consistent information about the state of all available networks in range and were given the means to control their network connectivity, then a whole range of possibilities would become available.

In order to optimize the use of available network resources, mobile nodes need to be able to collect information about numerous heterogeneous networks in a generic and standardized way, irrespective of the underlying network access technology. The collected information, both dynamic and static, can then be used by handover decision-making processes, such as, say, mobility managers. Mobility managers can be enhanced versions of *Mobile IP* (MIP)^[11, 12, 13], proprietary solutions, or other proposals stemming from recent research, such as [14]. Researchers in the area have proposed several cross-layer frameworks for enhancing the efficiency of handover decision makers (see [14, 15] and the references therein), but none of them has been formally standardized or is widely accepted so far. What is needed is a standard framework that can attract ample support from major vendors and operators, and can be deployed incrementally.

Introducing IEEE 802.21-2008

Figure 1 illustrates the progress toward the IEEE 802.21-2008 standard. The working group was initiated in 2004, and the latest draft version of the standard was accepted as a new standard by the IEEE-SA Standards Board in November 2008^[3]. The standard was published in January 2009. It is anticipated that actual deployment of the standard will take place at the earliest in late 2009–2010.

Figure 1: Timeline of the IEEE 802.21-2008 Standardization Effort



IEEE 802.21-2008, also known as *Media-Independent Handover Services*, features a broad set of properties that meet the requirements of effective heterogeneous handovers. It allows for transparent service continuity during handovers by specifying mechanisms to gather and distribute information from various link types to a handover decision maker. The collected information comprises timely and consistent notifications about changes in link conditions and available access networks.

Note that the scope of IEEE 802.21-2008 is restricted to access technology-independent handovers. Intratechnology handovers, handover policies, security mechanisms, media-specific link layer enhancements to support IEEE 802.21-2008, and *Layer 3* (L3) and upper-layer enhancements are outside the scope of IEEE 802.21-2008. This article summarizes the salient points of [3], which henceforth is referred to as IEEE 802.21.

The IEEE 802.21 Reference Model

IEEE 802.21 facilitates a variety of handover methods, including both *hard handovers* and *soft handovers*. A hard handover, also known as “break-before-make” handover, typically implies an abrupt switch between two access points, base stations, or, generally speaking, PoAs. Soft handovers require the establishment of a connection with the target PoA while still routing traffic through the serving PoA. In soft (“make-before-break”) handovers, mobile nodes remain briefly connected with two PoAs. Note, however, that depending on service requirements and application traffic patterns, hard handovers may often go unnoticed. For example, web browsing and audio/video streaming with prebuffering can be accommodated when handing over between different PoAs in the range of one network by employing mechanisms that allow transferring the node connection context from one PoA to another quickly.

The main design elements of IEEE 802.21 can be classified into three categories: a framework for enabling transparent service continuity while handing over between heterogeneous access technologies; a set of handover-enabling functions; and a set of *Service Access Points* (SAPs).

Transparent Service Continuity

IEEE 802.21 specifies a framework that enables transparent service continuity while a mobile node switches between heterogeneous access technologies. The consequences of a particular handover need to be communicated and considered early in the process and, clearly, before the handover execution. In soft handovers, it is crucial that service continuity, during and after the handover, is ensured without any user intervention. To this end, IEEE 802.21 specifies essential mechanisms to gather all necessary information required for an affiliation with a new access point before breaking up the currently used connection. Interactive applications, such as VoIP, are typically the most demanding in terms of handover delays, and high-quality VoIP calls can be served only by soft handovers. On the other hand, video streaming can accommodate hard handovers, as long as the vertical break-before-make handover delay does not exceed the application buffer interval delay. In the case of hard handovers, handover preparation signaling can initiate the connection context transfer from the serving PoA to the target PoA beforehand.

For instance, lack of the required level of QoS support or low available capacity in a candidate access network may lead the network selecting entity to prevent a planned handover. On the other hand, for example, increasing delay, jitter, or packet-loss rates in the currently serving network may degrade the perceived QoS throughout the network, or only for a particular application, triggering the mobility manager to start assessing the potential of candidate target access networks and subsequently initiate an IEEE 802.21-assisted handover.

IEEE 802.21 also allows the reception of dynamic information about the performance of the serving network and other networks in range. In other words, IEEE 802.21 provides methods for continuous monitoring of available access conditions. However, IEEE 802.21 does not specify any methods for collecting this dynamic information at the link layer.

Handover-Enabling Functions

IEEE 802.21 defines a set of handover-enabling functions, which are specified with respect to existing network elements in the protocol stack, and introduces a new logical entity called *Media-Independent Handover Function* (MIHF). The MIHF logically resides between the link layer and the network layer. It provides, among others, abstracted services to entities residing at the network layer and above, called *MIH Users* (MIHUs). MIHUs are anticipated to make handover and link-selection decisions based on their internal policies, context, and the information received from the MIHF. To this end, the primary role of the MIHF is to assist in handovers and handover decision making by providing all necessary information to the network selector or mobility management entities. The latter are responsible for handover decisions regardless of the entity position in the network. The MIHF is not meant to make any decisions with respect to network selection.

Service Access Points

SAPs with associated primitives between the MIHF and MIHUs (MIH_SAP) give MIHUs access to the following services that the MIHF provides:

- The *Media-Independent Event Service* (MIES) provides event reporting about, for example, dynamic changes in link conditions, link status, and link quality. Events can be both local and remote. Remote events are obtained from a peer MIHF entity.
- The *Media-Independent Command Service* (MICS) enables MIHUs to manage and control the parameters related to link behavior and handovers. MICS provides a set of commands for accomplishing that, as we will see later in this article. Commands can be both local and remote. The information obtained with MICS is dynamic.
- The *Media-Independent Information Service* (MIIS) allows MIHUs to receive static information about the characteristics and services of the serving network and other available networks in range. This information can be used to assist in making a decision about which handover target to choose and to make preliminary preparations for a handover.

Figure 2 illustrates the general reference model of IEEE 802.21. The scope of IEEE 802.21 includes only the operation of MIHF and the primitives associated with the interfaces between MIHF and other entities. A single media-independent interface between MIHF and MIHU (MIH_SAP) is sufficient.

On the other hand, there is a need for defining a separate technology-dependent interface, which is specific to the corresponding media type supported, between the MIHF and the lower layers (MIH_LINK_SAP).

The primitives associated with the MIH_LINK_SAP enable MIHF to receive timely and consistent link information and control link operation during handovers. For example, the currently supported link layers include wired and wireless media types from the IEEE family of standards (for example, 802.3, 802.11, 802.15, and 802.16), as well as those defined by the *Third-Generation Partnership Project* (3GPP) and *Third-Generation Partnership Project 2* (3GPP2). Besides these, IEEE 802.21 specifies a media-independent SAP (MIH_NET_SAP), which provides transport services for Layer 2 (L2) and Layer 3 (L3) MIH message exchange with remote MIHFs. Functions over the LLC_SAP are not specified in IEEE 802.21.

Figure 2: The IEEE 802.21-2008 Reference Model

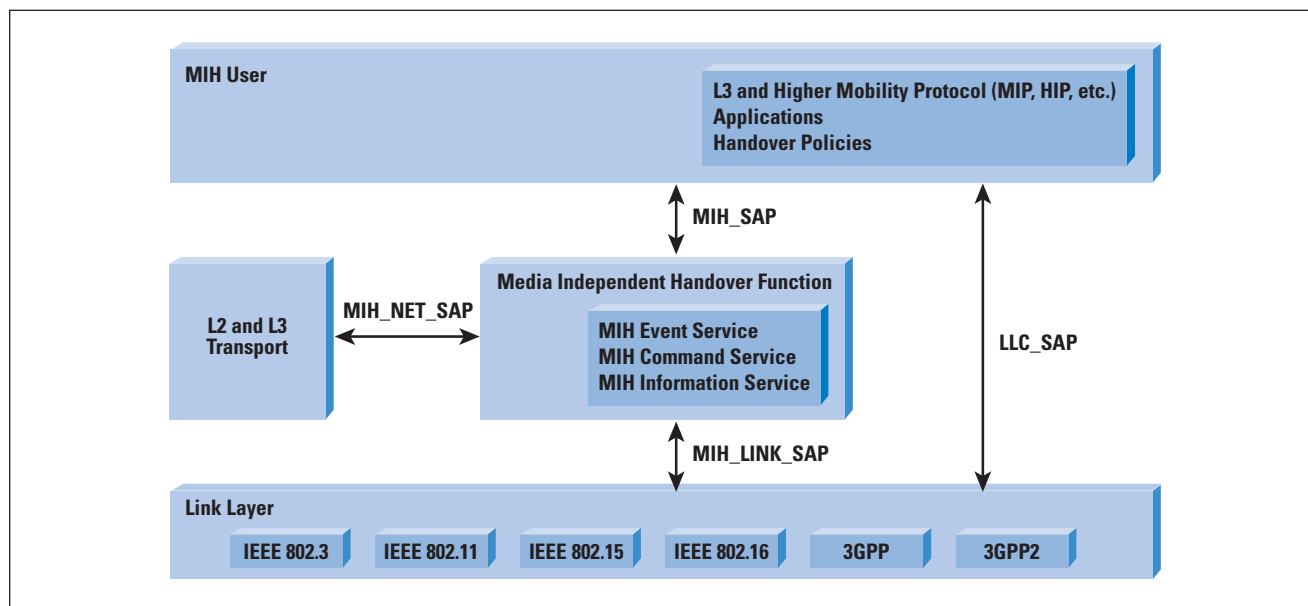
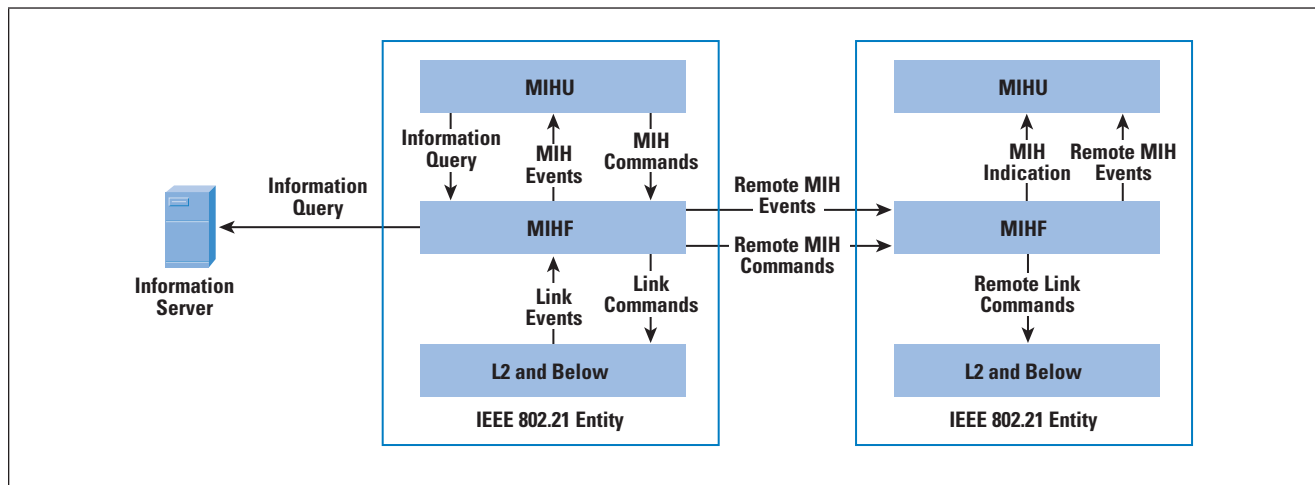


Figure 3 presents the messages directions of each MIHF service class, including both local and remote events and commands. The MIHF can subscribe to particular sets of events from a peer MIHF. Remote commands are initiated by local MIHUs and are conveyed to the peer MIHF through the local MIHF. Finally, MIIS information can be obtained through queries to the local database and to remote Information Servers.

Figure 3: MIHF Services



IEEE 802.21 Illustrated

Figure 4 illustrates an example topology where different wireless networks overlap. Imagine that the multiaccess mobile device user watches a high-bitrate IPTV channel as she moves in this area. Three wireless access technologies are considered in this example: Wi-Fi (IEEE 802.11), WiMAX (IEEE 802.16), and 3G/UMTS (3GPP). In this example, we assume that all networks and the mobile device are IEEE 802.21-compatible and that the Wi-Fi area is covered by several 802.11 PoAs, as would be the case in a campus- or citywide deployment.

Figure 4: Example Topology with Heterogeneous Overlapping Wireless Access Networks

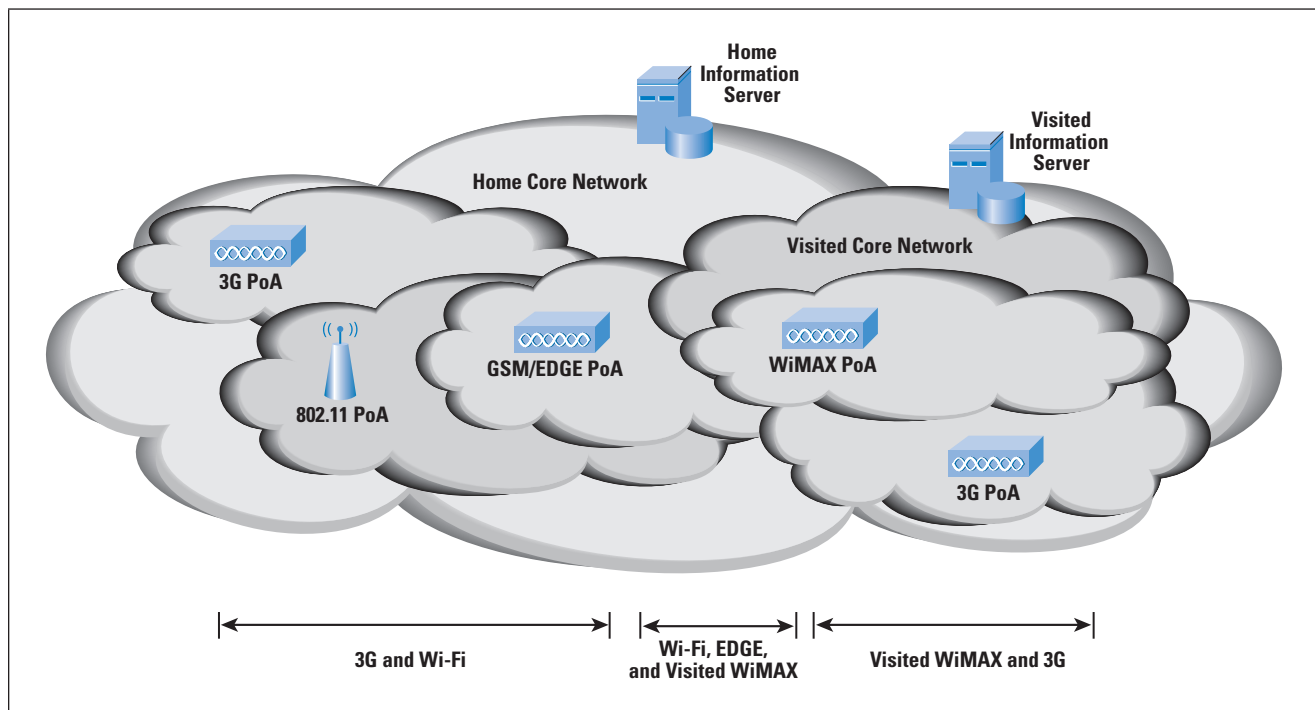
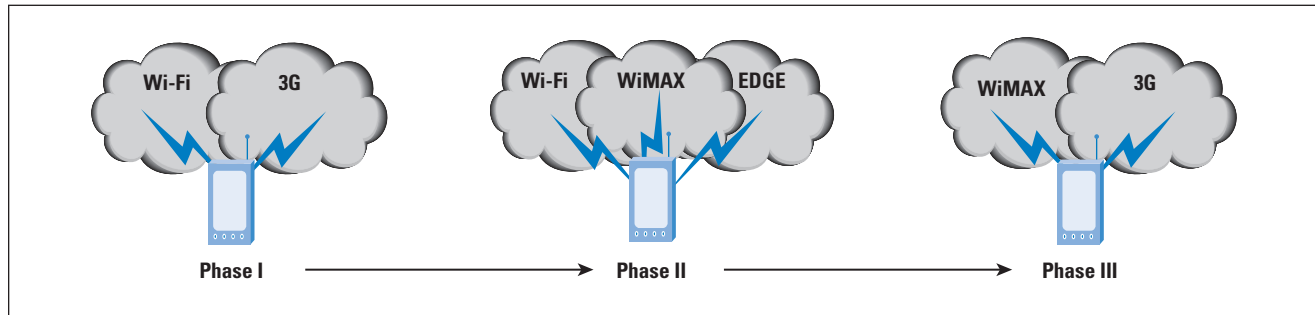


Figure 5 illustrates the network access environment as perceived by a mobile device in the area. The figure depicts three snapshots, indicating the overlapping networks in range at different locations. In order to deliver the IPTV stream transparently, for each of the available access networks we need to consider their effective available bandwidth, the associated cost per traffic unit, the terminal speed, the cell coverage area, the level of QoS support it can provide, and so on. Using information made available through the MIHF, we can determine which should be the next target access network.

Figure 5: Example Network Environment in Different Locations



In Phase I, the mobile node has two network access options. It can use a free and open Wi-Fi network or connect to the cellular operator's 3G/UMTS network. Note that opting to use the latter may, for instance, depend on the charging scheme of the operator. If subscribers pay based on traffic volume, one would assume that the free Wi-Fi network is a better option. On the other hand, as flat-rate plans become more popular, 3G may be a better option with its extended coverage and QoS guarantees. The IEEE 802.21 MIIS can provide this type of information, allowing for automation in dynamic access selection.

In Phase II, as the user moves, the device goes through a cellular technology handover from 3G/UMTS to *Enhanced Data rates for GSM Evolution* (EDGE)^[8]. At the same place, the public Wi-Fi network is still available and a new WiMAX network has just been detected. Assume that EDGE is not sufficient for delivering the IPTV stream. If in Phase I the network selection process opted for using the cellular network, then in Phase II the client application will experience significant degradation in service if it continues to use the EDGE access network. A vertical handover to the Wi-Fi or the WiMAX network should be considered. In contrast, if the mobile node first chose to stream the IPTV channel over the Wi-Fi access network, then it may need to reassess the situation based on events and link parameter reports using MIES and MICS, as we explain in the following sections. For example, an information query can reveal whether the WiMAX network is operated by a partner *Internet Service Provider* (ISP), and what the roaming cost would be.

Finally, in Phase III, the coverage area of the public Wi-Fi network ends. Through IEEE 802.21 services we find out that the only available networks are the roaming partner WiMAX and the home cellular network that is now offering 3G service.

The environment with several overlapping networks described previously and illustrated in Figures 4 and 5 is already a reality today in many places, and it is widely anticipated to be prevalent in the future. Next, we examine the three services defined by IEEE 802.21, namely MIES, MICS, and MIIS.

Media-Independent Event Service

Events indicate or predict changes in the state and transmission behavior of physical, data link, and logical link layers. In general, events are triggers for initiating candidate network discovery and handover procedures. The events defined in IEEE 802.21 are categorized as either *Link Events* or *MIH Events*, depending on their origin. Link events emanate from the link layers, whereas MIH events emanate from the MIHF and can be both remote and local. Local events propagate from lower layers to upper layers through the MIHF. Remote events occur at the protocol stack of another network entity and are transmitted from a peer MIHF to the local MIHF, as illustrated in Figure 3.

The *Media-Independent Event Service* (MIES) currently supports five types of events: MAC and PHY State Change events, Link Parameter events, Predictive events, Link Handover events, and Link Transmission events. A short introduction to the event types and corresponding events follows.

MAC and PHY State Change events correspond to state changes in MAC and *physical* (PHY) layers. The most characteristic events in this category are *Link_Up* and *Link_Down* events, which are generated when a Layer 2 connection with an access point is established or is torn down, respectively. Another event, called *Link_Detected*, indicates that a PoA has been detected but no affiliation is established yet.

Link Parameter events relate to changes in Layer 2 parameters. A *Link_Parameters_Report* can be sent when a MIHU has set thresholds for certain parameters. For example, a MIHU can set thresholds for the *Received Signal Strength Indicator* (RSSI) on IEEE 802.11 links, so that when a threshold is crossed proper action can be taken. A *Link_Parameters_Report* is also used for issuing periodical notifications about link conditions. Based on Link Parameter events, a MIHU can initiate the handover candidate discovery process, or trigger applications to adapt to changing link conditions.

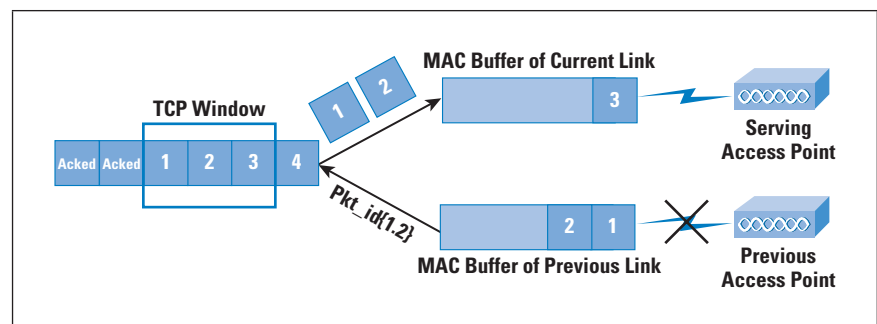
Predictive events inform about the probability of dramatic (negative) changes in link characteristics in the near future. For example, if strong decay in signal strength is observed, this decay may indicate imminent loss of link connectivity. Predictive events may include temporal information about when the actual event is expected to occur and what its presumed likelihood is. A *Link_Going_Down* event, for instance, may trigger a MIHU to consider possibilities for handing over to other available networks in range.

Link Handover events indicate the occurrence of Layer 2 handovers. The *Link_Handover_Imminent* event serves as a notification for an imminent handover, whereas a *Link_Handover_Complete* event reports the successful change of PoA. These events emanate from the link layer and are based solely on local Layer 2 information.

Link Transmission events show the transmission status of individual higher-layer *Protocol Data Units* (PDUs) at the link layer. Upper layers can, for example, adapt to data loss during a handover by improving buffer management based on Link Transmission events. These events may allow future upper-layer implementations to identify lost packets and recover without waiting for the expiration of retransmission timers.

Currently, for example, in the case of an ongoing session over TCP, the occurrence of a handover may have dramatic effects in performance. With IEEE 802.21, MIHUs can be informed about individual packets that have already been delivered to the sending buffer of the MAC layer but were not successfully transmitted before the handover occurred. In other words, the MAC layer outgoing buffer may contain TCP segments that cannot be delivered through the wireless network to the peer at the other end of the TCP connection. These segments were not successfully delivered from the local *Automatic Repeat-reQuest* (ARQ) module over the first hop, but are still buffered and cannot be transmitted because there is no link connectivity. In this case, TCP could use the information from Link Transmission events that identifies which packets need to be resent through the new access network, as illustrated in Figure 6 for packet numbers 1 and 2. Note, however, that IEEE 802.21 does not define any identifier for reliable packet identification, only the size of the packet ID (2 bytes), and it is up to the implementer to determine how different messages will be locally identified.

Figure 6: Link Transmission Event Indicating Undelivered Packets



Media-Independent Command Service

The *Media-Independent Command Service* (MICS) enables higher layers to control the stream of events originating from lower layers. Commands can originate from MIHUs (MIH commands) or from the MIHF (Link commands) and the destination can be the MIHF or any lower layer, respectively, as shown in Figure 3. The responses to Link commands are sent to MIHUs as indications. MIHUs can use command services to determine the status of different links in a uniform way, and control each interface accordingly, aiming for optimal connectivity. MICS defines the following set of commands that enable MIHUs to configure, control, and get information from the lower layers:

- *MIH commands* can be directed to lower layers residing at both local and remote MIHF entities. They originate from the upper layers and are directed to the MIHF. Similarly with MIH events, MIH commands can be both remote and local. MIH commands are typically used for network selection and handover management because they allow upper layers to initialize, prepare for, and execute handovers. MIH commands are also used to configure custom thresholds for link parameters. As mentioned previously, when set thresholds are crossed, MIHUs get the corresponding notifications through Link Parameter events.
- *Link commands* originate from the MIHF and are sent to lower layers in order to control their operation. Link commands can be issued only locally. Nevertheless, Link commands can be executed on behalf of local MIHUs, which could act on information received from a remote peer. Link commands are often initiated by MIHUs. For example, an MIHU can issue the *MIH_Get_Link_Parameters* MIH command, which when received by the local MIHF will lead to the generation of a remote *Link_Get_Parameters* Link command, as shown in Figure 3. This way, the MIHF can acquire the current parameter values of active link(s) for MIHU, and then deliver this information to the requesting MIHU. Note that MICS provides dynamic information about different link parameters, in contrast with MIIS, described next, which can report only static information.

Media-Independent Information Service

The *Media-Independent Information Service* (MIIS) facilitates handovers through a unified set of mechanisms that the MIHF can use to discover and obtain static (or rarely changing) information about networks in the vicinity of a multiaccess node. In other words, MIIS allows mobile nodes to check for available networks in range while using their currently active access network. MIIS information exchange occurs at the link layer (Layer 2) or network layer (Layer 3), so that all necessary information related to link layer or higher-layer services is collected before a mobile node authenticates with a new PoA.

MIIS defines a set of *Information Elements* (IEs) that are indispensable for network selection, classified into three groups: General Information and Access Network-Specific Information; PoA-Specific Information; and Other Information, which includes vendor- and network-specific details. The types of information handled by MIIS are solely related to handover decisions and conformance to the affiliation with the new PoA. Information relevant for assessing candidate networks by the handover machinery includes connection establishment details, such as PoA address and location; which security mechanisms are supported in a given access network; and what QoS guarantees can be provided.

General Information Elements and *Access Network-Specific Information Elements* give a general overview of neighboring networks. Information Elements may include, for instance, a list of available networks and their associated operators, roaming agreements and costs, and security and QoS support. For instance, user policies, defined at higher layers, may dictate that if a given access network operator charges users based on their traffic volume, then the network selector entity should not consider the corresponding access when a high-bitrate service, such as IPTV, is active.

PoA-Specific Information Elements refer to each PoA available in the access network and report PoA location and addressing information, supported data rates, PHY and MAC layer types, and channel parameters that can optimize link layer connectivity. Some additional information related to higher-layer services and individual capabilities of particular PoAs may be included as well. For instance, an advanced mobility manager on the mobile node can use the information about the geographical position of a PoA and compare it with the current or expected node location based on its mobility patterns. With careful planning and by taking advantage of this information, mobile nodes may be able to reduce the number of handovers and optimize the use of network resources.

MIIS provides mechanisms for issuing and responding to queries for Information Elements. Such information may reside in a separate server or in a local information database at the mobile node (see Figure 3). An MIHF could have access to an information server in its IEEE 802.21-enabled *Point-of-Service* (PoS) range from which it can obtain information regarding the home PoS and possibly other PoSs, such as those of roaming partners. If the home information server is not able to provide any information regarding the visited network, an MIIS query can be directed to the peer MIHF, residing in the visited PoS, which can access the visited PoS information server. Information queries can often be answered locally, based on information gathered from previous queries and by preprovisioning, for example, from the information server.

Information Elements and their relationships are captured in an Information Service schema which, in turn, defines the information structure. IEEE 802.21 specifies that information that is to be presented across different technologies should be in a standardized, common, and open format, such as XML or *Type Length Value* (TLV).

Service Management

In order to use and provide MIHF services, MIHF entities need to be configured appropriately. IEEE 802.21 defines three service management functions: MIH capability discovery, MIH registration, and MIH event subscription.

MIHF may discover other MIHF entities and their capabilities using the MIH capability discovery procedure. Depending on the information obtained from this procedure, the local MIHF can determine which peer MIHFs it should register with. The MIH capability discovery function uses the MIH protocol (introduced in the following section) at Layer 2 or Layer 3, and media-specific Layer 2 broadcast messages are allowed. For example, an MIHF can listen to media-specific broadcast messages, such as IEEE 802.11 beacons, or media-independent Layer 2 *MIH_Capability_Discover* broadcast messages, because an MIHF entity residing in the network may announce its existence and capabilities periodically. MIHF can also send *MIH_Capability_Discover* request messages using multicast or unicast to detect peer MIHFs in a solicited way. For instance, MIHF can send a request by unicast for obtaining the capabilities of a specific IEEE 802.21 network entity. In this case, only the IEEE 802.21 network entity addressed should respond to these request messages.

MIH registration is a symmetric procedure by which two peer MIHFs authenticate and can then communicate with each other in a more trusted manner. After MIH registration is completed, the two peer MIHF entities can symmetrically request services from their registered peer. Note that MIH registration is not necessary for obtaining some level of support from a peer MIHF. However, by registering and authenticating, peer MIHFs typically will get access to much more extensive information. That is, although the MIHF residing on the mobile node may be able to access information services from the network-side MIHFs without registration and authentication, the available information may be only a subset of that provided after authenticating.

Finally, MIH event subscription enables MIHUs to subscribe to a particular set of events provided by MIES from the local or peer MIHF. Event subscription from a peer MIHF requires registration and knowledge about its capabilities. The subscription contains only the list of events the MIHU is interested in. Note that event sources may not be necessarily capable of providing all events that the subscriber is interested in subscribing to. Each subscription request is matched by a confirmation message from the event source indicating the events approved for subscription.

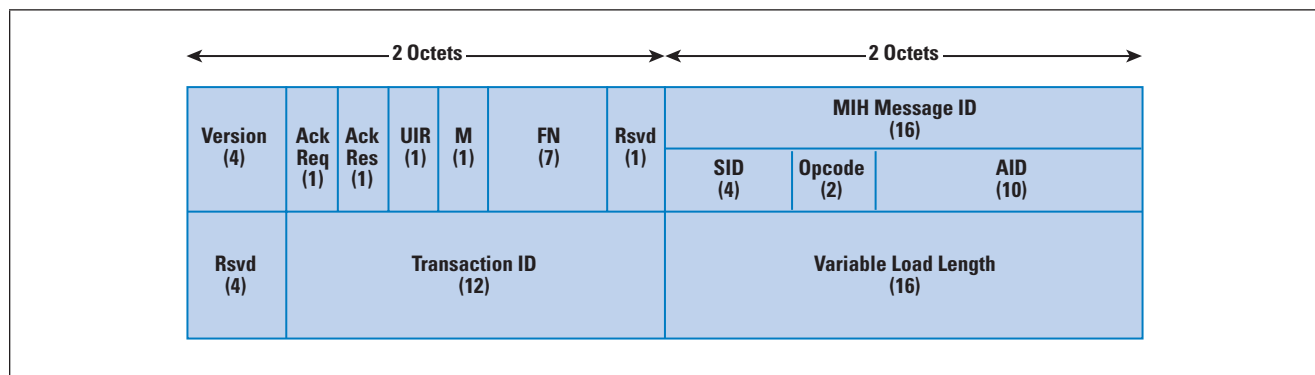
Media-Independent Handover Protocol

The *Media-Independent Handover Protocol* (MIHP) specifies the rules and services for unified communication between peer MIHFs. The protocol defines the message format, header, and encoding format and is meant to be used solely for communicating with peer MIHF entities. For internal communication no particular encoding is dictated.

MIH protocol messages can be carried over Layer 2 management frames, Layer 2 data frames, or over Layer 3/IP transport. Note that cellular technologies do not provide Layer 2 transport without changes in their protocol stack.

The MIH protocol messages, or frames, comprise a header part and a TLV-encoded payload part. The MIHF frame header consists of eight octets. Figure 7 illustrates the MIH protocol header indicating the corresponding bit length for each field in parentheses.

Figure 7: MIH Protocol Header



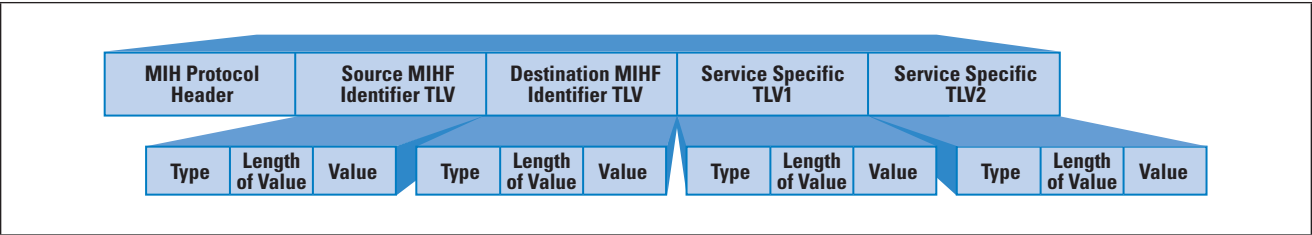
The *Version* field in the MIH frame header specifies the version of the MIH protocol used. The two *Ack* fields are for acknowledgement purposes and are discussed later in the article. The *Unauthenticated Information Request* (UIR) flag indicates that the response message may be sent with a limited length because of the nature of unauthenticated message exchange. Recall that when an MIHF issues requests without registering first with its peer, it may receive less information than if it had registered earlier. If this flag is set, then the information included in the response message may not reflect the complete information available to registered MIHFs. The *More Fragments* (M) and *Fragment Number* (FN) fields are used in message fragmentation.

The *MIH Message ID* field comprises three subfields. The *Service Identifier* (SID) field indicates the MIHF service class (MIES, MICS, MIIS, or Service Management) that this message belongs to. The *Operation code* (Opcode) specifies whether the message is a request, response, or indication. The *Action Identifier* (AID) is related with and scoped by the SID. For instance, if the SID indicates MIES, AID points to the actual event type. The *Variable Load Length* field contains the total length of the variable, TLV-encoded payload carried by this message frame.

The MIH protocol messages use the *Transaction ID* and *MIHF ID* fields as identifiers, but only the former is included in the header. The Transaction ID field is an identifier that helps to match each request, response, or indication message with its acknowledgement.

The payload part contains service-specific messages encoded in TLV format. The first two TLVs in the payload part (not shown in Figure 7) should be the *Source Identifier* and *Destination Identifier*, which are both the same data type as the MIHF ID. Every MIHF must have a unique MIHF ID, which may be assigned to it at configuration time. The MIHF ID shall be invariant and could be, for example, a *Fully Qualified Domain Name* (FQDN) or *Network Access Identifier* (NAI). The MIHF ID is used during the MIH registration phase and is appended to the payload part of every message requiring endpoint identification. In broadcast messages, the Destination Identifier TLV is defined as zero length. Figure 8 shows the message structure consisting of the MIH Protocol header, source and destination identifiers, and service-specific TLVs. In TLV encoding, the Type field (1 octet) denotes the parameter type, the Length field (variable octets) indicates the length of the Value field, and the Value field (variable octets) carries the actual value of the parameter.

Figure 8: MIH Protocol Frame Structure



Acknowledging MIH messages is not mandatory. Still, the MIH protocol does support the use of acknowledgements to ensure reliable message exchange. The sender MIHF can set the *ACK-Req* field to instruct the receiver to return an acknowledgement with *ACK-Rsp* bit set. The *MIH Message ID* and *Transaction ID* must be the same in the request message and its acknowledgement. An acknowledgement message may carry no payload. Note, however, that despite employing these two ID fields, the MIH protocol does not specify any further mechanisms for reliable authentication or shielding message exchanges from third parties.

MIH Communication Model

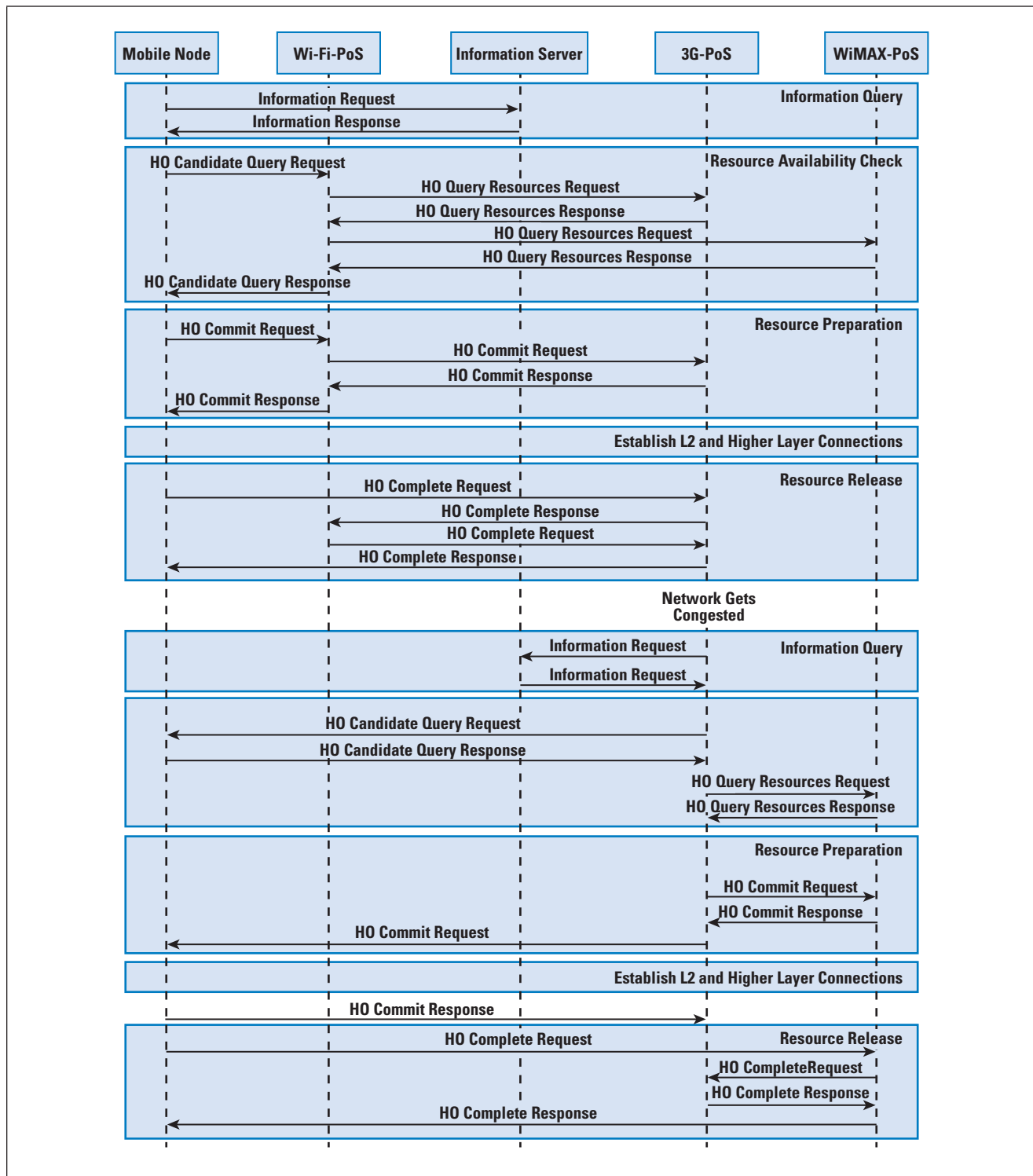
The MIHF communication model specifies different MIHF roles and their communication relationships, such as supported transport mechanisms and service classes. The assigned MIHF roles depend on their location in the network. For example, an MIHF on a mobile node can communicate directly with network-side entities called *MIH PoSs* using Layer 2 or Layer 3 communication. MIH PoSs may include the serving PoA or candidate PoAs. Network-side MIHFs can communicate with each other at Layer 3 or above using the MIH protocol, introduced in the previous section.

Let us revisit the example use case of IEEE 802.21 illustrated in Figures 4 and 5. Figure 9 presents the IEEE 802.21 message exchanges in mobile- and network-initiated handover procedures in the case where the mobile node hands over from a Wi-Fi to the 3G cellular network (between Phase II and Phase III in Figure 5) and then hands over to a WiMAX network (Phase III in Figure 5). First, during the discovery of handover candidate PoAs, the mobile node MIHF employs MIIS to gather static information about the surrounding networks. The request is issued over the currently used Wi-Fi access. This information is obtained from the information server that may reside in a different network than the one currently in use.

After receiving the response to its Information Request, the mobile node initiates the handover process by querying about the availability of resources in the networks it is interested in. These requests are sent through the serving PoS (*Wi-Fi-PoS* in Figure 9), which disseminates the requests to the MIH PoSs of the candidate networks (*3G-PoS* and *WiMAX-PoS* in Figure 9). The response indicating the capabilities of the two candidate networks is returned to the mobile node MIHF from the serving PoS. After receiving this information, an MIHU on the mobile node decides which network to hand over to, based on policies and the output of its network selection algorithms. Then a *Handover Commit Request* message is sent, and after the candidate network has made its final commitment for the handover (and the appropriate resources are reserved successfully), the mobile node establishes a Layer 2 connection with the PoA in the area of the candidate PoS, that is, the *3G-PoS* in our example case. Following this successful intertechnology handover, the resources used in the previous link can optionally be released. In the case where no resources are explicitly reserved, this step is skipped.

As we progress in the timeline of our example case, the network-side MIHU initiates a handover to the WiMAX network. This handover could be, for example, the result of observing congestion in the cellular network that indicates that a new PoS should be found for the mobile node. The serving PoS (*3G-PoS*) collects information about networks in the range of the mobile node from the Information Server. Upon determining that a suitable WiMAX candidate network that can serve the mobile node exists, the *3G-PoS* triggers a network-initiated handover. First, the serving PoS requests permission from the mobile node to proceed with the handover. If the mobile node does not object, the serving PoS proceeds with the rest of the handover procedure, which is similar to the mobile-initiated handover described previously except that it is handled by a network entity.

Figure 9: IEEE 802.21-Assisted Handover Message Sequence Diagram



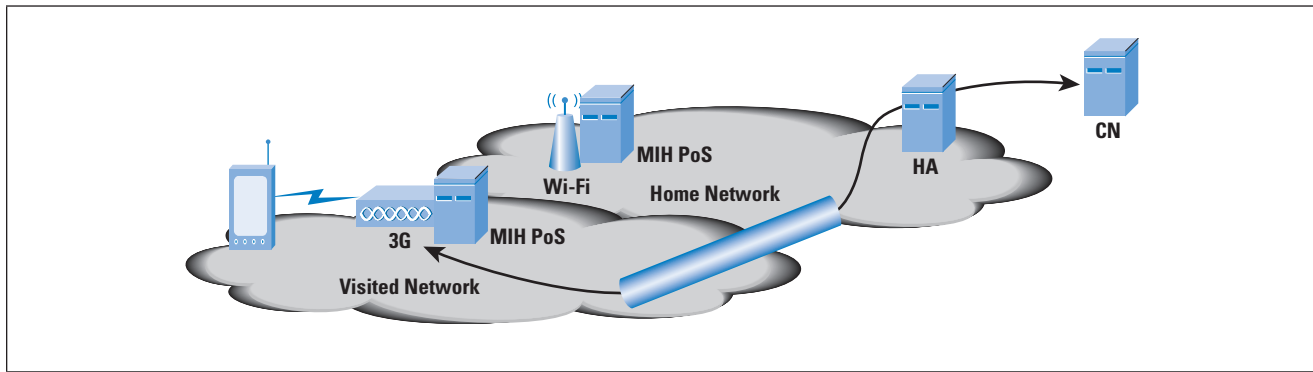
Handover Execution

As illustrated in the example, the handover decision and target assessment constitute a multiphase process where the assistance of IEEE 802.21 is essential. However, the actual handover execution is outside the scope of the standard. This section briefly describes how handovers can be carried out by MIP with the cooperation of IEEE 802.21. After choosing the target network by capitalizing on the IEEE 802.21 services, the mobile node establishes a new connection with the handover target network while still routing traffic through the currently serving network. The mobile node obtains a *Care-of Address* (CoA) for this new link from the IP address space of the target network. The CoA is an IP address assigned to the new link of the mobile node and is used while connected to the visiting network^[11]. With MIPv4, the CoA is provided by a *Foreign Agent* (FA) in the visited network, which also acts as a router for the mobile node^[12]. With MIPv6, the Foreign Agent is not needed^[13] and the CoA is obtained directly, say, for example, from a *Dynamic Host Configuration Protocol* (DHCP) server. The mobile node can obtain the IP address of the DHCP server in the target network through the IEEE 802.21MIIS.

In MIP, each mobile node has a *Home Agent* (HA), which routes the traffic of the mobile node. After successfully affiliating with a PoA in the target network, the mobile node notifies the Home Agent of the CoA by performing a binding update. In a bidirectional tunnel mode, the Home Agent establishes an IP-IP tunnel between the Home Agent and the Foreign Agent (MIPv4) or the Home Agent and the mobile node CoA (MIPv6). This mode does not require any binding updates on the *Correspondent Node* (CN). In other modes, either the uplink traffic of the mobile node is sent directly to the Correspondent Node using the CoA as source address, or all bidirectional communication between the Correspondent Node and the mobile node uses the CoA only. In the first case, traffic from the Correspondent Node to the mobile node travels through the Home Agent, but in the latter case there is no need for the Home Agent detour. However, these modes need address binding at the Correspondent Node and are in practice less frequently used than the bidirectional tunnel mode.

Figure 10 illustrates a situation where a link with the Wi-Fi PoA is broken down by the mobile node and the IPv6 traffic between the Correspondent Node and the mobile node, now employing IEEE 802.21-enabled 3G network, travels through the tunnel between Home Agent and the mobile node.

Figure 10: Mobile IPv6 Tunnel



Layer 3 handover executions based on RFC 3344^[12] and RFC 3775^[13] may often exceed the typical handover delay budgets, thus introducing gaps in connectivity that are perceptible at the application layer. Recent standardization efforts have focused on decreasing handover delays by enhancing MIP so that it can provide for transparent mobility management for both IPv4^[16] and IPv6^[17, 18]. The proposed enhancements either reduce the amount of signaling or allow the mobile node to configure the new Layer 3 connection before reassociating with the new network. In this context, IEEE 802.21 can provide the essential information for preestablishing the connection based on media-independent Layer 2 link detection events as well as static address information from the target network.

Summary and Outlook

We presented an overview of the IEEE 802.21 Media-Independent Handover Services standard. We anticipate that its adoption in the near future will allow for better network resource usage and permit multiaccess devices to select the network access best suited for their communication needs. After motivating the needs for a standard to cope with heterogeneous network handovers, we introduced the IEEE 802.21 Reference Model and the MIH Services. We briefly presented the MIH Protocol, although a more thorough description calls for a separate overview article. Finally, we illustrated network operation when IEEE 802.21 is adopted using example use cases featuring both network- and terminal-initiated intertechnology (or vertical) handovers.

We expect that in the future, when IEEE 802.21-2008 is widely deployed, there will be significant efforts to further amend and extend it in order to provide for even better services. In fact, because security mechanisms are outside the scope of the base IEEE 802.21 standard, the work on defining a security-related extension to IEEE 802.21 (IEEE P802.21a) has already begun. Moreover, another amendment (IEEE P802.21b) that deals with handovers with downlink-only technologies, such as *Digital Video Broadcasting* (DVB), has also been introduced (see www.ieee802.org/21 for more information about the amendments). Nevertheless, it remains uncertain whether vendors will stand by this promising standard and incorporate it in future products and solutions.

References

- [1] T. Sridhar, “Wi-Fi, Bluetooth and WiMAX,” *The Internet Protocol Journal*, Volume 11, No. 4, December 2008.
- [2] E. Gustafsson and A. Jonsson, “Always Best Connected,” *IEEE Wireless Communications*, Volume 10, No. 1, February 2003.
- [3] IEEE Std 802.21-2008, *IEEE Standard for Local and Metropolitan Area Networks—Part 21: Media Independent Handover Services*, IEEE, January 2009.
- [4] H. Kaaranen, S. Naghian, L. Laitinen, A. Ahtiainen, and V. Niemi, *UMTS Networks: Architecture, Mobility and Services*, 2nd Edition, John Wiley & Sons, 2005.
- [5] V. Vanghi, A. Damnjanovic, and B. Vojcic, *The cdma2000 System for Mobile Communications: 3G Wireless Evolution*, Prentice Hall, 2004.
- [6] Y. Mälarstig, O. Holmström, and P. Davidsson, *Svensk telemarknad 2007*, PTS-ER-2008:15, ISSN 1650-9862, June 2008.
- [7] J. Pinola and K. Pentikousis, “Mobile WiMAX,” *The Internet Protocol Journal*, Volume 11, No. 2, June 2008.
- [8] K. Pentikousis, “Wireless Data Networks,” *The Internet Protocol Journal*, Volume 8, No. 1, March 2005.
- [9] M. Balazinska and P. Castro, “Characterizing Mobility and Network Usage in a Corporate Wireless Local-Area Network,” *Proc. First International Conference on Mobile Systems, Applications, and Services (MobiSys)*, San Francisco, California, USA, May 2003, pp. 303–316.
- [10] T. Henderson, D. Kotz, and I. Abyzov, “The Changing Usage of a Mature Campus-wide Wireless Network,” *Computer Networks*, Volume 52, No. 14, October 2008, pp. 2690–2712.
- [11] W. Stallings, “Mobile IP,” *The Internet Protocol Journal*, Volume 4, No. 2, June 2001.
- [12] C. Perkins (Ed.), “IP Mobility Support for IPv4,” RFC 3344, August 2002.
- [13] D. Johnson, C. Perkins, and J. Arkko, “Mobility Support in IPv6,” RFC 3775, June 2004.

- [14] K. Pentikousis, R. Agüero, J. Gebert, J. A. Galache, O. Blume, and P. Pääkkönen, "The Ambient Networks Heterogeneous Access Selection Architecture," *Proc. First Ambient Networks Workshop on Mobility, Multiaccess, and Network Management (M2NM)*, Sydney, Australia, October 2007, pp. 49–54.
- [15] J. Mäkelä and K. Pentikousis, "Trigger Management Mechanisms," *Proc. Second International Symposium on Wireless Pervasive Computing (ISWPC)*, San Juan, Puerto Rico, February 2007, pp. 378–383.
- [16] K. El Malki (Ed.), "Low Latency Handoffs in Mobile IPv4," RFC 4881, June 2007.
- [17] H. Soliman, C. Castelluccia, K. El Malki, and L. Bellier, "Hierarchical Mobile IPv6 Mobility Management," RFC 4140, August 2005.
- [18] R. Koodli (Ed.), "Mobile IPv6 Fast Handovers," RFC 4068, July 2005.

ESA PIRI received his M.Sc. from the University of Oulu, Oulu, Finland, in 2008. During his studies, he specialized in information networks systems and wrote his Master's thesis on mobility management issues in heterogeneous networks. Currently he is working as a Research Scientist at VTT Technical Research Centre of Finland in Oulu, Finland. He can be contacted by e-mail at: **esa.piri@vtt.fi**

KOSTAS PENTIKOUSIS studied computer science at Aristotle University of Thessaloniki, Greece (B.Sc. 1996, summa cum laude) and State University of New York at Stony Brook, USA (M.Sc. 2000, Ph.D. 2004). He has published internationally in several areas, including mobile computing; transport protocols; applications; network traffic measurements and analysis; and simulation and modeling. Dr. Pentikousis is a Senior Research Scientist at VTT Technical Research Centre of Finland. Visit <http://ipv6.willab.fi/kostas> for more information and contact details.

Book Review

Geeks Bearing Gifts

Geeks Bearing Gifts v1.1: How the computer world got this way, by Ted Nelson, ISBN: 978-0-578-00438-9, Published by Mindful Press, 2009, distributed through Lulu.Com, <http://www.lulu.com>

In a short but interesting book, computer pioneer Ted Nelson takes a very broad look at the origins and evolution of many of the basic ideas that underpin today's computer industry. The emphasis is on concepts and technologies rather than the success of individuals, the companies they founded, and the shape of the computer industry. This approach differentiates the book from other accounts, such as Robert X. Cringley's *Accidental Empires* and Martin Campbell-Kelly's *From Airline Reservations to Sonic the Hedgehog*.

Although the book is suitable for a fairly broad readership, an appreciation of the current makeup of the industry is helpful in understanding the significance of some of Nelson's ideas.

Organization

Geeks Bearing Gifts is divided into 60 short chapters, arranged in chronological order from the time the ideas originated, rather than when they appeared in fully developed form (indeed many are still developing). In the initial chapters Nelson covers topics such as language, alphabets, and encryption before moving on to examine the origins of computing. He then examines the contribution of pioneers from both inside and outside the United States, giving more credibility to contributors from outside of the United States than is normal.

As would be expected, Nelson deals in some detail with the topic of information presentation, in particular the origins of hypertext and associated developments such as *Xanadu* and the World Wide Web. He discusses the differences between these technologies, spending some time reflecting on his attempts to develop *Xanadu* at Brown University; he suggests that many of the deficiencies of the Web come from misdirection of that phase of the project.

Nelson next examines a wide selection of topics ranging from networks (both local and the Internet), object-orientated programming, and early desktop machines, before reaching the pivot point of his book: the UNIX operating system. He chose UNIX as the fulcrum of his analysis because he believes "so much led into it and so much has resulted from it."

Nelson next considers PUI (the PARC user interface), PCs, the role of the Microsoft and Apple operating systems and their evolution, the influence of the spreadsheet, the Internet, browsers, the Internet crash, and the current major companies in computing. He explores the promise, hype, and reality of the Web 2.0 model and its likely influence. (PARC stands for the Xerox Palo Alto Research Center.)

The last two chapters are summaries and thought guides. The first of these suggests that it is people and ideas rather than technology that advance the computer industry and that the myth of technological necessity has stifled imagination. The final chapter illustrates what the book is about—the disagreements and decisions that have made the technical world what it is today.

Synopsis

Nelson captures most of the important developments in the computer industry, although he acknowledges that in 199 pages it is possible to tell the reader only a little of where the software ideas come from and what they are. He sets out to show how varied and conflicting the initiatives that have propelled the evolution of computer technology have been, exposing the “ideas, disagreements, manoeuvres, forgotten possibilities, and politics.”

The book reads like a collection of themed essays, rather than a coherent sequence of stories. Nonetheless it is both informative and thought-provoking.

The Author

Ted Nelson is considered to be a radical thinker; he is one of the pioneers of the computer industry initiating the Xanadu project, which was started in the early 1960s with the objective of developing a computer network with a simple user interface. He is credited with inventing the term “hypertext.”

He holds a first degree in philosophy, a Masters in sociology, and a Doctorate in Media and Governance. Among his honors are a visiting fellowship at the Oxford Internet Institute and a Fellowship of Wadham College, Oxford; in addition, France has knighted him as “Officier des Arts et Lettres.” Visit:

http://en.wikipedia.org/wiki/Ted_Nelson

and

<http://www.ibiblio.org/pioneers/nelson.html>

...for more information.

—Edward Smith, BT, UK

edward.a.smith@btinternet.com

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. Contact us at ipj@cisco.com for more information.

Fragments

RIPE Announces IPv6 Website

The RIPE NCC recently announced the launch of the *IPv6 Act Now!* website. Available at www.IPv6ActNow.org, the website explains IPv6 in terms that everyone can understand and provides a variety of useful information aimed at promoting the global adoption of IPv6. The site is designed for anyone with an interest in IPv6, including network engineers, company directors, law enforcement agencies, government representatives and civil society. The content is regularly updated and includes:

- Education, advice and opinions from the experts
- Latest IPv6-related news stories
- Videos and articles from Internet community leaders
- Current IPv4 exhaustion and IPv6 uptake statistics
- The RIPE community's statement on IPv6 deployment
- Information on community-developed IPv6 distribution policies
- Useful links to other sources of information about IPv6
- A forum for everyone to share experiences, ask questions and find answers

The site also includes contributions from other *Regional Internet Registries* (RIRs) and industry partners. If you have and comments or suggestions about the website, please contact:

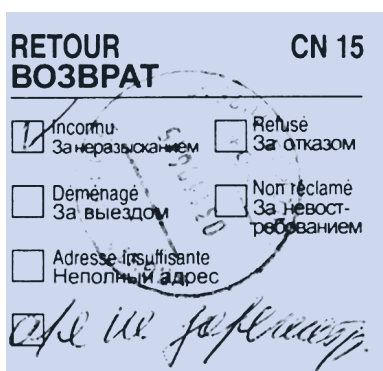
ipv6actnow@ripe.net

Four-byte AS numbers from APNIC

From July 1, 2009, the *Asia Pacific Network Information Centre* (APNIC) will assign four-byte *Autonomous System* (AS) numbers by default when receiving requests. Two-byte AS numbers will only be assigned if the applicant can demonstrate that a four-byte only AS number is unsuitable. This change marks the next phase of the transition to four-byte AS numbers. The final phase begins in January 2010, when APNIC will cease to make any distinction between two-byte and four-byte AS numbers, and will operate AS number assignments from an undifferentiated four-byte AS number pool. For more information please see: <http://icons.apnic.net/asn>

Please Tell Us When You Move

We receive large quantities of undeliverable copies of *The Internet Protocol Journal*. For international mailings, the returned mail piece usually includes a standard CN 15 label, an example of which is shown here. We have an extensive collection of CN 15 labels from all over the world, but we would much rather ensure that your journal is delivered to the correct address. So, if you're moving your home or office, please use the online subscription system to update your details, or just send an e-mail message to ipj@cisco.com with the new information. You can also suspend paper delivery and read IPJ online if you wish.



Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ contains standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2009 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

September 2009

Volume 12, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Cloud Computing.....	2
End-to-End Security.....	20
Letter to the Editor	27
Fragments	28

FROM THE EDITOR

This journal has covered numerous emerging technologies since we started publishing in June 1998. It would be an interesting exercise to look at which of these technologies have been successfully deployed, which ones have been rejected, and which ones are still emerging or slowly being deployed. In this issue we examine another emerging technology, or perhaps “a new concept” would be a better term, because a collection of new and old technologies are coming together to form what is collectively known as *Cloud Computing*. In a two-part article on cloud computing, T. Sridhar gives an overview of the concepts underlying this area of development. Part 1 of the article is subtitled “Models and Technologies.” It will be followed by Part 2: “Infrastructure and Implementation Topics,” which will be published in our next issue.

In the last year, I have had one of my credit cards “compromised” (unauthorized charges posted to the account) and subsequently replaced twice. This situation is always annoying and worrisome. Most likely, these breaches resulted from the card information being captured through an online purchase transaction. I am sure I will never know the full story, and luckily the credit card companies are pretty good about detecting fraudulent charges and quickly resolving the matter. When you start thinking about the number of network and server elements involved in a typical e-commerce transaction, it isn’t entirely surprising that someone with criminal intentions could exploit a weakness in the overall system. Our second article, by Michael Behringer, explores the topic of “end-to-end security” in more detail.

Those of you who have been subscribers to this journal for several years have probably noticed that your subscription has been “auto-renewed” once a year without requiring any renewal action on your part. Starting with the December 2009 issue, we will no longer extend your subscription when it expires unless you renew it by visiting the IPJ “Subscriber Services” webpage. You will need to use your e-mail address and Subscription ID in order to gain access to your record, where you can renew, update your delivery address, or change delivery method. IPJ is available on paper, as well as online in both HTML and PDF formats. You can also contact us at ipj@cisco.com regarding your renewal. The expiration date and Subscription ID are printed on the back of the journal for subscribers in the United States, and on the envelope for our international subscribers. We believe that this new renewal policy will result in fewer undeliverable or unwanted copies being mailed out—a plus for the environment.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

Cloud Computing—A Primer

Part 1: Models and Technologies

by T. Sridhar

Cloud computing is an emerging area that affects IT infrastructure, network services, and applications. Part 1 of this article introduces various aspects of cloud computing, including the rationale, underlying models, and infrastructures. Part 2 will provide more details about some of the specific technologies and scenarios.

The term “cloud computing” has different connotations for IT professionals, depending upon their point of view and often their own products and offerings. As with all emerging areas, real-world deployments and customer success stories will generate a better understanding of the term. This discussion starts with the *National Institute of Standards and Technology* (NIST) definition:

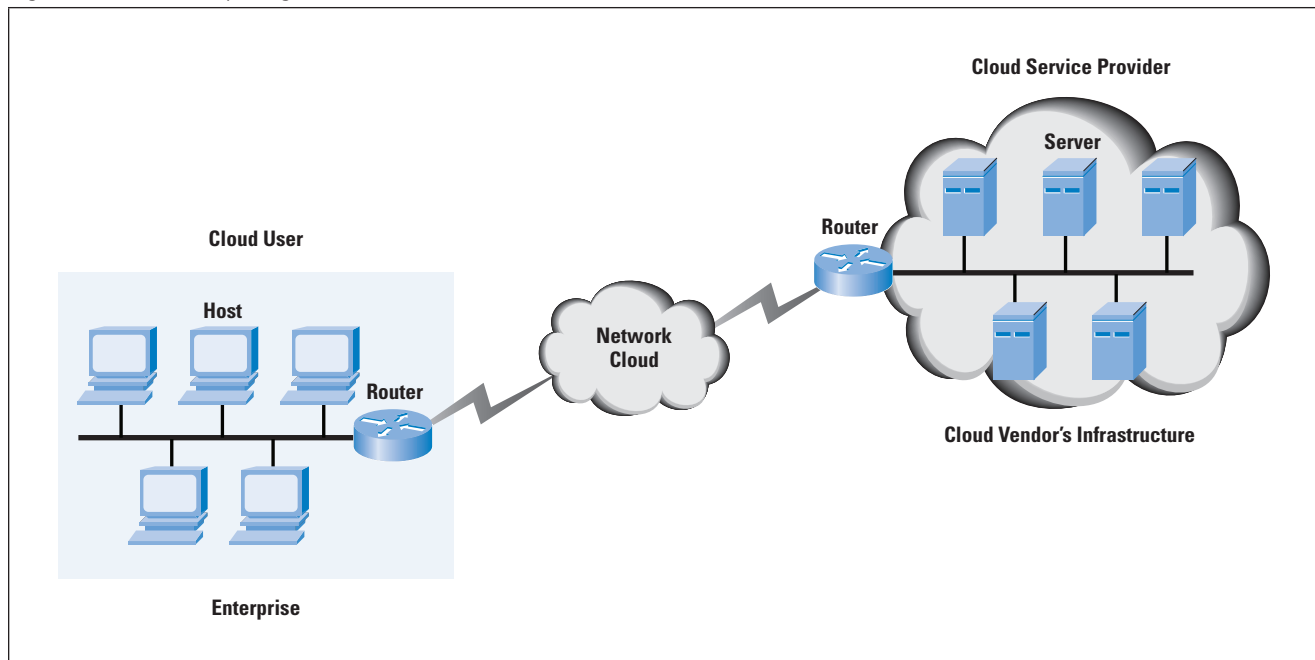
“Cloud computing is a model for enabling convenient, on-demand network access to a shared pool of configurable computing resources (for example, networks, servers, storage, applications, and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction.”

The following is a list of characteristics of a cloud-computing environment. Not all characteristics may be present in a specific cloud solution.

- *Elasticity and scalability:* Cloud computing gives you the ability to expand and reduce resources according to your specific service requirement. For example, you may need a large number of server resources for the duration of a specific task. You can then release these server resources after you complete your task.
- *Pay-per-use:* You pay for cloud services only when you use them, either for the short term (for example, for CPU time) or for a longer duration (for example, for cloud-based storage or vault services).
- *On demand:* Because you invoke cloud services only when you need them, they are not permanent parts of your IT infrastructure—a significant advantage for cloud use as opposed to internal IT services. With cloud services there is no need to have dedicated resources waiting to be used, as is the case with internal services.
- *Resiliency:* The resiliency of a cloud service offering can completely isolate the failure of server and storage resources from cloud users. Work is migrated to a different physical resource in the cloud with or without user awareness and intervention.
- *Multitenancy:* Public cloud services providers often can host the cloud services for multiple users within the same infrastructure. Server and storage isolation may be physical or virtual—depending upon the specific user requirements.

- *Workload movement*: This characteristic is related to resiliency and cost considerations. Here, cloud-computing providers can migrate workloads across servers—both inside the data center and across data centers (even in a different geographic area). This migration might be necessitated by cost (less expensive to run a workload in a data center in another country based on time of day or power requirements) or efficiency considerations (for example, network bandwidth). A third reason could be regulatory considerations for certain types of workloads.

Figure 1: Cloud Computing Context



Cloud computing involves shifting the bulk of the costs from *capital expenditures* (CapEx), or buying and installing servers, storage, networking, and related infrastructure) to an *operating expense* (OpEx) model, where you pay for usage of these types of resources. Figure 1 provides a context diagram for the cloud.

How Is Cloud Computing Different from Hosted Services?

From an infrastructure perspective, cloud computing is very similar to *hosted services*—a model established several years ago. In hosted services, servers, storage, and networking infrastructure are shared across multiple tenants and over a remote connection with the ability to scale (although scaling is done manually by calling or e-mailing the hosting provider). Cloud computing is different in that it offers a pay-per-use model and rapid (and automatic) scaling up or down of resources along with workload migration. Interestingly, some analysts group all hosted services under cloud computing for their market numbers.

Virtualization and Its Effect on Cloud Computing

It can be argued to good effect that cloud computing has accelerated because of the popularity and adoption of virtualization, specifically server virtualization. So what is virtualization? Here, virtualization software is used to run multiple *Virtual Machines* (VMs) on a single physical server to provide the same functions as multiple physical machines. Known as a *hypervisor*, the virtualization software performs the abstraction of the hardware to the individual VMs.

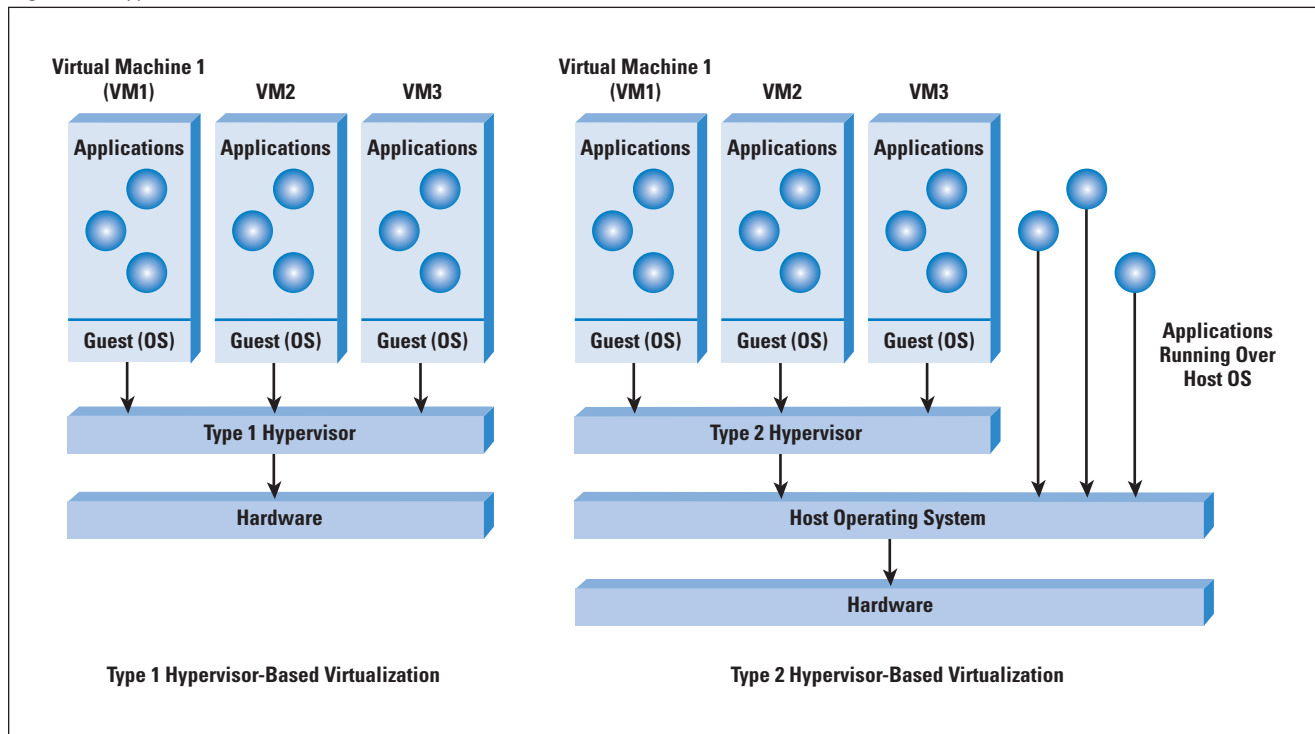
Virtualization is not new—it was first invented and popularized by IBM in the 1960s for running multiple software contexts on its main-frame computers. It regained popularity in the past decade in data centers because of server usage concerns. Data centers and web farms consisted of multiple physical servers. Measurement studies on these server farms noted that individual server usage was often as low as 15 percent for various reasons, including traffic loads and the nature of the applications (available, not always used fully), among others. The consequence of this server sprawl with low usage was large financial outlays for both CapEx and OpEx—extra machines and related power and cooling infrastructure and real estate.

Enter virtualization. A hypervisor is implemented on a server either directly running over the hardware (a *Type 1 hypervisor*) or running over an *operating system* (OS) (a *Type 2 hypervisor*). The hypervisor supports the running of multiple VMs and schedules the VMs along with providing them a unified and consistent access to the CPU, memory, and I/O resources on the physical machine. A VM typically runs an operating system and applications. The applications are not aware that they are running in a virtualized environment, so they do not need to be changed to run in such an environment. Figure 2 depicts these scenarios. The OS inside the VM may be virtualization-aware and require modifications to run over a hypervisor—a scheme known as paravirtualization (as opposed to *full virtualization*).

VM Migration: An Advantage of Virtualization

Some vendors have implemented VM migration in their virtualization solution—a big advantage for application uptime in a data center. What is VM migration? Consider the case of a server with a hypervisor and several VMs, each running an OS and applications. If you need to bring down the server for maintenance (say, adding more memory to the server), you have to shut down the software components and restart them after the maintenance window—significantly affecting application availability. VM migration allows you to move an entire VM (with its contained operating system and applications) from one machine to another and continue operation of the VM on the second machine. This advantage is unique to virtualized environments because you can take down physical servers for maintenance with minimal effect on running applications.

Figure 2: Hypervisors in Virtualization



You can perform this migration after suspending the VM on the source machine, moving its attendant information to the target machine and starting it on the target machine. To lower the downtime, you can perform this migration while the VM is running (hence the name “live migration”) and resuming its operation on the target machine after all the state is migrated.

The following are some of the benefits of virtualization in a cloud-computing environment:

- *Elasticity and scalability:* Firing up and shutting down VMs involves less effort as opposed to bringing servers up or down.
- *Workload migration:* Through facilities such as live VM migration, you can carry out workload migration with much less effort as compared to workload migration across physical servers at different locations.
- *Resiliency:* You can isolate physical-server failure from user services through migration of VMs.

It must be clarified that virtualization is not a prerequisite for cloud computing. In fact, there are examples of large cloud service providers using only commodity hardware servers (with no virtualization) to realize their infrastructure. However, virtualization provides a valuable toolkit and enables significant flexibility in cloud-computing deployments.

Major Models in Cloud Computing

This section discusses some popular models of cloud computing that are offered today as services. Although there is broad agreement on these models, there are variations based on specific vendor offerings—not surprising during these early days of cloud computing.

Software as a Service

Consider the case of an enterprise with its set of software licenses for the various applications it uses. These applications could be in human resources, finance, or customer relationship management, to name a few. Instead of obtaining desktop and server licenses for software products it uses, an enterprise can obtain the same functions through a hosted service from a provider through a network connection. The interface to the software is usually through a web browser. This common cloud-computing model is known as *Software as a Service* (SaaS) or a hosted software model; the provider is known as the *SaaS Provider*.

SaaS saves the complexity of software installation, maintenance, upgrades, and patches (for example, for security fixes) for the IT team within the enterprise, because the software is now managed centrally at the SaaS provider's facilities. Also, the SaaS provider can provide this service to multiple customers and enterprises, resulting in a multitenant model. The pricing of such a SaaS service is typically on a per-user basis for a fixed bandwidth and storage. Monitoring application-delivery performance is the responsibility of the SaaS provider. **Salesforce.com** is an example of a SaaS provider. The company was founded to provide hosted software services, unlike some of the software vendors that have hosted versions of their conventional offerings.

Platform as a Service

Unlike the fixed functions offered by SaaS, *Platform as a Service* (PaaS) provides a software platform on which users can build their own applications and host them on the PaaS provider's infrastructure. The software platform is used as a development framework to build, debug, and deploy applications. It often provides middleware-style services such as database and component services for use by applications. PaaS is a true cloud model in that applications do not need to worry about the scalability of the underlying platform (hardware and software). When enterprises write their application to run over the PaaS provider's software platform, the elasticity and scalability is guaranteed transparently by the PaaS platform.

The platforms offered by PaaS vendors like Google (with its *App-Engine*) or **Force.com** (the PaaS offering from **Salesforce.com**) require the applications to follow their own *Application Programming Interface* (API) and be written in a specific language. This situation is likely to change but is a cause for concerns about lock-in. Also, it is not easy to migrate existing applications to a PaaS environment. Consequently, PaaS sees the most success with new applications being developed specifically for the cloud. Monitoring application-delivery performance is the responsibility of the PaaS provider. Pricing for PaaS can be on a per-application developer license and on a hosted-seats basis. Note that PaaS has a greater degree of user control than SaaS.

Infrastructure as a Service

Amazon is arguably the first major proponent of *Infrastructure as a Service* (IaaS) through its *Elastic Computing Cloud* (EC2) service. An IaaS provider offers you “raw” computing, storage, and network infrastructure so that you can load your own software, including operating systems and applications, on to this infrastructure. This scenario is equivalent to a hosting provider provisioning physical servers and storage and letting you install your own OS, web services, and database applications over the provisioned machines. Amazon lets you rent servers with a certain CPU speed, memory, and disk capacity along with the OS and applications that you need to have installed on them (Amazon provides some “canned” software for the OS and applications known as *Amazon Machine Images* [AMIs], so that is one starting point). However, you can also install your own OSs (or no OS) and applications over this server infrastructure.

IaaS offers you the greatest degree of control of the three models. You need to know the resource requirements for your specific application to exploit IaaS well. Scaling and elasticity are your—not the provider’s—responsibility. In fact, it is a mini do-it-yourself data center that you have to configure to get the job done. Interestingly, Amazon uses virtualization as a critical underpinning of its EC2 service, so you actually get a VM when you ask for a specific machine configuration, though VMs are not a prerequisite for IaaS. Pricing for the IaaS can be on a usage or subscription basis. CPU time, storage space, and network bandwidth (related to data movement) are some of the resources that can be billed on a usage basis.

In summary, these are three of the more common models for cloud computing. They have variations and add-ons, including *Data Storage as a Service* (providing disk access on the cloud), communications as a service (for example, a universal phone number through the cloud), and so on.

Public, Private, and Internal Clouds

We have focused on cloud service providers whose data centers are external to the users of the service (businesses or individuals). These clouds are known as *public clouds*—both the infrastructure and control of these clouds is with the service provider. A variation on this scenario is the *private cloud*. Here, the cloud provider is responsible only for the infrastructure and not for the control. This setup is equivalent to a section of a shared data center being partitioned for use by a specific customer. Note that the private cloud can offer SaaS, PaaS, or IaaS services, though IaaS might appear to be a more natural fit.

An *internal cloud* is a relatively new term applied to cloud services provided by the IT department of an enterprise from the company's own data centers. This setup might seem counterintuitive at first—why would a company run cloud services for its internal users when public clouds are available? Doesn't this setup negate the advantages of elasticity and scalability by moving this service to inside the enterprise?

It turns out that the internal cloud model is very useful for enterprises. The biggest concerns for enterprises to move to an external cloud provider are security and control. CIOs are naturally cautious about moving their entire application infrastructure and data to an external cloud provider, especially when they have several person-years of investment in their applications and infrastructure as well as elaborate security safeguards around their data. However, the advantages of the cloud—resiliency, scalability, and workload migration—are useful to have in the company's own data centers. IT can use per-usage billing to monitor individual business unit or department usage of the IT resources and charge them back. Controlling server sprawl through virtualization and moving workloads to geographies and locations in the world with lower power and infrastructure costs are of value in a cloud-computing environment. Internal clouds can provide all these benefits.

This classification of clouds as public, private, and internal is not universally accepted. Some researchers see the distinction between private and internal clouds to be a matter of semantics. In fact, the NIST draft definition considers a private cloud to be the same as an internal cloud. However, the concepts are still valid and being realized in service provider and enterprise IT environments today.

When Does Cloud Computing Make Sense?

Outsourcing your entire IT infrastructure to a cloud provider makes sense if your deployment is a “green field” one, especially in the case of a startup. Here, you can focus on your core business without having to set up and provision your IT infrastructure, especially if it primarily involves basic elements such as e-mail, word processing, collaboration tools, and so on. As your company grows, the cloud-provided IT environment can scale along with it.

Another scenario for cloud usage is when an IT department needs to “burst” to access additional IT resources to fulfill a short-term requirement. Examples include testing of an internally developed application to determine scalability, prototyping of “nonstandard” software to evaluate suitability, execution of a one-time task with an exponential demand on IT resources, and so on. The term *cloud bursting* is sometimes used to describe this scenario. The cloud resources may be loosely or tightly coupled with the internal IT resources for the duration of the cloud bursting. In an extremely loosely coupled scenario, only the results of the cloud bursting are provided to the internal IT department. In the tightly coupled scenario, the cloud resources and internal IT resources are working on the same problem and require frequent communication and data sharing.

In some situations cloud computing does not make sense for an enterprise. Regulation and legal considerations may dictate that the enterprise house, secure, and control data in a specific location or geographical area. Access to the data might need to be restricted to a limited set of applications, all of which need to be internal. Another situation where cloud computing is not always the best choice is when application response time is critical. Internal IT departments can plan their server infrastructure and the network infrastructure to accommodate the response-time requirements. Although some cloud providers provide high-bandwidth links and can specify *Service-Level Agreements* (SLAs) (especially in the case of SaaS) for their offerings, companies might be better off keeping such demanding applications in house.

An interesting variation of these scenarios is when companies outsource their web front ends to a cloud provider and keep their application and database servers internal to the enterprise. This setup is useful when the company is ramping up its offerings on the web but is not completely certain about the demand. It can start with a small number of web servers and scale up or down according to the demand. Also, acceleration devices such as *Application Delivery Controllers* (ADCs) can be placed in front of the web servers to ensure performance. These devices provide server load balancing, *Secure Sockets Layer* (SSL) front ends, caching, and compression. The deployment of these devices and the associated front-end infrastructure can be completely transparent to the company; it only needs to focus on the availability and response time of its application behind the web servers.

Cloud Computing Infrastructure

The most significant infrastructure discussion is related to the data center, the interconnection of data centers, and their connectivity to the users (enterprises and consumers) of the cloud service.

A simple view of the cloud data center is that it is similar to a corporate data center but at a different scale because it has to support multiple tenants and provide scalability and elasticity. In addition, the applications hosted in the cloud as well as virtualization (when it is used) also play a part.

A case in point is the *MapReduce* computing paradigm that Google implements to provide some of its services (other companies have their own implementations of MapReduce). Put simply, the MapReduce scheme takes a set of input key-value pairs, processes it, and produces a set of output key-value pairs. To realize the implementation, Google has an infrastructure of commodity servers running Linux interconnected by Ethernet switches. Storage is local through inexpensive *Integrated Drive Electronics* (IDE) disks attached to each server.

Jobs, which consist of a set of tasks, are scheduled and mapped to the available machine set. The scheme is implemented through a *Master* machine and *Worker* machines. The latter are scheduled by the Master to implement Map and Reduce tasks, which themselves operate on chunks of the input data set stored locally. The topology and task distribution among the servers is optimized for the application (MapReduce in this case). Although Google has not made public the details of how the back-end infrastructure is implemented for Google Apps and Gmail, we can assume that the physical and logical organization is optimized for the tasks that need to be carried out, in a manner similar to what is done for MapReduce.

SaaS vendors can partition their cloud data center according to load, tenant, and type of application that they will offer as a service. In some cases they might have to redirect the traffic to a different data center, based on the load in the default data center. IaaS provides the greatest degree of control for the user, as discussed earlier. Even here, the topology and load assignment can be based on the number and type of servers that are allocated.

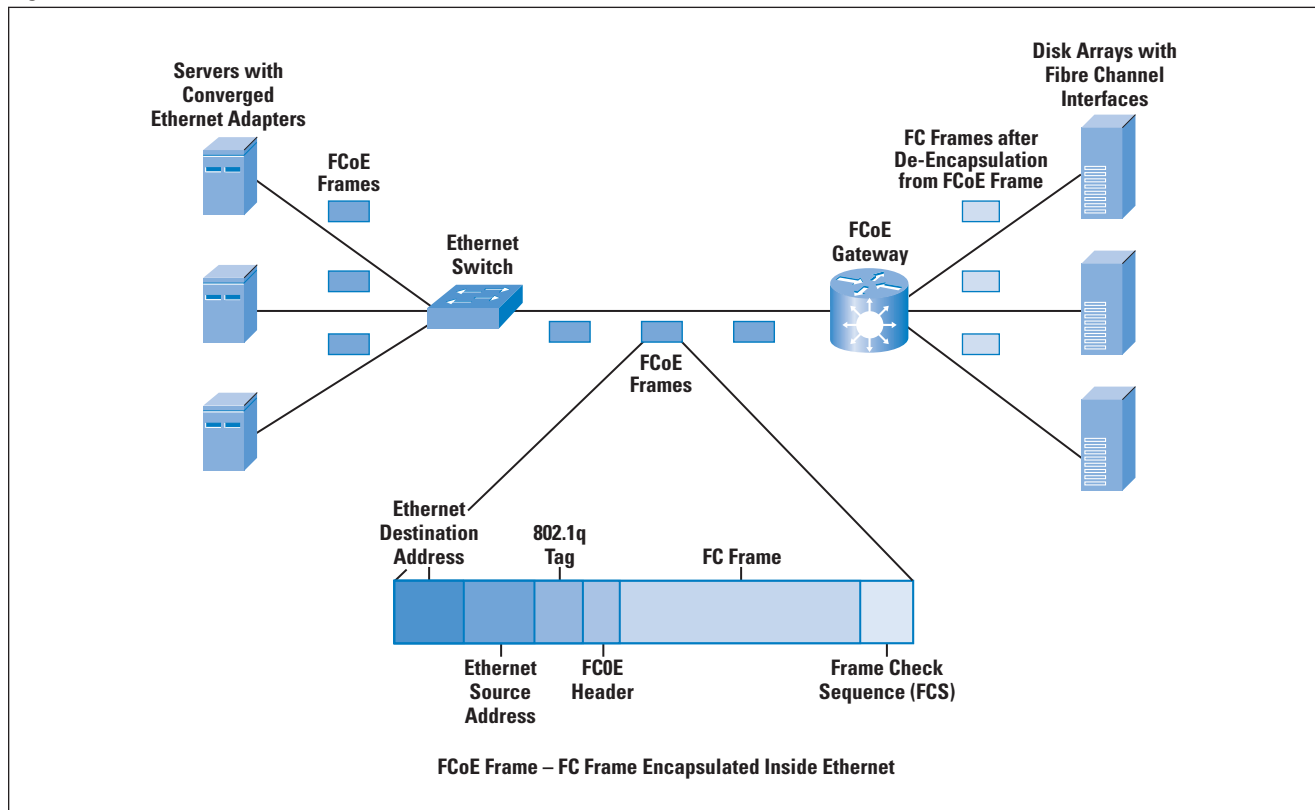
Storage Infrastructure

Storage plays a major part in the data center and for cloud services, especially in environments with virtualization. Storage can be locally attached or accessible through a network—the most popular storage network technologies being *Fibre Channel* and Ethernet. For such network access of storage, servers are equipped with Fibre Channel or Ethernet adapters through which they connect to a Fibre Channel or Ethernet switch. The switch provides the connectivity to storage arrays. Fibre Channel is more popular, though *Network Attached Storage* (NAS) devices with Ethernet interfaces also have a strong presence in the data center. Another Ethernet-based storage option is the *Internet Small Computer System Interface* (iSCSI), which is quite popular among smaller data centers and enterprises because of the cost benefits. This technology involves running the SCSI protocol on a TCP/IP-over-Ethernet connection.

Fibre Channel connections to the storage network necessitate two types of network technologies in the data center: Ethernet for server-to-server and server-to-client connectivity and Fibre Channel for server-to-storage connectivity. A recent initiative in data-center technology is a converged network, which involves the transport of *Fibre Channel over Ethernet* (FCoE). FCoE removes the need for each server to have a Fibre Channel adapter to connect to storage. Instead, Fibre Channel traffic is encapsulated inside an Ethernet frame and sent across to a FCoE gateway that provides Ethernet-to-FCoE termination to connect to Fibre Channel storage arrays (refer to Figure 3). Some storage products provide FCoE functions, so the Ethernet frame can be carried all the way to the storage array. An adapter on the server that provides both “classical” Ethernet and FCoE functions is known as a *Converged Network Adapter* (CNA). Cloud-computing environments can reduce the data-center network complexity and cost through this converged network environment.

Another area in which storage is important is in virtualization and live migration. When a VM migrates to a different physical machine, it is important that the data used by the VM is accessible to both the source and the target machines. Alternatively, if the VM is migrated to a remote data center, the stored data needs to be migrated to the remote data center too. Also, in a virtualized environment, the Fibre Channel, Ethernet, or converged adapter driver should support multiple VMs and interleave its storage traffic to the storage devices. This interleaving is done in consonance with the hypervisor and a designated VM (paravirtualized environments often use this tool), as appropriate.

Figure 3: FCoE in a Cloud Data-Center Environment



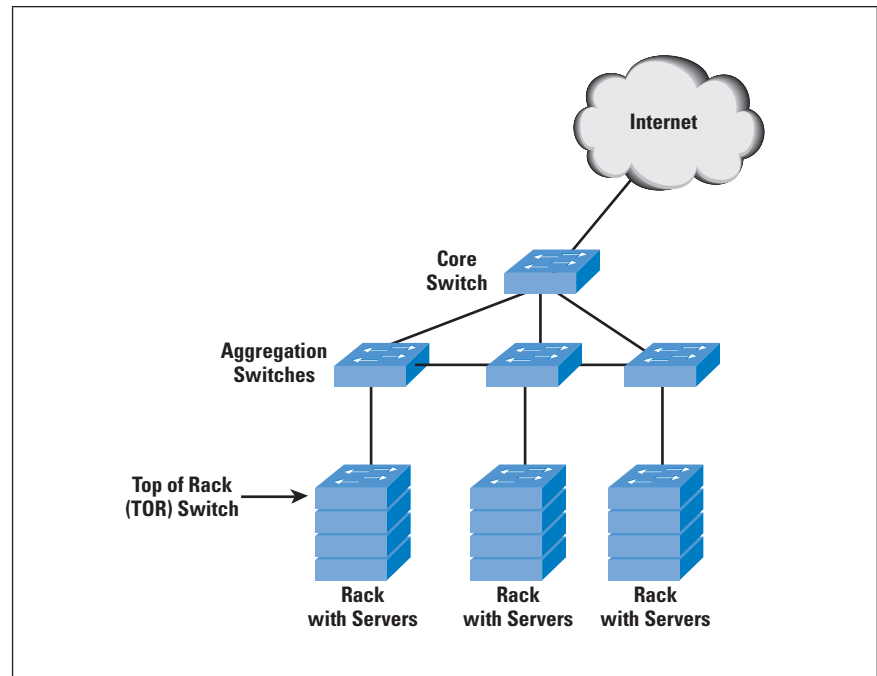
Cloud Computing: Effect on the Network

The previous discussion indicated that the network is a big part of cloud computing. A cloud user connects to the network to access the cloud resources, as indicated earlier in Figure 1. The cloud is accessible through a public network (the Internet) or through a private network (dedicated lines or *Multiprotocol Label Switching* [MPLS] infrastructure, for example). Response-time guarantees depend upon this connectivity. Some cloud vendors offer dedicated links to their data centers and provide appropriate SLAs for uptime or response time and charge for such SLAs. Others might implement a best-effort scheme but provide tools for monitoring and characterizing application performance and response time, so that users can plan their bandwidth needs.

The most significant effect on the network is in the data center, as indicated previously. Let us start with the network architecture or topology. The most common network architecture for enterprises is the three-layer architecture with access, aggregation or distribution, and core switches. The data center requires a slightly different variation to this layering, as proposed by some vendors. The data center consists mainly of servers in racks interconnected through a *Top-of-Rack* (TOR) Ethernet switch which, in turn, connects to an aggregation switch, sometimes known as an *End-of-Rack* (EOR) switch (Figure 4).

The aggregation switch connects to other aggregation switches and through these switches to other servers in the data center. A core switch connects to the various aggregation switches and provides connectivity to the outside world, typically through Layer 3 (IP). It can be argued that most of intra-data center traffic traverses only the TOR and the aggregation switches. Hence the links between these switches and the bandwidth of those links need to account for the traffic patterns. Some vendors have proposed a fat-tree or a leaf-spine topology to address this anomaly, though this is not the only way to design the data-center network. Incidentally, the fat-tree topology is not new—it has been used in *Infiniband* networks in the data center.

Figure 4: Example Data-Center Switch Network Architecture



The presence of virtualized servers adds an extra dimension. Network connections to physical servers will need to involve “fatter pipes” because traffic for multiple VMs will be multiplexed onto the same physical Ethernet connection. This result is to be expected because you have effectively collapsed multiple physical servers into a single physical server with VMs. It is quite common to have servers with 10-Gbps Ethernet cards in this scenario.

New Protocols for Data-Center Networking

Numerous initiatives and standards bodies are addressing the standards related to cloud computing. From the networking side, the IEEE is working on new protocols and the enhancement of existing protocols for data centers. These enhancements are particularly useful in data centers with converged networks—the area is often known as *Convergence Enhanced Ethernet* (CEE).

A previous section indicated the importance of FCoE for converged storage network environments. The IEEE is working to enable FCoE guarantees (because Fibre Channel is a reliable protocol as compared to best-effort Ethernet) through an Ethernet link in what is known as “Lossless Ethernet.” FCoE is enabled through a *Priority Flow Control* (PFC) mechanism in the 802.1Qbb activities in the IEEE. In addition, draft IEEE 802.1Qau provides end-to-end congestion notification through a signaling mechanism propagating up to the ingress port, that is, the port connected to the server *Network Interface Card* (NIC). This feature is useful in a data-center topology.

A third draft IEEE 802.1aq defines shortest-path bridging. This work is similar to the work being done in the IETF TRILL (*Transparent Interconnect of Lots of Links*) working group. The key motivation behind this work is the relatively flat nature of the data-center topology and the requirement to forward packets across the shortest path between the endpoints (servers) to reduce latency, rather than a root bridge or priority mechanism normally used in the *Spanning Tree Protocol* (STP). The shortest-path bridging initiative in IEEE 802.1aq is an incremental advance to the *Multiple Spanning Tree Protocol* (MSTP), which uses the *Intermediate System-to-Intermediate System* (IS-IS) link-state protocol to share learned topologies between switches and to determine the shortest path between endpoints.

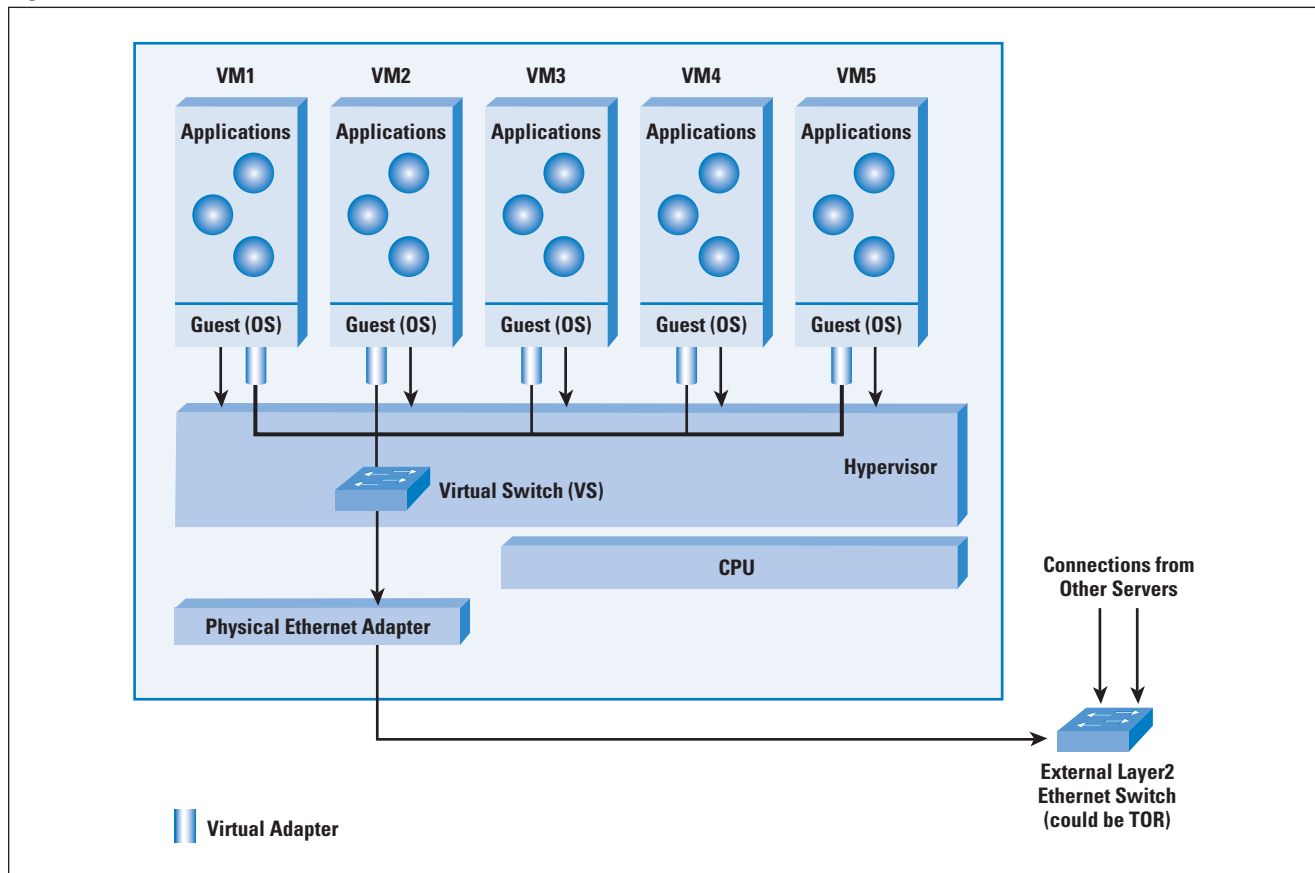
The fourth draft 802.1Qaz is also known as *Enhanced Transmission Selection* (ETS). It allows lower-priority traffic to burst and use the unused bandwidth from the higher-priority traffic queues, thus providing greater flexibility.

Virtualized Network Equipment Functions

Though cloud computing does not depend upon virtualization, several cloud infrastructures are built with virtualized servers. In an environment with physical servers, switches are used to connect servers to other servers. Firewalls and application-delivery controllers are other types of equipment that you can use in a data center on the connection to external clients. With a virtualized environment, you can move some or all of these functions to reside inside a server.

Consider the case of the software-based *Virtual Switch* as shown in Figure 5. You can use the Virtual Switch to switch between VMs inside the same physical server and aggregate the traffic for connection to the external switch. The Virtual Switch is often implemented as a plug-in to the hypervisor. The VMs have virtual Ethernet adapters that connect to the Virtual Switch, which in turn connects to the physical Ethernet adapter on the server and to the external Ethernet switch. To the network manager, the virtual switch can appear as a part of the network. Unlike physical switches, the Virtual Switch does not necessarily have to run network protocols for its operation, nor does it need to treat all its ports the same because it knows that some of them are connected to virtual Ethernet ports (for example, it can avoid destination address learning on the ports connected to the VMs). It can function through appropriate configuration from an external management entity.

Figure 5: Virtual Ethernet Switch in a Virtualized Server Environment



It is possible to implement a virtualized firewall as a VM instead of as a plug-in to the hypervisor. These VMs are self-contained, with an operating system along with the firewall software. The complete package is known as a *firewall virtual appliance*. These VMs can be loaded and configured so that network packets destined for any of the VMs pass through the firewall VM, where they are validated before being passed to the other VMs. Another use of the firewall VM is as a front end to the physical servers in the data center. The disadvantage of a virtual appliance is the performance hit due to its implementation as a software function in a virtualized environment.

Management

Management has several facets in a cloud-computing environment: billing, application-response monitoring, configuring network resources (virtual and physical), and workload migration. In a private cloud or tightly coupled environment, management of the applications may have to be shared between the internal cloud and the private cloud.

You can manage cloud-computing environments in several ways, depending upon the specific area. You can manage the network equipment (physical and virtual) through the *Simple Network Management Protocol* (SNMP) and a network management console. In a virtualized environment, the virtualization vendor often offers a framework to manage and monitor VMs, so this is another part of the equation. Several vendors offer products to act as management front ends for public clouds; for example, Amazon, whose products act as brokers and management consoles for your application deployed over the Amazon cloud offering.

It is clear that this area of management for cloud computing is still evolving and needs to be tied together for a unified management view.

Cloud Computing: Common Myths

Thus far, we have considered the important technologies, terminology, and developments in cloud computing. This section outlines some common myths about cloud computing.

- *Myth: Cloud computing should satisfy all the requirements specified: scalability, on demand, pay per use, resilience, multitenancy, and workload migration.*

In fact, cloud-computing deployments seldom satisfy all the requirements. Depending upon the type of service offered (SaaS, IaaS, or PaaS), the service can satisfy specific subsets of these requirements. There is, however, value in trying to satisfy most of these requirements when you are building a cloud service.

- *Myth: Cloud computing is useful only if you are outsourcing your IT functions to an external service provider.*

Not true. You can use cloud computing in your own IT department for on-demand, scalable, and pay-per-use deployments. Several vendors offer software tools that you can use to build clouds within your enterprise's own data center.

- *Myth: Cloud computing requires virtualization.*

Although virtualization brings some benefits to cloud computing, including aspects such as efficient use of servers and workload migration, it is not a requirement for cloud computing. However, virtualization is likely to see increased usage in cloud deployments.

- *Myth: Cloud computing requires you to expose your data to the outside world.*

With internal clouds you will never need to expose your data to the outside world. If data security and privacy are concerns, you can develop a cloud model where web front ends are in the cloud and back-end data always resides in your company's premises.

- *Myth: Converged networks are essential to cloud computing.*

Although converged networks (with FCoE, for example) have benefits and will see increased adoption in data centers in the future, cloud computing is possible without converged networks. In fact, some cloud vendors use only Fibre Channel for all their storage needs today. Use of converged networks in the future will result in cost efficiencies, but it is not a requirement today.

Cloud Computing: Gaps and Concerns

Cloud-computing technology is still evolving. Various companies, standards bodies, and alliances are addressing several remaining gaps and concerns. Some of these concerns follow:

- *Security:* Security is a significant concern for enterprise IT managers when they consider using a cloud service provider. Physical security through isolation is a critical requirement for private clouds, but not all cloud users need this level of investment. For those users, the cloud provider must guarantee data isolation and application security (and availability) through isolation across multiple tenants. In addition, authentication and authorization of cloud users and encryption of the “network pipe” from the cloud user to the service provider application are other factors to be considered.
- *Network concerns:* When cloud bursting is involved, should the servers in the cloud be on the same Layer 2 network as the servers in the enterprise? Or, should a Layer 3 topology be involved because the cloud servers are on a network outside the enterprise? In addition, how would this work across multiple cloud data centers?
- *Cloud-to-cloud and Federation concerns:* Consider a case where an enterprise uses two separate cloud service providers. Compute and storage resource sharing along with common authentication (or migration of authentication information) are some of the problems with having the clouds “interoperate.” For virtualized cloud services, VM migration is another factor to be considered in federation.
- *Legal and regulatory concerns:* These factors become important especially in those cases involving storing data in the cloud. It could be that the laws governing the data are not the laws of the jurisdiction where the company is located.

Conclusion

This article introduced the still-evolving area of cloud computing, including the technologies and some deployment concerns. Definitions and standardization in this area are a work in progress, but there is clear value in cloud computing as a solution for several IT requirements. In Part 2 we will provide a more detailed look at some of the technologies and scenarios for cloud computing.

For Further Reading

- [1] Draft NIST Working Definition of Cloud Computing,
<http://csrc.nist.gov/groups/SNS/cloud-computing/index.html>
- [2] “Identifying Applications for Public and Private Clouds,” Tom Nolle, Searchcloudcomputing,
http://searchcloudcomputing.techtarget.com/tip/0,289483,sid201_gci1358701,00.html?track=NL-1329&ad=710605&asrc=EM_NLT_7835341&uid=8788654
- [3] “The Wisdom of Clouds,” James Urquhart’s blog on Cloud Computing,
<http://news.cnet.com/the-wisdom-of-clouds/>
- [4] “Virtualization – State of the Art,” SCOPE Alliance,
<http://www.scope-alliance.org/sites/default/files/documents/SCOPE-Virtualization-StateofTheArt-Version-1.0.pdf>
- [5] “Live Migration of Virtual Machines,” Clark, et al.,
<http://www.cl.cam.ac.uk/research/srg/netos/papers/2005-migration-nsdi-pre.pdf>
- [6] “MapReduce: Simplified Data Processing on Large Clusters,” Dean & Ghemawat,
<http://labs.google.com/papers/mapreduce.html>
- [7] “Cloud Computing Drives New Networking Requirements,” *The Lippis Report*, 120,
<http://lippisreport.com/2009/02/lippis-report-120-cloud-computing-drives-new-networking-requirements/>
- [8] “A New Approach to Network Design When You Are in the Cloud,” *The Lippis Report*, 121,
<http://lippisreport.com/2009/03/a-new-approach-to-network-design-in-the-cloud/>
- [9] “Unified Fabric Options Are Finally Here,” *The Lippis Report*, 126,
<http://lippisreport.com/2009/05/lippis-report-126-unified-fabric-options-are-finally-here/>
- [10] “Virtualization with Hyper-V,” Microsoft,
<http://www.microsoft.com/windowsserver2008/en/us/hyperv-overview.aspx>
- [11] “Citrix XenServer,” Citrix,
<http://www.citrix.com/English/ps2/products/feature.asp?contentID=1686939>

- [12] “VMware Virtual Networking Concepts,” VMware,
http://www.vmware.com/files/pdf/virtual_networking_concepts.pdf
- [13] “Cisco Nexus 1000v Virtual Ethernet Switch,” Cisco Systems,
http://www.cisco.com/en/US/prod/collateral/switches/ps9441/ps9902/data_sheet_c78-492971.html
- [14] “Application Delivery Challenge,” Layland Consulting,
http://www.edge-delivery.org/dl/whitepapers/Application_Delivery_Challenge.pdf
- [15] “Cloud Networking: Design Patterns for ‘Cloud-Centric’ Application Environments,”
<http://www.aristanetworks.com/en/CloudCentricDesignPatterns.pdf>
- [16] IEEE 802.1Qaz – Enhanced Transmission Selection,
<http://www.ieee802.org/1/pages/802.1az.html>
- [17] IEEE 802.1Qau – Congestion Notification,
<http://www.ieee802.org/1/pages/802.1au.html>
- [18] IEEE 802.1Qbb – Priority Flow Control,
<http://www.ieee802.org/1/pages/802.1bb.html>
- [19] IEEE 802.1aq - Shortest Path Bridging,
<http://www.ieee802.org/1/pages/802.1aq.html>
- [20] IETF Transparent Interconnection of Lots of Links (trill) Working Group,
<http://www.ietf.org/dyn/wg/charter/trill-charter.html>

T. SRIDHAR received his BE in Electronics and Communications Engineering from the College of Engineering, Guindy, Anna University, Madras, India, and his Master of Science in Electrical and Computer Engineering from the University of Texas at Austin. He can be reached at TSridhar@leitnet.com

Why End-to-End Security Is Necessary But Not Sufficient

by Michael H. Behringer, Cisco Systems

End-to-end security relies on protocols and mechanisms that are implemented exclusively on the endpoints of a connection. The most typical example is an HTTPS connection (based, for example, on *Transport Layer Security* (TLS)^[1]) to a web server; *IP Security* (IPsec)^[2] can also be used for end-to-end security, as was initially proposed as a default connection mechanism for IPv6.

There is a perception that end-to-end security is sufficient as a security solution, and that network-based security is obsolete in the presence of end-to-end security. This article outlines why in practice end-to-end security alone is not sufficient, and why network-based security is also required.

Defining “End”

The traditional definition of an endpoint is a client or server. In this definition end-to-end security starts on the client and ends on the server. Given the multitude of applications running in parallel on an operating system, and given increasing virtualization, this definition is usually no longer precise enough. The operating system can establish a security association on either the session or application level. It can also be terminated on a front end, on behalf of numerous servers, as is the case in many TLS^[1] deployments.

Because the main goal of this article is to understand why the network has a role to play in security, the precise definition of an endpoint is not relevant here. Abstractly seen, an endpoint is an entity that communicates over a network with another entity. This definition, albeit vague, is sufficient for the discussion at hand.

End-to-End Security Is Fundamental

Security on the endpoints (client-server, or client-client for peer-to-peer) is an absolute requirement for secure communications. Such a solution contains the following components:

- *Identity*: This component encompasses known and verifiable entity identities on both ends; note that an identity can be temporary for a connection. For example, a user often is identified by username and password, whereas a server may be identified through a server certificate.
- *Protocols* (for example, TLS [1] and IPsec [2]): Protocols are used to dynamically negotiate session keys, and to provide the required security functions (for example, encryption and integrity verification) for a connection. Protocols use algorithms to implement these functions.

- *Algorithms* (for example, *Advanced Encryption Standard* [AES]^[3], *Triple Digital Encryption Standard* [3DES]^[4], and *Secure Hash Algorithm* [SHA-1]^[5]): These algorithms use the previously mentioned session keys to protect data in transit, for example through encryption or integrity checks.
- *Secure implementation*: The endpoint (client or server) that runs one of these protocols mentioned previously must be free of bugs that could compromise security. Web browser security is relevant here. Also malware can compromise security, for example by logging key strokes on a PC.
- *Secure operation*: Users and operators have to understand the security mechanisms, and how to deal with exceptions. For example, web browsers warn about invalid server certificates, but users can override the warning and still make the connection. This concern is a nontechnical one, but is of critical concern today.

For full end-to-end security, all of these components must be secure. In networks with end-to-end security, both ends can typically (depending on the protocols and algorithms used) rely on the fact that their communication is not visible to anyone else, and that no one else can modify the data in transit. End-to-end security is used successfully today, for example, in online banking applications. Correct and complete end-to-end security is required; without it, many applications such as online banking would not be possible.

However, a single security problem in any of the components can compromise the overall security for a connection. Today, most critical are implementation problems on endpoints, as well as human errors, specifically in handling exception cases.

Practical Shortcomings of End-to-End Security

Solutions that rely exclusively on end-to-end security have many potential problems, which fall into two broad categories: those that affect the end user and those that affect the network operator (the service provider, or the enterprise network operator, for example).

The End-User View

As reports on online crime and fraud demonstrate very clearly, even in the perceived presence of end-to-end security it is difficult to ensure that none of the components mentioned previously is “broken.” Although protocols and algorithms in use tend to be secure and reliable, the main problems lie in the two main areas of endpoint security (secure implementation component) and lack of user education (secure operation component).

Endpoint security concerns include the presence of malware, as well as bugs in software. Even security professionals have difficulty determining whether a PC contains malware. Such malware can control the connection before it is secured, thereby achieving the ability to see the data, as well as potentially change it in real time. Although endpoint security software such as antivirus solutions as well as zero-day prevention solutions provides good security, they are not always installed, and antivirus software is often not up-to-date. Users also can temporarily disable the solutions. Therefore, the presence of malware remains a security concern. Bugs in software are also relevant, for example in the web browser or the operating system.

The lack of user education is the other important concern on the endpoint: Users must know how to identify a secured connection, for example by the little padlock in a web browser (although not even this security mechanism is completely secure). They must also know how to deal with exceptions such as expired or invalid certificates. Most average users do not entirely understand all these details, leading to breaches of security.

The Network Operator View

In the early days of IPv6 it was postulated that the protocol would come with IPsec end-to-end security built in and always “on,” thereby eliminating all security problems. This assumption turned out to be wrong, because many problems remain on the network side—for example, general problems with end-to-end security—and they apply to all variants, such as IPsec, TLS, or *Secure Sockets Layer* (SSL).

Today, most enterprise network operators as well as service providers are skeptical about the ubiquitous use of end-to-end security solutions. The fundamental concern is that the endpoints generally cannot be trusted. The network operator, whether enterprise, university, or service provider, has an obligation to enforce certain policies on the endpoint, for example, to ensure that it does not spread worms, send spam mail, or attack servers. If, however, network operators cannot “see” the traffic of an endpoint because it is end-to-end secured, then they cannot comply with their obligations to control the endpoints.

From a network operator’s perspective it is therefore not generally desirable to use end-to-end security for all communications, but only for those that really need it.

Why Network-Based Security Is Essential

There are many examples where network-based security is essential, and where end-to-end security solutions not only do not help, but may actually present an additional problem. In all those cases it is essential to have strong network-based security solutions in place. Some examples explain this in more detail.

The Service Provider with DSL Customers

A service provider with DSL customers needs to control its users' traffic in various ways. However, the provider has no control over the endpoints, because those are the customers' property. Because they also cannot force their customers to use appropriate security software, there is always a certain percentage of infected PCs on any given service provider's network. Critical service provider concerns follow:

- *Control of PCs infected with malware:* Such PCs (also referred to as “bots” or “zombies”) can infect other PCs and participate in illegal activities, such as spam mail, click fraud^[12], Denial-of-Service (DoS) attacks, etc. There is a strong, often legal requirement for providers to identify such infected PCs, to isolate them, and to alert their owners and help them to “disinfect” the PC. Network-based security mechanisms are required, essentially because security on the endpoint has failed.
- *Attacks from the users:* Even in the absence of malware, a service provider's user can participate in illegal activities, such as DoS attacks, or intrusions on web servers or routers. Network-based methods are required to detect such attempts, beginning with simple forms such as IP spoofing [6], and to prevent or block them. One example is network-based solutions against DoS attacks^[7,8].
- *Control of bandwidth:* Many service providers need to enforce bandwidth limits on some applications or users because they violate service agreements. Also here, applications are necessary to control the PCs, and to limit their usage of the service to remain within contracted boundaries. Service providers today employ a large number of network-based security mechanisms, ranging from visibility solutions to enforcement of certain policies. Endpoint security does not solve these problems, because the PC is not under control of the service provider, and is typically untrusted.
- *Services:* Service providers also try to differentiate themselves from their competition by offering managed services, for example managed security services^[9]. Those services are also network-based, and they complement endpoint security solutions that their customers use.

The Service Provider with Customers Under Attack

Service providers may also be required to help their customers when they are under attack. DoS attacks illustrate why endpoint security may not be sufficient, and network-based security is required. Under a DoS attack, a web server, for example, may receive more traffic than it can handle. Such attacks can also overload network resources, such as subscriber lines or routers; therefore, endpoint security is not able to solve such attacks. Massive overprovisioning would be the only way to handle DoS attacks, but this approach is commercially not generally feasible. Network-based solutions based on flow analysis and selective discard of flows are required to help in such situations.

The Enterprise Network

At first glance it seems that enterprises should have full control over the PCs in the enterprise. In such a case, it would be possible to rely completely on end-to-end security. However, this assumption is unrealistic. Numerous current shortcomings make this approach impractical today:

- Enterprise PCs can also get infected with malware, leading to the same problem as for service providers described previously: the need to monitor and control the behavior of a PC in the network. Solutions to control endpoints are themselves network-based; for example, network endpoint assessment^[10] and user authentication (802.1x)^[11].
- Attacks from users, or against services within the enterprise, also exist in an enterprise environment, as explained previously for service providers. Solutions are network-based.
- The enforcement of *Quality of Service* (QoS) is also a security concern: Users could wrongly classify all their traffic as “high-priority.” In the absence of full application control on the PC (which is impractical today), the network needs to control flows from the PC, and potentially enforce a QoS policy. If all flows were encrypted end-to-end, this control would be “blind,” probably leading to undesired results. Network security mechanisms are required to control the QoS policy.
- Scale: In an enterprise with several offices that are connected over an untrusted network (for example, the Internet), it may be impractical today to roll out full end-to-end security across the entire enterprise. The currently used approach in most enterprises is to connect the offices with IPsec gateways, and leave traffic within an office in the clear. This scenario increases manageability and scalability of the network. Again, this solution is network-based security solution.
- Although PCs can theoretically be equipped with IPsec (for example) for all communications, many end devices in an enterprise do not support the security mechanisms required. Printers, faxes, and scanners are examples. Full end-to-end security, however, would require all endpoints to support a common mechanism, such as IPsec or TLS. Until all such devices have this support, network-based mechanisms are required to secure communications with them.

Summary

End-to-end security protocols and solutions are an essential cornerstone in network security. We cannot live without them. However, it is unrealistic in today's networks to assume that end-to-end security solutions alone will suffice. The fundamental underlying problem is that typically the network operator, where a PC is attached, has a need and often an obligation to monitor the behavior of the endpoint, and to control malicious activities emerging from that PC. All solutions to control endpoints, however, are by definition network-based. Therefore, network-based security mechanisms are also an essential component of overall network security: Overall security requires both endpoint security and network-based security.

References

- [1] T. Dierks, et al., "The Transport Layer Security (TLS) Protocol Version 1.2," RFC 5246, August 2008.
- [2] S. Kent, et al., "Security Architecture for the Internet Protocol," RFC 4301, December 2005.
- [3] Joan Daemen and Vincent Rijmen, *The Design of Rijndael: AES—The Advanced Encryption Standard*, Springer-Verlag, 2002. ISBN 3-540-42580-2.
- [4] ANSI X9.52:1998, "Triple Data Encryption Algorithm Modes of Operation," July 1998.
- [5] FIPS 180-2, "Secure Hash Standard (SHS)," February 2004.
- [6] F. Ali, "IP Spoofing," *The Internet Protocol Journal*, Volume 10, No. 4, December 2007.
- [7] W. Eddy, "Defenses Against TCP SYN Flooding Attacks," *The Internet Protocol Journal*, Volume 9, No. 4, December 2006.,
- [8] C. Patrikakis, et al., "Distributed Denial of Service Attacks," *The Internet Protocol Journal*, Volume 7, No. 4, December 2004.
- [9] K. Trivedi and D. Holloway, "Secure Multivendor Networks," *The Internet Protocol Journal*, Volume 10, No. 3, September 2007.
- [10] P. Sangster, et al., "Network Endpoint Assessment (NEA): Overview and Requirements," RFC 5209, June 2008.
- [11] IEEE 802.1X "Port-Based Network Access Control,"
<http://www.ieee802.org/1/pages/802.1x.html>

- [12] According to Wikipedia: “Click fraud is a type of Internet crime that occurs in pay per click online advertising when a person, automated script or computer program imitates a legitimate user of a web browser clicking on an ad, for the purpose of generating a charge per click without having actual interest in the target of the ad’s link. Click fraud is the subject of some controversy and increasing litigation due to the advertising networks being a key beneficiary of the fraud.

Use of a computer to commit this type of Internet fraud is a felony in many jurisdictions, for example, as covered by *Penal Code 502* in California, USA, and the *Computer Misuse Act 1990* in the United Kingdom. There have been arrests relating to click fraud with regard to malicious clicking in order to deplete a competitor’s advertising budget.”

http://en.wikipedia.org/wiki/Click_fraud

MICHAEL H. BEHRINGER works at Cisco Systems as a distinguished engineer, where he focuses on core security problems, such as MPLS security, multicast security, and Denial-of-Service attack prevention. Michael holds a diploma in computer science from the Technical University of Munich. He is an active member of the IETF, and has published several papers, RFCs, and a book about MPLS VPN security. E-mail: mbehring@cisco.com

Letter to the Editor

End of Eternity

Dear Ole,

In their “The End of Eternity” articles, (IPJ Volume 11, No. 4 and Volume 12, No. 1) Niall Murphy and David Wilson provide a detailed and compelling description of the lasting harm that could result from the exhaustion of unallocated IPv4 addresses—harm to Internet users and aspiring new entrants, to technical-coordination and fault-management mechanisms, and to the likely irreplaceable cooperative decision-making and consensus-development mechanisms that distinguish the Internet from every other important transnational sphere of activity in human history. Thankfully, the authors foresee a potential happy ending—or at least yet another chapter in the story—in “an IPv6 Internet, or at least enough of one to keep off address scarcity for a workable subset of the industry.”

However, having foreshadowed how they expect the IP addressing cliffhanger to be resolved, the authors go on to detail a variety of interesting but considerably less persuasive assumptions and predictions, all based on the *stipulation* that establishing IPv4 address markets would represent the best means to “shorten the gap” between the end of IPv4 and the return to a “normal” state of Internet growth and development, that is, one that is unconstrained by IP address-related scarcity (or at least no more constrained than it has been over the last decade-plus of CIDR and hierarchical interdomain routing).

I believe that it is worth highlighting here the logic that binds these two engaging and well-written articles together into something that is, unfortunately, substantially less than the sum of its parts. If the authors are to be taken at their word that “an IPv6 Internet” represents the only currently feasible and also *satisfactory* conclusion to “the IPv4 end game,” then that conclusion does not by itself entail that IPv4 markets are the only, or most obvious or effective—or even *workable*—candidate mechanisms for coordinating the distribution of IP addressing in the run-up to more widespread IPv6 adoption. And yet, that postulate is offered, without explanation or defense, as the grounding justification for an investigation of various optional features and collateral effects that the foretold IPv4 address market might have.

Many observers have committed untold pages and pixels to the exploration of hypothetical IPv4 address markets, both in IPJ and elsewhere, going back as far as RFC 1744 (1994). The two articles by Murphy and Wilson represent valuable additions to that growing corpus. However, to my knowledge, no other writings in this area have built on the proposition that IPv6 is indispensable; therefore, IPv4 addresses should be privately traded. To put it in the most generous possible terms, this claim is highly contestable. As separate and independent analyses, IPJ readers may derive many useful insights from these two articles, but attributing any special relevance to those insights based on any presumptive connection between IPv4 markets and the future necessity or viability of IPv6 would be a mistake.

—Tom Vest, Consultant
tvest@eyeeconomics.com

CSNET Receives 2009 Postel Service Award

The *Internet Society* (ISOC) has awarded the *Jonathan B. Postel Service Award* for 2009 to CSNET, the *Computer Science Network*, a research networking effort that during the early 1980s provided the critical bridge from the original research undertaken through the ARPANET to the modern Internet.

The award recognizes the pioneering work of the four principal investigators that conceived and later led the building of CSNET—Peter J. Denning, David Farber, Anthony C. Hearn and Lawrence Landweber—and the U.S. National Science Foundation program officer and visionary responsible for encouraging and funding CSNET—Kent Curtis.

Stephen Wolff, a past recipient of the Postel Award, said, “CSNET was a critical link in the transition from the research-oriented ARPANET to today’s global Internet. CSNET also helped lead the way by sharing technologies, fostering connections, and nurturing the worldwide community that provided a foundation for the global expansion of the Internet.”

ISOC presented the award, including a US\$20,000 honorarium and a crystal engraved globe, during the 75th meeting of the *Internet Engineering Task Force* (IETF) in Stockholm, Sweden. The awardees have requested that the ISOC present the honorarium to non-profit organizations they believe support the spirit of the award.

Lynn St. Amour, President and CEO of the ISOC, said “In many ways, CSNET helped set the stage for the Internet that today reaches more than 1 billion people. CSNET’s community-driven, self-sustaining governance structure was an early example of the model that helps ensure that even as today’s Internet grows and evolves, it remains an open platform for innovation around the world.”

CSNET began in 1981 with a five-year grant from the U.S. *National Science Foundation* (NSF). Five years later, CSNET connected more than 165 academic, government and industrial computer research groups comprised of more than 50,000 researchers, educators and students across the United States and around the world. It had concluded a seminal resource sharing agreement with the ARPANET and was self-governing and self-supporting. Open to all computer researchers, it demonstrated that researchers valued the kind of informal collaboration it made possible. CSNET’s success was critical to the decision by NSF in 1986 to adopt the Internet technology for NSFNET, the network backbone to connect its supercomputing centers and their research communities. CSNET provided software, policies, and experienced alumni to the NSFNET teams. NSFNET became the first backbone of the modern Internet.

The CSNET architecture supported the Internet standards, SMTP and TCP/IP, and a variety of connection protocols including telephone dialup, X.25, and ARPANET. This architecture, along with strong technical support, enabled participants of differing means and skill levels to all join the community. CSNET pioneered the model of university, industry, government partnerships that were key to the pre-commercial Internet.

The CSNET proposal was assembled by a lengthy community consensus process that began in 1979. The four principal investigators, who led this effort and served as the project's management committee, were:

Peter Denning was head of the computer science department at Purdue University. His team included professor Douglas Comer, who was responsible for the software that ran TCP/IP over the GTE Telenet X.25 commercial packet network.

David Farber was a professor of electrical engineering at University of Delaware. His team included then graduate student David Crocker, who was responsible for Phonenet, dial-in telephone connections to relay servers for e-mail exchange.

Anthony Hearn was head of the information sciences department at RAND. His team included Michael O'Brien, who was responsible for the relays connecting CSNET and ARPANET.

Lawrence Landweber was a professor of computer science at the University of Wisconsin. His team included professor Marvin Solomon and Michael Litzkow who were responsible for the name server, a precursor of modern Directory Services.

At the NSF, the late *Kent Curtis* helped conceive the entire effort and, with assistance from Bill Kearn, saw it through its formative years. He was recognized for his pivotal role by the Computing Research Association's first distinguished service award in 1988.

The *Jonathan B. Postel Service Award* was established by the Internet Society to honor individuals or organizations that, like Jon Postel, have made outstanding contributions in service to the data communications community. The award is focused on sustained and substantial technical contributions, service to the community, and leadership. With respect to leadership, the nominating committee places particular emphasis on candidates who have supported and enabled others in addition to their own specific actions. Previous recipients of the Postel Award include Jon himself (posthumously and accepted by his mother), Scott Bradner, Daniel Karrenberg, Stephen Wolff, Peter Kirstein, Phill Gross, Jun Murai, Bob Braden and Joyce K. Reynolds (jointly), Nii Quaynor, and La Fundación Escuela Latinoamericana de Redes (EsLaRed). The award consists of an engraved crystal globe and a US\$20,000 honorarium. For more information about the award, visit: <http://www.isoc.org/postel>

ISOC is a non-profit organization founded in 1992 to provide leadership in Internet related standards, education, and policy. ISOC is dedicated to ensuring the open development, evolution, and use of the Internet for the benefit of people throughout the world. More information is available at: <http://www.isoc.org>

NRO Declaration on RPKI

The *Number Resource Organization* (NRO) recently declared: “Over several years, a set of mechanisms has been under development for digital certification of Internet number resources, through a so-called *Resource Public Key Infrastructure*, or “RPKI.” Like other PKIs, the RPKI requires one or more root authorities, to act as so-called *trust anchors* for one or more certification hierarchies.^[1]

The RPKI architecture has been designed to allow a number of trust anchor configurations involving: either a single trust anchor located at the root of a single certification hierarchy; a set of independent trust anchors to be located at the roots of several independent hierarchies; or a hybrid of these. The alternative models may have advantages and disadvantages in various dimensions including: operational efficiency; alignment with resource allocation hierarchies; centralisation vs distribution of functions; recognised global or regional authority; and, operational capacity of the respective host organisations.

The *Regional Internet Registries* (RIRs) believe that the optimal eventual RPKI configuration involves a single authoritative trust anchor. That configuration may not be achievable in the short-term and the details and timelines for its implementation will depend among other things on discussions within the RIRs’ communities and dialogues with others including the *Internet Architecture Board* (IAB) and the *Internet Engineering Task Force* (IETF).

In the meantime, the RIRs have agreed to undertake pragmatic implementations of RPKI services based on interim trust anchor models, such as, self-signed trust anchors. All such implementations will comply with the overall RPKI architecture. The implementations will also have the ability to evolve into a single trust anchor model and to provide robust and fully operational (and inter-operational) services for those who wish to use them. The objective is for all RIRs to be ready to start issuing certificates by no later than January 1, 2011.

The RIRs will continue working with and receiving feedback from their respective communities and industry partners to ensure effective ongoing evolution of the RPKI system.”

For more information about the NRO, see <http://www.nro.net/>

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

ARIN Hosts 4-byte ASN Wiki

The *American Registry for Internet Numbers* (ARIN) has created a wiki to focus on issues related to 4-byte *Autonomous System Numbers* (ASNs)^[2]. This wiki provides a central repository for ongoing discussion and information exchange associated with 4-byte ASN topics and issues. The wiki can be found at: www.get4byteasn.info

Ongoing Internet growth is rapidly depleting the existing pool of 2-byte ASNs (65,536 numbers in total). As a result, the IETF has approved the expansion of AS Numbers from 2-bytes to 4-bytes, to include over 4 billion ASNs. Following a globally coordinated policy, ARIN and the other RIRs began assigning 4-byte ASNs by request in January 2007 and by default in January 2009. However, some routers do not support the use of these 4-byte ASNs.

ARIN has set up this wiki to help educate the community about 4-byte ASN operational issues, to help vendors understand how to provide 4-byte ASN support in their products and to help network operators find those products. A wide range of community stakeholders will be able to share and benefit from information contributed to the wiki. ARIN looks forward to participation from everyone, including users, ISPs, and vendors, with interest in this topic.

Upcoming Events

The *North American Network Operators' Group* (NANOG) will meet in Dearborn, Michigan, October 18–21. Following the NANOG meeting, the *American Registry for Internet Numbers* (ARIN) will meet in the same venue October 21–23. For more information see: <http://nanog.org> and <http://arin.net>

The *Internet Engineering Task Force* (IETF) will meet in Hiroshima, Japan, November 8–13, 2009 and in Anaheim, California, March 21–26, 2010. For more information see: <http://www.ietf.org/meeting/>

The *Internet Corporation for Assigned Names and Numbers* (ICANN) will meet in Seoul, Korea, October 25–30, 2009 and Nairobi, Kenya, March 7–12, 2010, and in Brussels, Belgium, June 21–25, 2010. For more information, see: <http://icann.org/>

The *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) will meet in Kuala Lumpur, Malaysia, February 23–March 5, 2010. For more information see: <http://www.apricot2010.net/>

References

- [1] Huston, Geoff, “Resource Certification,” *The Internet Protocol Journal*, Volume 12, No. 1, March 2009.
- [2] Huston, Geoff, “Exploring Autonomous Systems Numbers,” *The Internet Protocol Journal*, Volume 9, No. 1, March 2006.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2009 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

December 2009

Volume 12, Number 4

A Quarterly Technical Publication for
Internet and Intranet Professionals

In This Issue

From the Editor	1
Cloud Computing.....	2
SSH.....	18
Book Review.....	30
Fragments	38

FROM THE EDITOR

In our last issue we brought you Part 1 of a two-part article on *Cloud Computing*. T. Sridhar introduced various aspects of cloud computing, including the rationale, underlying models, and infrastructures. Part 2, subtitled “Infrastructure and Implementation Topics,” is included in the current issue. Cloud computing has received a great deal of press in recent months and continues to be an area of rapid development. I’m confident that we will have more articles about this topic in future editions of IPJ.

With this issue we start a new series of articles under the general heading “Protocol Basics.” The idea is to present a series of in-depth tutorials on numerous protocols that are used every day on the Internet and in enterprise networks. The articles will cover protocol details as well as implementation, deployment, and usage scenarios. In some cases the articles will also summarize the “lessons learned” and present “best-practice” guidelines. To start the series, we asked Bill Stallings to give us an overview of the *Secure Shell* (SSH) Protocol. We invite you to send us suggestions for other protocols that you’d like to see covered in this series.

Today’s Internet is a result of many years of technological development and innovative uses of the resulting infrastructure. Of equal importance has been many *policy* choices made over the years, ranging from what protocols to use to how to allocate finite resources such as the IPv4 address space. A new book, *Protocol Politics: The Globalization of Internet Governance*, explores some of this history. The book is examined in an extended review by Tom Vest.

Let me remind you that we will no longer be automatically extending your subscription when it expires. Please take a moment to check your expiration date (printed on the back of the journal for subscribers in the United States, and on the envelope for our international subscribers). Visit the “Subscriber Services” section of our webpage at www.cisco.com/ipj to update or renew your subscription. You can also contact us by e-mail to ipj@cisco.com regarding any aspect of your subscription.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

Cloud Computing—A Primer

Part 2: Infrastructure and Implementation Topics

by T. Sridhar

Cloud computing is an emerging area that affects IT infrastructure, network services, and applications. In Part 1^[0] of this two-part article, we introduced various aspects of cloud computing, including the rationale, underlying models, and infrastructures. In Part 2 we discuss specific infrastructure aspects of cloud computing in detail, specifically:

- Network Infrastructure
- Cloud-to-Cloud and Federation Considerations
- Security

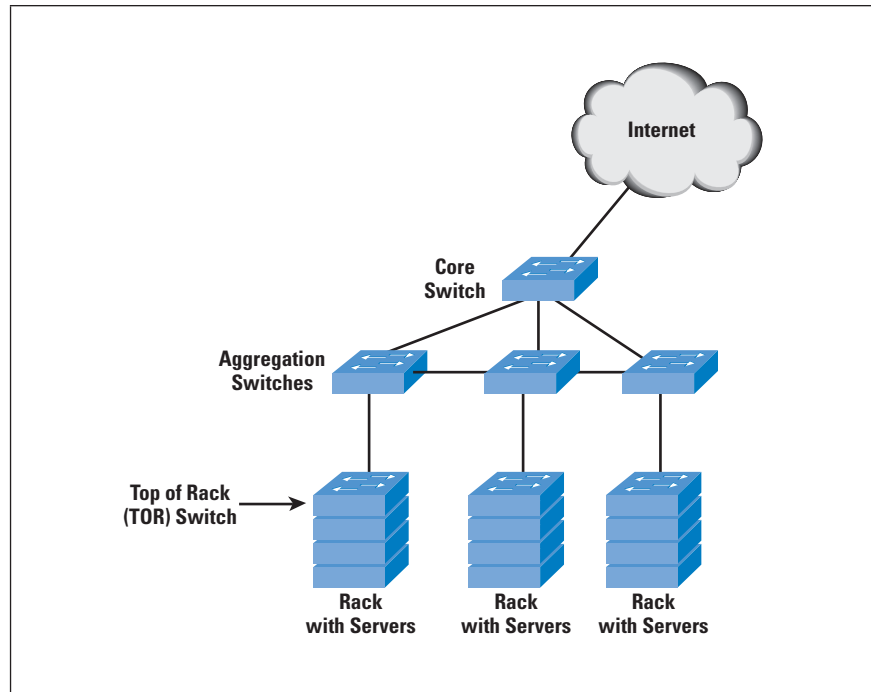
In addition, we will provide some perspective on select topics in cloud computing that have garnered interest. Remember that cloud computing is an emerging area where approaches to some of these topics are still evolving. In addition, although cloud computing is not intrinsically dependent upon virtualization, there is common agreement that virtualization (specifically, server virtualization) will be an integral part of cloud-computing solutions of the future. Consider the discussion in the following sections in this context.

Network Infrastructure

In a limited sense, the cloud can be treated as a large data center run by an external entity providing the capability for elasticity, on-demand resources, and per-usage billing. Data-center architecture often follows the common three-layer network topology of access, aggregation, and core networks with enabling networking elements (switches and routers). Consider the topology shown in Figure 4 of Part 1, reproduced here as Figure 0. The servers can be connected through a 1-Gbps link to a *Top of Rack* (TOR) switch, which in turn is connected through one or more 10-Gbps links to an aggregation *End of Row* (EOR) switch. The EOR switch is used for interserver connectivity across racks. The aggregation switches themselves are connected to core switches for connectivity outside the data center.

From a functional perspective, data-center server organization has often adopted a three-tier architecture (a specific case of an N-tier architecture). The three-tier functional architecture has a web or *Presentation Tier* on the front end, an *Application Tier* to perform the application and business-processing logic, and finally a *Database Tier* (to run the database management system), which is accessed by the Application Tier for its tasks (refer to Figure 1 on page 4).

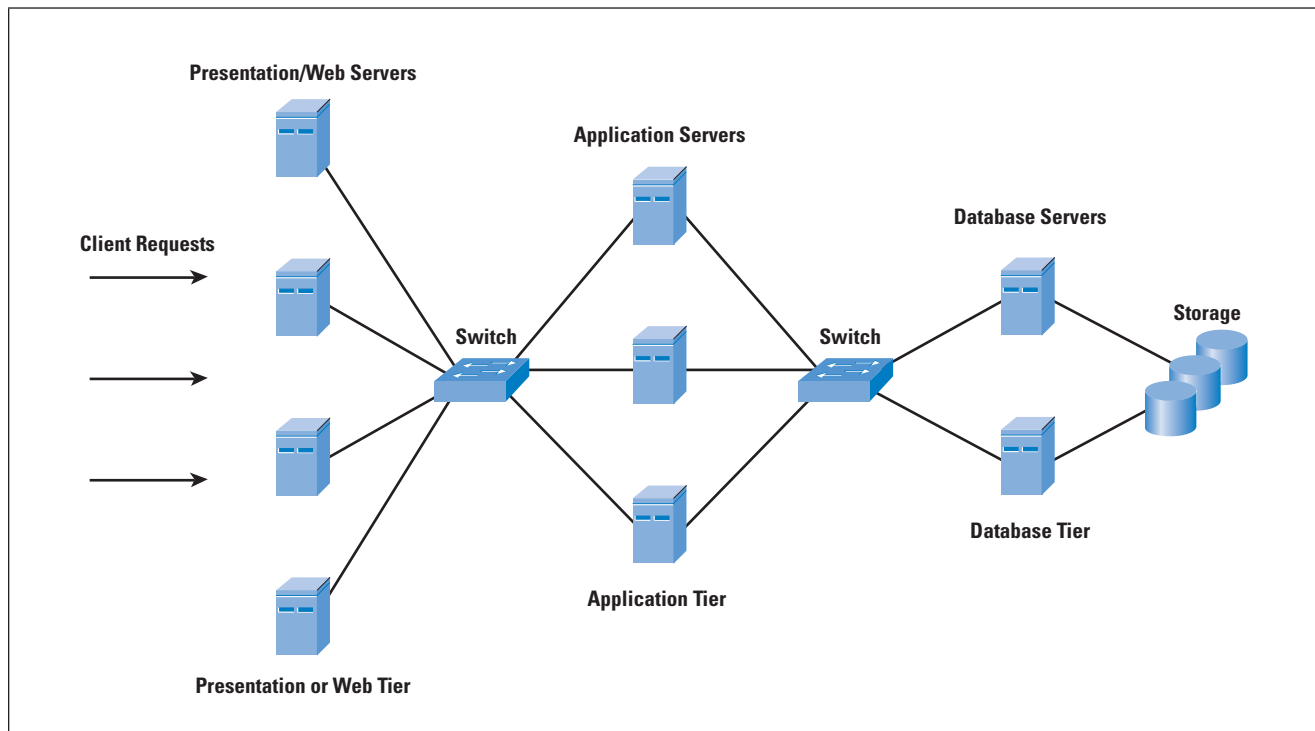
Figure 0: Example Data-Center
Switch Network Architecture
(from Part 1)



Although it is not necessary for each tier to be represented by its own physical servers (for example, you could have the Application and Database functions mapped into a single physical server), it is a common representation. The reason for this multitiered design is to control the connections and interactions, as well as for scaling and security. It is not uncommon for the Presentation Tier to be in a *Demilitarized Zone* (DMZ) while the other tiers are located deep inside the data center. Although all tiers could connect to storage for performing their functions, the Database Tier is the one with the maximum storage bandwidth requirements.

It follows that the server connectivity and the network topology for the cloud data centers might follow a similar organization. If you are an enterprise, you can perform the same business functions as before, but by using the external cloud. The choice of servers, software loads, and their interconnection will depend upon what you need to accomplish. In the following sections, we discuss how this design is handled in *Infrastructure as a Service* (IaaS), *Platform as a Service* (PaaS), and *Software as a Service* (SaaS).

Figure 1: Three-Tier Functional Server Architecture



Data-Center Infrastructure Extension – IaaS

If the cloud is thus seen as an extension of the existing data center, IaaS as outlined in Part 1 is a natural fit. Here, you would specify the number of servers in each tier, load the appropriate server image (web, business logic, or database manager), and “connect” them (through a menu or *Application Programming Interface* [API] provided by the IaaS provider) by specifying the links between them. You can also specify the network connectivity at this time (more on this later). For an enterprise IT administrator, this model provides the greatest degree of control and, to an extent, a familiar operating topology. The cloud provider handles the elasticity by ensuring that the number of servers and switches is adequate for you to configure and connect in the specified topology. Per-use billing and on-demand resource addition and removal are also provided by the cloud provider. Note that if you have complete control, you also are responsible for security, application usage, and resource management.

PaaS and SaaS Infrastructure

In the case of PaaS, you transfer more control to your cloud service provider. The platform used to build the service you require can scale transparently without any of your involvement other than at the time of configuration. You do not need to understand the tier connectivity, bandwidth requirements, or how it all functions under the hood.

Cloud service providers can realize this function—often with a three-tier topology similar to that for traditional data centers. However, some of them have innovated to perform parts of the function differently. For example, the database functions may rely upon a model of *scaling out* (splitting the database across multiple servers) instead of *scaling up* (increasing the capability of the machine running the database servers). Their claim is that with clouds involving large amounts of data that you can partition and work on, it is easier to scale out than scale up. According to some cloud service providers, traditional relational databases are not suitable candidates for scale-out. Hence, some cloud vendors have provided their own database models and implementations—a common one being the type known as the *Key-Value database*.

SaaS vendors have the highest degree of control among the three models. The realization of the network topology can be similar to existing data centers and scale up or down according to the number of users that are added. However, because they offer a specific set of applications to the cloud users, their server and network topology is quite straightforward.

For the following discussions, we will use IaaS as the representative cloud service model, with a primary consideration being “cloud bursting”—how an existing IT infrastructure can take advantage of the power of the cloud when it needs additional resources. Note that some of the discussion might also be relevant for internal clouds. In addition, we will assume a virtualized server infrastructure for the IaaS cloud because this infrastructure provides a greater degree of flexibility for cloud service providers (Amazon being a key example).

Virtualization and Its Demands on Switching

In Part 1, we provided the context for a virtual switch within a physical server containing multiple virtual machines. There are some addressing and control factors to consider in this model. Consider a data center with 100 servers, each with 16 virtual machines but with one physical 10-Gbps Ethernet connection to the external switch from each physical machine. If we were to carry forward the model where each physical server is replaced with its virtual equivalent but still needs to be addressable (through a *Media Access Control* [MAC] layer address and an IP address), you would need 16 MAC and IP addresses for the virtual servers that now reside “on top” of the single physical link, for a total of 1600 addresses across all servers. This problem is exacerbated when you increase the number of VMs per server. Switching between MAC addresses belonging to the virtual machines is done by the virtual switch inside the server.

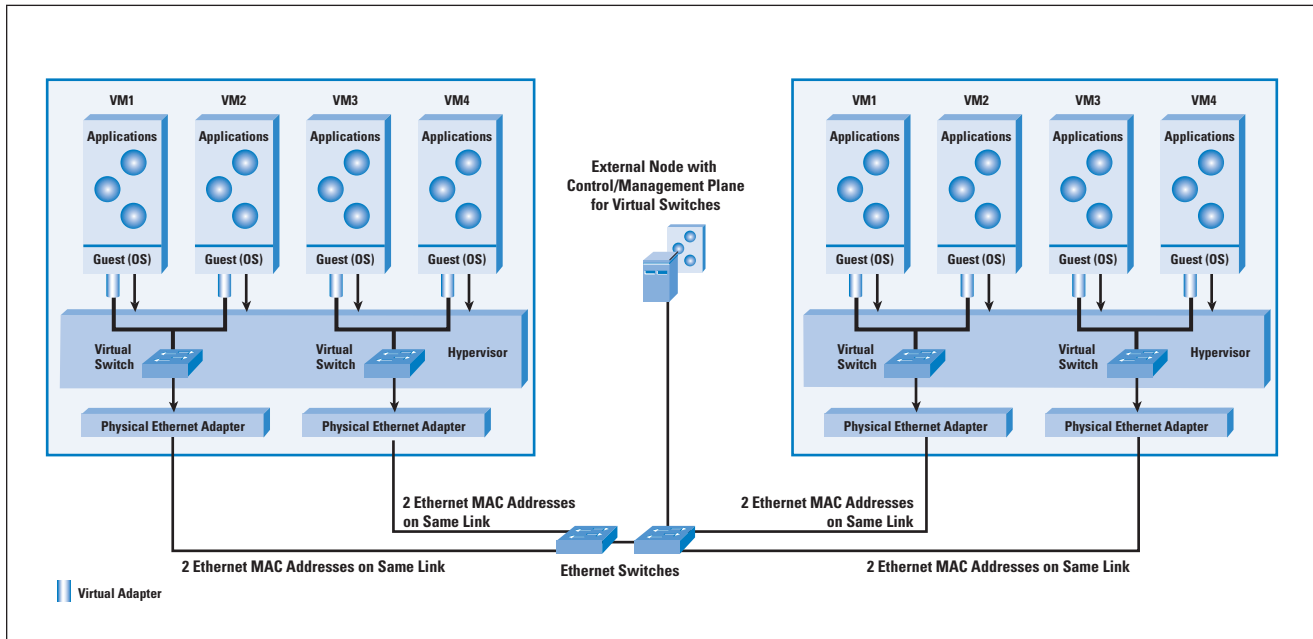
Consider the topology in Figure 2. The virtual switch treats the physical link as an uplink to the external physical switch. This intra-machine *Virtual Machine* (VM) switch with an uplink to the external switch is completely in line with access and aggregation switch topologies where the access layer is subsumed inside the server. Note that each physical host can have more than one virtual switch to support greater logical segmentation. In such cases, it is common for each of the virtual switches to have its own physical uplink to the external Ethernet switch.

The virtual switch does not need to learn MAC addresses like a traditional switch—it assumes that all destination-unknown frames should be forwarded over the physical link (or uplink to the physical switch). In addition, it switches traffic between the intramachine VMs according to policy. For example, you could prohibit two VMs on the same machine from communicating with each other by configuring an access control list on the virtual switch. The VMs may all be on the same or on different VLANs. Broadcasts and intra-VLAN traffic are forwarded according to the rules for each VLAN. In effect, the virtual switch is a simple function that is used for aggregation and access control within a physical server containing VMs.

Management of these virtual switches can follow an aggregation model—where multiple virtual switches are managed through an external node (physical machine or VM), as shown in Figure 2. This external node provides the management view on behalf of the switches. Often, the external node can run control-plane protocols for Layer 2/3 functions, in effect appearing like a control or management plane with multiple data-plane instances (the virtual switches). When VMs need to be migrated to other physical servers, this separation of control- or management-plane functions permits easier migration of policy and access lists.

Virtual switches do have some disadvantages. Inter-VM traffic within the same machine is not visible to the network and cannot be subject to appropriate monitoring by network administrators. The IEEE is discussing approaches to providing external network switches the visibility into the intra-VM traffic. The options include “hair pinning,” where inter-VM traffic would still be carried over to an external switch and brought back to the same physical server.

Figure 2: Virtual Switch Aggregation and Management by External Node



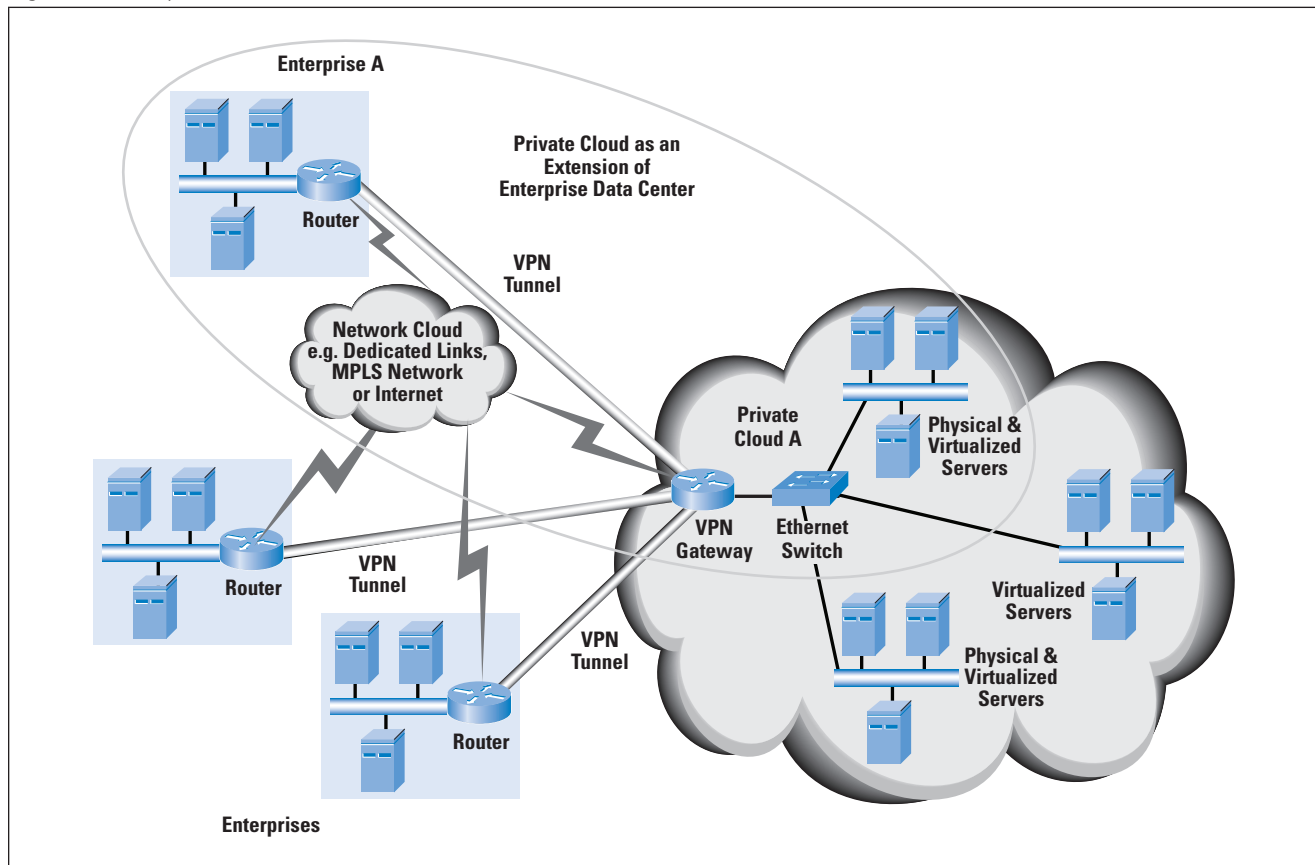
IaaS Private Clouds

Consider an IaaS cloud to which an enterprise connects to augment its server capacity for a limited period of time. Assume that the enterprise uses a $10.x.x.x$ private addressing scheme for all its servers because they are internal to the enterprise. It would be ideal if the additional servers provided by the IaaS cloud were part of the same addressing scheme (the $10.x.x.x$ scheme). As shown in Figure 3, the IaaS cloud service provider has partitioned a portion of its public cloud to realize a private cloud for enterprise A. The private cloud is reachable as a LAN extension to the servers in enterprise A's data center.

How is this reachability realized? A secure *Virtual Private Network* (VPN) tunnel is first established between the enterprise data center and the public cloud. This tunnel uses public IP addresses to establish the site-to-site VPN connection. The VPN gateway on the cloud service provider side uses multiple contexts—each context corresponding to a specific private cloud. Traffic from enterprise A is decrypted and forwarded over to an Ethernet switch to the private cloud for enterprise A. A server on enterprise A's internal data center sees a server on private cloud A to be on the same network.

In practice, data-center servers might be segmented into their own VLANs or IP networks according to policy and applications. The configuration and forwarding policies on the private cloud end would reflect this segmentation as well.

Figure 3: Example of Private Clouds



The following are some possible evolution scenarios for this scheme:

- *Automation of the VPN connection between the enterprise and cloud service provider:* This automation can be done through a management system responsible for the cloud bursting and server augmentation. The system sets up the VPN tunnels and configures the servers on the cloud service provider end. The management system is set up and operated by the cloud service provider.
- *Integration of the VPN functions with the site-to-site VPN network functions from service providers:* For example, service providers offer MPLS Layer 3 VPNs and Layer 2 VPNs (also known as *Virtual Private LAN Service*, or VPLS) as part of their offerings. Enterprise and cloud service providers could be set up to use these network services.
- *Cloud service providers using multiple data centers:* In such a situation, a VPLS-like service can be used to bridge the individual data centers, providing complete transparency from the enterprise side about the location of the cloud servers.

CloudNet is an example of a framework being developed by AT&T Labs and the University of Massachusetts at Amherst to address the latter two scenarios.

Layer 2 versus Layer 3 Connectivity for Cloud Networks

Enterprises and vendors follow some guidelines regarding where to use Layer 2 (switching) and Layer 3 (routing) in the network. Layer 2 is the simpler mode, where the Ethernet MAC address and *Virtual LAN* (VLAN) information are used for forwarding. The disadvantage of Layer 2 networks is scalability. When we use Layer 2 addressing and connectivity in the manner specified previously for IaaS clouds, we end up with a flat topology, which is not ideal when there are a large number of nodes. The option is to use routing and subnets—to provide segmentation for the appropriate functions at the cost of forwarding performance and network complexity.

VM migration introduces its own set of problems. The most common scenario is when a VM is migrated to a different host on the same Layer 2 topology (with the appropriate VLAN configuration). Consider the case where a VM with open *Transmission Control Protocol* (TCP) connections is migrated. If live migration is used, TCP connections will not see any downtime except for a short “hiccup.” However, after the migration, IP and TCP packets destined for the VM will need to be resolved to a different MAC address or the same MAC address but now connected to a different physical switch in the network so that the connections can be continued without disruption. Proposed solutions include an unsolicited *Address Resolution Protocol* (ARP) request from the migrated VM so that the switch tables can be updated, a pseudo-MAC address for the VM that is externally managed (defined in research work being done at the University of California at San Diego), and so on.

With VPLS and similar Layer 2 approaches, VM migration can proceed as before—across the same Layer 2 network. Alternatively, it may be less complex to “freeze” the VM and move it across either a Layer 2 or Layer 3 network with the TCP connections having to be torn down by the counterpart(s) communicating with the VM. This scenario is not a desired one from an application availability consideration, but it can lower complexity.

Cloud Federation

Thus far we have considered the situation of data centers that are owned or run by the same cloud services provider. Connectivity between the data centers to provide the vision of “one cloud” is completely within the control of the cloud service provider.

There may be situations where an organization or enterprise needs to be able to work with multiple cloud providers because of migration from one cloud service to another, merger of companies working with different cloud providers, cloud providers who provide best-of-class services, and so on. Cloud interoperability and the ability to share various types of information between clouds become important in such scenarios. Although cloud service providers might see less urgency for any interoperability, enterprise customers will see a need to push them in that direction.

This broad area of cloud interoperability is sometimes known as *cloud federation*. One definition of cloud federation as proposed by Reuven Cohen of Enomaly follows:

“Cloud federation manages consistency and access controls when two or more independent geographically distributed clouds share either authentication, files, computing resources, command and control, or access to storage resources.”

The following are some of the considerations in cloud federation:

- An enterprise user wishing to access multiple cloud services would be better served if there were just a single sign-on scheme. This scheme may be implemented through an authentication server maintained by an enterprise that provides the appropriate credentials to the cloud service providers. Alternatively, a central trusted authentication server to which all the cloud services interface could be used.
- Computing and storage resources may be orchestrated through the individual enterprise or through an interoperability scheme established between the cloud providers (through a federation agreement, for example). Files may need to be transferred, services invoked, and computing resources added or removed in a useful and transparent manner. A related area is VM migration and how it can be done transparently and reliably. The *Desktop Management Task Force* (DMTF) has released a specification called the *Open Virtualization Format* (OVF) for describing a VM. It can be reasonably assumed that the payload for VM migration will be in the OVF format so that it can be interpreted across multiple vendor offerings. In effect, cloud federation has to provide transparent workload orchestration between the clouds on behalf of the enterprise user.
- Connectivity between clouds includes Layer 2 versus Layer 3 considerations and secure tunnel technologies that need to be agreed upon. Consistency and a common understanding are required irrespective of the model or technologies.
- An often-ignored concern for cloud confederation is charging or billing and reconciliation. Management and billing systems need to work together for cloud federation to be a viable option. This reality is underlined by the fact that clouds rely on per-use billing. Cloud service providers might need to look closely at telecom service provider business models for peering arrangements as a possible starting point.

Cloud federation is a relatively new area in cloud computing. It is likely that standards bodies will first need to agree upon a set of requirements before the service interfaces can be defined and subsequently realized. Provider and vendor innovation will also significantly affect this area—in fact, cloud service operators are likely to establish peering relationships and start addressing this area even before the standards bodies.

Security

As indicated in Part 1, the biggest deterrent for IT managers from venturing into cloud computing is the problem of security and loss of control. Before considering a move to a cloud service provider, enterprises need to consider some of the following security topics:

- The cloud service provider's security processes will need to be as good as or better than the processes that the enterprise uses. An audit of the vendor's processes will need to be done periodically, possibly including patches and security updates for the individual components that are used. For example, in an IaaS scenario with some preconfigured images of operating systems and applications, the cloud service provider should have the latest patches applied on the individual components.
- Infrastructure and data isolation must be assured between multiple tenants of the cloud service provider. This requirement is complicated because it is closely intertwined with the business model used by the cloud provider. For example, an IaaS provider might provide multiple tenants with VMs running on the same physical machine. Depending upon the type of work that is to be executed on the cloud, this setup may or may not be acceptable to a cloud user. In such cases, the cloud service provider should have the ability to provide separate physical servers for specific customers (and bill appropriately).
- In cases where a hypervisor and VMs are used, the hypervisor should be treated as an operating system and have the latest security patches applied to it. Security patches and updates are also essential for paravirtualized operating systems used in the VMs.
- Security functions can run as virtual appliances over hypervisors in a cloud environment. Thus it is possible for cloud users in an IaaS environment to load and configure their own firewall or other security virtual appliance to run within the cloud. The software images used for these virtual appliances need to be managed and patched similar to the way the OS, hypervisor, and other applications are managed and patched.
- Logging and audit trails for applications are important for enterprises to understand both application performance and security gaps. Cloud services providers should enable access to their application monitoring and profiling tools, where applicable.
- Authentication mechanisms ("You are who you say you are") are required at both ends of the connection—at the cloud user and cloud service provider levels. The cloud user and operator must agree upon schemes such as authentication with digital certificates and certificate authorities.

- Configuration and updates to the network infrastructure must be audited and tracked. For example, incorrect VLAN configuration on the switches can result in undesired traffic patterns between physical machines and computing resources. It would be useful to log and audit the configuration records for proper security and uptime.
- Because the cloud service is exposed to the outside world, the cloud infrastructure should support security functions such as intrusion detection and prevention, firewalling to prevent disallowed traffic, and *Denial of Service* (DoS) prevention. The cloud service is vulnerable to *Distributed Denial of Service* (DDoS) attacks—which can effectively choke its access lines, resulting in cloud users being locked out of the cloud service. Network-based DDoS prevention is a possible solution—with one of the techniques involving distribution of the cloud infrastructure to specific geographic areas and the ability to redirect cloud users in case of DDoS lockouts.

Virtualization and Security

Two options are under discussion for security in the context of virtualization. Both are useful in building out security-enabled cloud infrastructures. One option involves plug-ins to the hypervisor so that packets destined to the VMs are captured and processed by the security plug-ins. This setup enables application of security functions to the packet before it gets to the VMs. A second option is to make a specific VM handle the security functions without changing or adding to the hypervisor. The hypervisor plug-in option has the advantage of performance and initial isolation, whereas the separate VM option has the advantage of keeping the hypervisor simple and extrapolating the model that exists in physical server infrastructure. Note that these options are not mutually exclusive.

VM migration is another area where security is an important consideration. The hypervisor is responsible for the two-way communication, with the hypervisor on the destination physical machine to accomplish the migration. It is important that the connection between the source and destination hypervisors is authenticated and encrypted during the course of this migration. In addition, VM migration introduces the possibility of a DoS attack because a rogue hypervisor could overwhelm a destination machine by migrating a large number of VMs to the destination machine. Policies and logic are required at the hypervisor level to ensure that these vulnerabilities are addressed. In addition, network-based throttling might be required so that live migration does not cause congestion, which might happen if a large number of VMs need to be migrated to a destination machine at the same time.

Standards Bodies Involved in Cloud Computing

Numerous standards bodies are involved in cloud computing, addressing aspects of interoperability, virtualization migration formats, and security. Some of the organizations involved have established liaisons with the other *Standards Development Organizations* (SDOs) so that there is no duplication of effort.

The *Desktop Management Task Force* (DMTF) has specified a portable format for packaging the software to run as a VM. Known as the *Open Virtualization Format* (OVF), this package format is seeing increased use. The VM can be written onto a disk or external storage and can be moved from one physical machine to another. The DMTF has also formed a group called the *Open Cloud Standards Incubator*, which focuses on standardizing the interactions between cloud environments, including the development of resource management, packaging formats, and security.

The *Cloud Security Alliance* (CSA) is a new group formed to address security aspects of cloud computing with a focus on security assessment and management. The initial part of the effort is on developing an *Audit, Assertion, Assessment and Assurance* (API) set (A6).

The *Organization for the Advancement of Structured Information Standards* (OASIS) sees cloud computing as an extension of the *Service-Oriented Architecture* (SOA) used today in IT environments. The areas for standardization include security and policy, content format control, registry and directory standards, as well other SOA methods.

The *Storage Networking Industries Association* (SNIA) has a *Cloud Storage Technical Working Group* (TWG) that works on storage-related problems related to implementation in a cloud. The TWG has developed an interface known as the *Cloud Data Management Interface* (CDMI), which clients will use for control and configuration of the cloud.

Some Perspectives on Cloud Computing

In this section we outline and provide some perspective on cloud-computing topics that have seen interest (and some heated discussion). This list is not intended to be comprehensive but to provide a quick snapshot. Though this section has a degree of subjectivity, it is directed only to providing a broader perspective.

- *Cloud computing and SOA*: Some view cloud computing as a specific deployment case of an SOA—and this view is more popular than the one that says that cloud computing is the evolution of SOA. David Linthicum outlines that these views are complementary in that cloud-computing services will most likely be defined through SOA. IaaS provides a new variant because you can now access raw compute and storage resources as a service. Independent of the argument that “We have seen this before,” there is value to defining and invoking available services in the cloud.

- *Server virtualization schemes:* Comparisons are sometimes made based on how vendor products approach virtualization—type 1 versus type 2—and full versus paravirtualization. These approaches have pros and cons. The final decision often hinges on total costs, so it might be useful to move forward from this debate. Incidentally, vendors provide several useful tools for VM backup, recovery, fault tolerance, load management, and so on, and these tools work equally well for the various approaches to virtualization. It may be argued that these tools and features such as VM migration and the associated costs are more useful areas for comparison.
- *Other types of virtualization:* This article has deliberately omitted discussion of other types of virtualization, including desktop, application, and presentation virtualization. Some of these schemes (server-hosted desktop virtualization is one example) are affected by the cloud, specifically in the areas of network connectivity, authentication, and quality of experience. In general, any thin-client experience is affected by the cloud or data center because most of the work is done at the servers. From a cloud perspective, these types of virtualization schemes are considered to be applications that need to run reliably and consistently.
- *Data transfer and network bandwidth:* IaaS has provided a flexible model, in which you are charged based on compute power usage, storage consumed, and the duration of usage. However, there is another important factor—data needs to be sent back and forth between the cloud user and cloud service provider. Several IaaS providers charge for the amount of data transferred over the link. These charges can quickly add up if your applications are very chatty and require a lot of back-and-forth data traffic. Another concern here is the amount of time the initial upload or download can consume—for example, when you want to move a large number of your files to the IaaS provider's storage, you can tie up the link for hours. In fact, one provider has a model where cloud users can send storage media through a postal or package service for upload to the cloud provider's storage arrays.
- *WAN acceleration for the cloud:* Continuing on the previous point, chatty protocols and applications can benefit from WAN acceleration devices that can be used on both ends of a WAN link to cache and locally serve enterprise applications. These devices are not specific to the cloud—they have been used for several years for application performance improvement when a WAN link is involved. Recently, virtual network appliances for WAN acceleration are seeing deployment—here the WAN acceleration is performed by an individual VM instead of a dedicated appliance.
- *VM migration:* This article outlined some of the concerns with VM migration with respect to Layer 2 and Layer 3 topologies. Another consideration is the amount of data that needs to be moved when a VM is migrated across a network. It can potentially be in the range of gigabytes, depending upon the VM and the included operating environment.

Live migration implements this transfer in an incremental fashion so that the demand on the network is spread out. However, snapshot migration (where a VM is suspended or frozen and migrated over the network in full) can cause a surge of data on the network, leading to application performance problems for other VMs and physical machines. Throttling the amount of data that can be sent in a specific period of time, bandwidth reservation and policing at the intermediate network devices is highly desirable in such situations.

- *Management:* The current management paradigms for the cloud components are quite discrete and provide a strong level of control. For example, it is possible to log in to the *Command-Line Interface* (CLI) of a specific switch in the data center for configuration and control of the switch parameters. Similarly, it is possible to use the management console provided by the virtualization vendor to configure individual parameters for the hypervisors and VMs (for example, when to initiate VM migration to a different physical machine). Efforts are being made to unify management schemes not just through partnerships between the individual vendors but also with machine-readable interfaces (*Extensible Markup Language* [XML] being a baseline) across the multiple types of equipment and software in the cloud. Enterprise users are unlikely to accept point solutions or tools that require extensive user interaction in the long term.
- *Energy considerations:* One of the benefits of virtualization is the use of a lower number of physical servers to realize a specific function. It follows that overall energy consumption would be reduced because you have fewer servers. Although this fact may indeed be true, it would be good to characterize and monitor the effective energy savings for a specific application (“Your mileage may vary”). For example, the load on each server and the associated I/O and storage traffic may lead to higher power requirements on an individual server basis. Other considerations include the hardware infrastructure of the cloud data center because the power and cooling assumptions per rack are based on average server load.
- *Legal and regulatory considerations:* James Urquhart has compiled a set of criteria for workload migration across multiple locations, one of which is “Follow the law.” Consider the case of a cloud services provider or operator that has data centers in two separate countries. The operator might use the data centers for workload migration as well as load balancing. A problem might arise if the laws in one of the countries impose limitations on what can and cannot be done at the data center. Scenarios include access to all data stored at this data center by authorities or the ability to examine all transactions on the wire at the data center. Workload migration policy statements have to be provided to cloud users so that they understand what they are signing up to. Alternatively, they might be provided the ability to set preferences for workload migration. This area is potentially worrisome, so it is important that cloud users are aware of their specific situation.

Conclusion

This article has served as a vendor-neutral primer to the area of cloud computing. In Part 1, we provided an introduction to the still-evolving area of cloud computing, including the technologies and some deployment concerns. In Part 2, we provided a more detailed look at the networking factors in the cloud, security aspects, and cloud federation. We also highlighted some areas that are seeing increased attention with cloud-computing proponents and vendors.

The area of cloud computing is very dynamic and offers scope for innovative technologies and business models. Ongoing work with respect to solutions is substantial, in the vendor research labs and product development organizations as well as in academia. It is clear that cloud computing will see significant advances and innovation in the next few years.

For Further Reading (see Part 1 for additional references)

- [0] T. Sridhar, “Cloud Computing: A Primer, Part 1: Models and Technologies,” *The Internet Protocol Journal*, Volume 12, No. 3, September 2009.
- [1] “Building Data-Centric n-Tier Enterprise Systems,” PowerVision white paper, http://www.powervision.com/html/news/n_tier_arch.pdf
- [2] “Networking in the (Storm) Clouds,” Michael Morris, <http://www.networkworld.com/community/node/43872>
- [3] “Is the Relational Database Doomed?” Tony Bain, <http://www.readwriteweb.com/enterprise/2009/02/is-the-relational-database-doomed.php>
- [4] “Cisco VN-Link: Virtualization-Aware Networking,” http://www.cisco.com/en/US/solutions/collateral/ns340/ns517/ns224/ns892/ns894/white_paper_c11-525307_ps9902_Products_White_Paper.html
- [5] IEEE work (in progress) on Virtual Ethernet Bridging—VEPA and VN-Tag approaches—search for “VEPA” and “VN-Tag” in the directory at: <http://www.ieee802.org/1/files/public/docs2009>
- [6] “PortLand: A Scalable Fault-Tolerant Layer 2 Data Center Network Fabric,” Mysore et al., <http://ccr.sigcomm.org/online/?q=node/503>
- [7] “VL2: A Scalable and Flexible Data Center Network,” Greenberg et al., <http://ccr.sigcomm.org/online/?q=node/502>
- [8] “The Case for Enterprise-Ready Virtual Private Clouds,” Wood et al., http://www.usenix.org/event/hotcloud09/tech/full_papers/wood.pdf

- [9] "Solving the Problem of Cloud Interoperability," Reuven Cohen,
<http://reuvencohen.sys-con.com/node/798504>
- [10] "Security Guidance for Critical Areas of Focus in Cloud Computing," Cloud Security Alliance,
<http://www.cloudsecurityalliance.org/guidance/csaguide.pdf>
- [11] Rational Survivability Blog, Chris Hoff's blog on various topics, including cloud security,
<http://www.rationalsurvivability.com/blog/>
- [12] "Hey, You, Get Off of My Cloud: Exploring Information Leakage in Third-Party Compute Clouds" by Ristenpart, et al,
<http://people.csail.mit.edu/tromer/papers/cloud-sec.pdf>
- [13] "Empirical Exploitation of Live Virtual Machine Migration," Oberheide et al.,
<http://www.eecs.umich.edu/techreports/cse/2007/CSE-TR-539-07.pdf>
- [14] List of and links to cloud standards organizations,
http://cloud-standards.org/wiki/index.php?title=Main_Page/
- [15] "Open Virtualization Format Specification," DMTF,
http://www.dmtf.org/standards/published_documents/DSP0243_1.0.0.pdf
- [16] "SOA cloud computing relationship leaves some folks in a fog," David Linthicum,
<http://www.gcn.com/Articles/2009/03/09/Guest-commentary-SOA-cloud.aspx>
- [17] "Is your data center ready for virtualization," Eaton white paper,
http://i.zdnet.com/whitepapers/Eaton_Is_your_data_center_ready_for_virtualization.pdf
- [18] "The great paradigm shift of cloud computing is not self-service," James Urquhart, http://news.cnet.com/8301-19413_3-10127654-240.html?tag=mncol;txt

T. SRIDHAR received his BE in Electronics and Communications Engineering from the College of Engineering, Guindy, Anna University, Madras, India, and his Master of Science in Electrical and Computer Engineering from the University of Texas at Austin. He can be reached at TSridhar@leitnet.com

Protocol Basics: Secure Shell Protocol

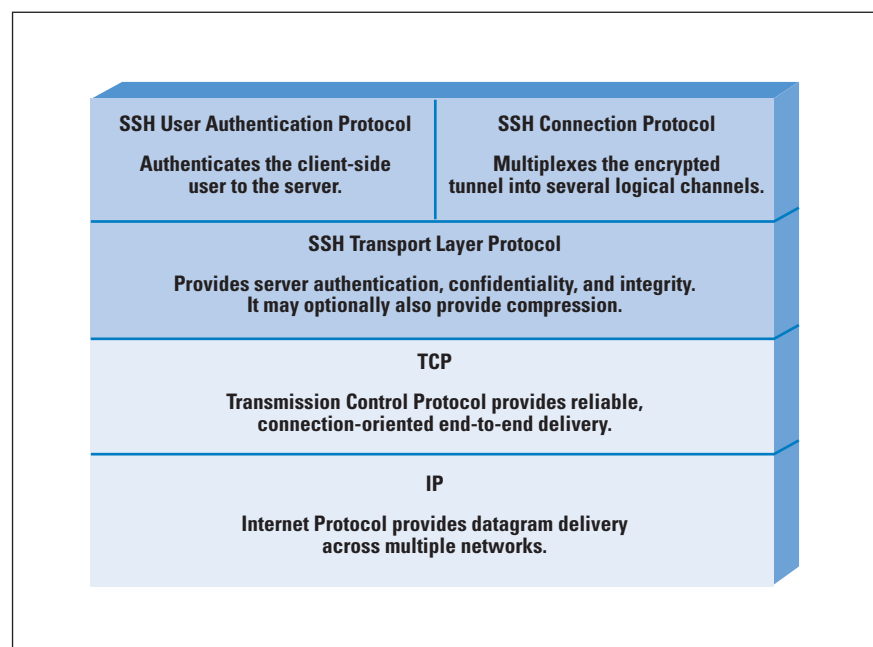
by William Stallings

Secure Shell (SSH) Protocol is a protocol for secure network communications designed to be relatively simple and inexpensive to implement. The initial version, SSH1, focused on providing a secure remote logon facility to replace Telnet and other remote logon schemes that provided no security^[4]. SSH also provides a more general client-server capability and can be used to secure such network functions as file transfer and e-mail. A new version, SSH2, provides a standardized definition of SSH and improves on SSH1 in numerous ways. SSH2 is documented as a proposed standard in RFCs 4250 through 4256^{[1-3], [5-8]}.

SSH client and server applications are widely available for most operating systems. It has become the method of choice for remote login and X tunneling and is rapidly becoming one of the most pervasive applications for encryption technology outside of embedded systems. SSH is organized as three protocols that typically run on top of TCP (Figure 1):

- *Transport Layer Protocol*: Provides server authentication, data confidentiality, and data integrity with forward secrecy (that is, if a key is compromised during one session, the knowledge does not affect the security of earlier sessions); the transport layer may optionally provide compression
- *User Authentication Protocol*: Authenticates the user to the server
- *Connection Protocol*: Multiplexes multiple logical communications channels over a single underlying SSH connection

Figure 1: SSH Protocol Stack



Transport Layer Protocol

Server authentication occurs at the transport layer, based on the server possessing a public-private key pair. A server may have multiple host keys using multiple different asymmetric encryption algorithms. Multiple hosts may share the same host key. In any case, the server host key is used during key exchange to authenticate the identity of the host. For this authentication to be possible, the client must have presumptive knowledge of the server public host key. RFC 4251 dictates two alternative trust models that can be used:

1. The client has a local database that associates each host name (as typed by the user) with the corresponding public host key. This method requires no centrally administered infrastructure and no third-party coordination. The downside is that the database of name-to-key associations may become burdensome to maintain.
2. The host name-to-key association is certified by a trusted *Certification Authority* (CA). The client knows only the CA root key and can verify the validity of all host keys certified by accepted CAs. This alternative eases the maintenance problem, because ideally only a single CA key needs to be securely stored on the client. On the other hand, each host key must be appropriately certified by a central authority before authorization is possible.

Figure 2: SSH Transport Layer Protocol Packet Exchanges

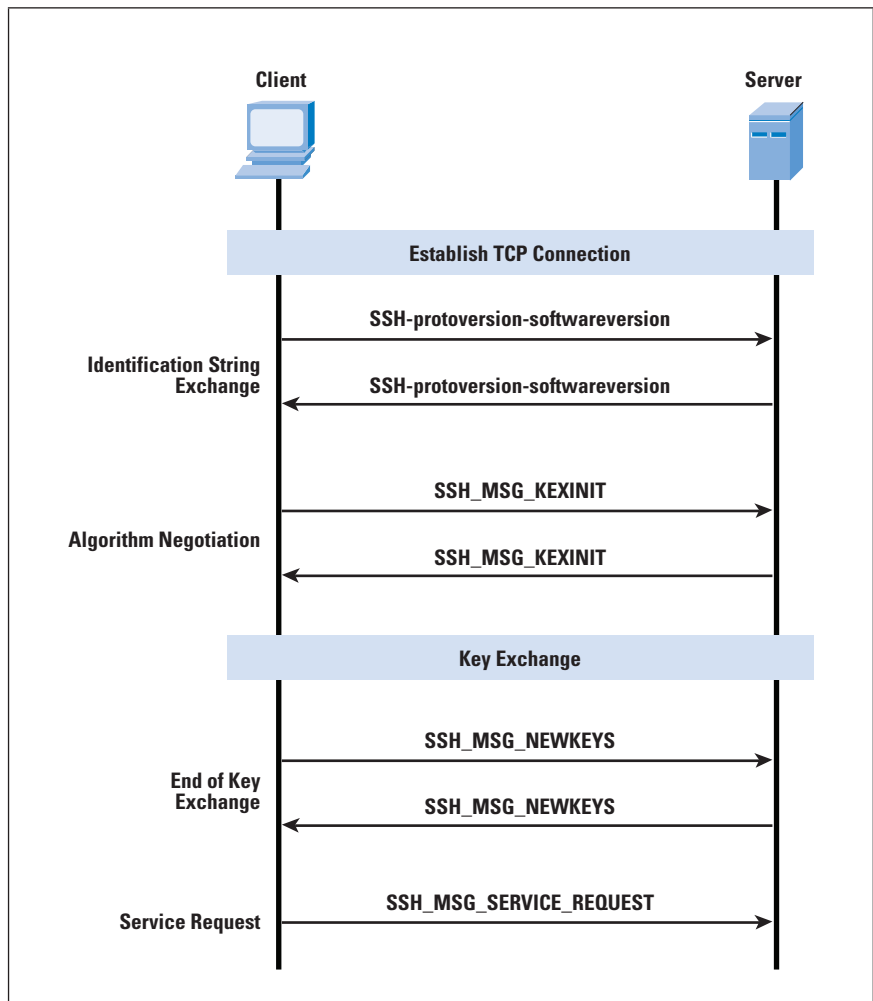
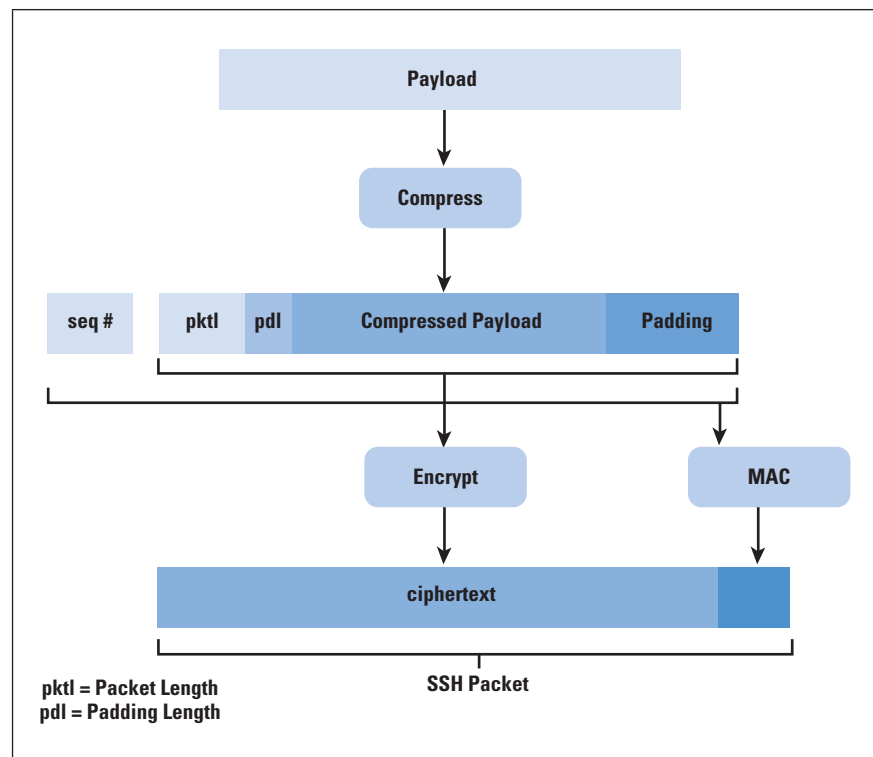


Figure 2 illustrates the sequence of events in the SSH Transport Layer Protocol. First, the client establishes a TCP connection to the server with the TCP protocol and is not part of the Transport Layer Protocol. When the connection is established, the client and server exchange data, referred to as packets, in the data field of a TCP segment. Each packet is in the following format (Figure 3):

- *Packet length*: Packet length is the length of the packet in bytes, not including the packet length and Message Authentication Code (MAC) fields.
- *Padding length*: Padding length is the length of the random padding field.
- *Payload*: Payload constitutes the useful contents of the packet. Prior to algorithm negotiation, this field is uncompressed. If compression is negotiated, then in subsequent packets this field is compressed.
- *Random padding*: After an encryption algorithm is negotiated, this field is added. It contains random bytes of padding so that that total length of the packet (excluding the MAC field) is a multiple of the cipher block size, or 8 bytes for a stream cipher.
- *Message Authentication Code (MAC)*: If message authentication has been negotiated, this field contains the MAC value. The MAC value is computed over the entire packet plus a sequence number, excluding the MAC field. The sequence number is an implicit 32-bit packet sequence that is initialized to zero for the first packet and incremented for every packet. The sequence number is not included in the packet sent over the TCP connection.

Figure 3: SSH Transport Layer Protocol Packet Formation



After an encryption algorithm is negotiated, the entire packet (excluding the MAC field) is encrypted after the MAC value is calculated.

The SSH Transport Layer packet exchange consists of a sequence of steps (Figure 2). The first step, the *identification string exchange*, begins with the client sending a packet with an identification string of the form:

SSH-protoversion-softwareversion SP comments CR LF

where SP, CR, and LF are space character, carriage return, and line feed, respectively. An example of a valid string is **SSH-2.0-b1llsSSH_3.6.3q3<CR><LF>**. The server responds with its own identification string. These strings are used in the Diffie–Hellman key exchange.

Next comes *algorithm negotiation*. Each side sends an **SSH_MSG_KEXINIT** containing lists of supported algorithms in the order of preference to the sender. Each type of cryptographic algorithm has one list. The algorithms include key exchange, encryption, MAC algorithm, and compression algorithm. Table 1 shows the allowable options for encryption, MAC, and compression. For each category, the algorithm chosen is the first algorithm on the client’s list that is also supported by the server.

Table 1: SSH Transport Layer Cryptographic Algorithms

Cipher	
3des-cbc*	Three-key Triple Digital Encryption Standard (3DES) in Cipher-Block-Chaining (CBC) mode
blowfish-cbc	Blowfish in CBC mode
twofish256-cbc	Twofish in CBC mode with a 256-bit key
twofish192-cbc	Twofish with a 192-bit key
twofish128-cbc	Twofish with a 128-bit key
aes256-cbc	Advanced Encryption Standard (AES) in CBC mode with a 256-bit key
aes192-cbc	AES with a 192-bit key
aes128-cbc**	AES with a 128-bit key
Serpent256-cbc	Serpent in CBC mode with a 256-bit key
Serpent192-cbc	Serpent with a 192-bit key
Serpent128-cbc	Serpent with a 128-bit key
arcfour	RC4 with a 128-bit key
cast128-cbc	CAST-128 in CBC mode

MAC Algorithm	
hmac-sha1*	HMAC-SHA1; Digest length = Key length = 20
hmac-sha1-96**	First 96 bits of HMAC-SHA1; Digest length = 12; Key length = 20
hmac-md5	HMAC-SHA1; Digest length = Key length = 16
hmac-md5-96	First 96 bits of HMAC-SHA1; Digest length = 12; Key length = 16

Compression Algorithm	
none*	No compression
zlib	Defined in RFCs 1950 and 1951

* = Required

** = Recommended

The next step is *key exchange*. The specification allows for alternative methods of key exchange, but at present only two versions of Diffie–Hellman key exchange are specified. Both versions are defined in RFC 2409 and require only one packet in each direction. The following steps are involved in the exchange. In this, C is the client; S is the server; p is a large safe prime; g is a generator for a subgroup of $GF(p)$; q is the order of the subgroup; V_S is the S identification string; V_C is the C identification string; K_S is the S public host key; I_C is the C `SSH_MSG_KEXINIT` message; and I_S is the S `SSH_MSG_KEXINIT` message that was exchanged before this part began. The values of p , g , and q are known to both client and server as a result of the algorithm selection negotiation. The hash function `hash()` is also decided during algorithm negotiation.

1. C generates a random number x ($1 < x < q$) and computes $e = g^x \bmod p$. C sends e to S.
2. S generates a random number y ($0 < y < q$) and computes $f = g^y \bmod p$. S receives e . It computes $K = e^y \bmod p$, $H = \text{hash}(V_C \parallel V_S \parallel I_C \parallel I_S \parallel K_S \parallel e \parallel f \parallel K)$, and signature s on H with its private host key. S sends $(K_S \parallel f \parallel s)$ to C. The signing operation may involve a second hashing operation.
3. C verifies that K_S really is the host key for S (for example, using certificates or a local database). C is also allowed to accept the key without verification; however, doing so will render the protocol insecure against active attacks (but may be desirable for practical reasons in the short term in many environments). C then computes $K = f^x \bmod p$, $H = \text{hash}(V_C \parallel V_S \parallel I_C \parallel I_S \parallel K_S \parallel e \parallel f \parallel K)$, and verifies the signature s on H .

As a result of these steps, the two sides now share a master key K . In addition, the server has been authenticated to the client, because the server has used its private key to sign its half of the Diffie–Hellman exchange. Finally, the hash value H serves as a session identifier for this connection. When computed, the session identifier is not changed, even if the key exchange is performed again for this connection to obtain fresh keys.

The *end of key exchange* is signaled by the exchange of `SSH_MSG_NEWKEYS` packets. At this point, both sides may start using the keys generated from K , as discussed subsequently.

The final step is *service request*. The client sends an `SSH_MSG_SERVICE_REQUEST` packet to request either the User Authentication or the Connection Protocol. Subsequent to this request, all data is exchanged as the payload of an SSH Transport Layer packet, protected by encryption and MAC.

The keys used for encryption and MAC (and any needed IVs) are generated from the shared secret key K , the hash value from the key exchange H , and the session identifier, which is equal to H unless there has been a subsequent key exchange after the initial key exchange. The values are computed as follows:

- Initial IV client to server: $\text{HASH}(K \parallel H \parallel \text{"A"} \parallel \text{session_id})$
- Initial IV server to client: $\text{HASH}(K \parallel H \parallel \text{"B"} \parallel \text{session_id})$
- Encryption key client to server: $\text{HASH}(K \parallel H \parallel \text{"C"} \parallel \text{session_id})$
- Encryption key server to client: $\text{HASH}(K \parallel H \parallel \text{"D"} \parallel \text{session_id})$
- Integrity key client to server: $\text{HASH}(K \parallel H \parallel \text{"E"} \parallel \text{session_id})$
- Integrity key server to client: $\text{HASH}(K \parallel H \parallel \text{"F"} \parallel \text{session_id})$

where $\text{HASH}()$ is the hash function determined during algorithm negotiation.

User Authentication Protocol

The *User Authentication Protocol* provides the means by which the client is authenticated to the server.

Three types of messages are always used in the User Authentication Protocol. Authentication requests from the client have the format:

```
byte    SSH_MSG_USERAUTH_REQUEST (50)
string  username
string  service name
string  method name
....    method-specific fields
```

where *username* is the authorization identity the client is claiming, *service name* is the facility to which the client is requesting access (typically the SSH Connection Protocol), and *method name* is the authentication method being used in this request. The first byte has decimal value 50, which is interpreted as **SSH_MSG_USERAUTH_REQUEST**.

If the server either rejects the authentication request or accepts the request but requires one or more additional authentication methods, the server sends a message with the format:

```
byte          SSH_MSG_USERAUTH_FAILURE (51)
name-list     authentications that can continue
boolean       partial success
```

where the *name-list* is a list of methods that may productively continue the dialog. If the server accepts authentication, it sends a single-byte message, **SSH_MSG_USERAUTH_SUCCESS (52)**.

The message exchange involves the following steps:

1. The client sends a **SSH_MSG_USERAUTH_REQUEST** with a requested method of none.
2. The server checks to determine if the username is valid. If not, the server returns **SSH_MSG_USERAUTH_FAILURE** with the partial success value of false. If the username is valid, the server proceeds to step 3.
3. The server returns **SSH_MSG_USERAUTH_FAILURE** with a list of one or more authentication methods to be used.
4. The client selects one of the acceptable authentication methods and sends a **SSH_MSG_USERAUTH_REQUEST** with that method name and the required method-specific fields. At this point, there may be a sequence of exchanges to perform the method.
5. If the authentication succeeds and more authentication methods are required, the server proceeds to step 3, using a partial success value of true. If the authentication fails, the server proceeds to step 3, using a partial success value of false.
6. When all required authentication methods succeed, the server sends a **SSH_MSG_USERAUTH_SUCCESS** message, and the Authentication Protocol is over.

The server may require one or more of the following authentication methods:

- *publickey*: The details of this method depend on the public-key algorithm chosen. In essence, the client sends a message to the server that contains the client's public key, with the message signed by the client's private key. When the server receives this message, it checks to see whether the supplied key is acceptable for authentication and, if so, it checks to see whether the signature is correct.
- *password*: The client sends a message containing a plaintext password, which is protected by encryption by the Transport Layer Protocol.
- *hostbased*: Authentication is performed on the client's host rather than the client itself. Thus, a host that supports multiple clients would provide authentication for all its clients. This method works by having the client send a signature created with the private key of the client host. Thus, rather than directly verifying the user's identity, the SSH server verifies the identity of the client host—and then believes the host when it says the user has already authenticated on the client side.

Connection Protocol

The SSH Connection Protocol runs on top of the SSH Transport Layer Protocol and assumes that a secure authentication connection is in use. That secure authentication connection, referred to as a *tunnel*, is used by the Connection Protocol to multiplex a number of logical channels.

RFC 4254, “The Secure Shell (SSH) Connection Protocol,” states that the Connection Protocol runs on top of the Transport Layer Protocol and the User Authentication Protocol. RFC 4251, “SSH Protocol Architecture,” states that the Connection Protocol runs over the User Authentication Protocol. In fact, the Connection Protocol runs over the Transport Layer Protocol, but assumes that the User Authentication Protocol has been previously invoked.

All types of communication using SSH, such as a terminal session, are supported using separate channels. Either side may open a channel. For each channel, each side associates a unique channel number, which need not be the same on both ends. Channels are flow-controlled using a window mechanism. No data may be sent to a channel until a message is received to indicate that window space is available.

The life of a channel progresses through three stages: opening a channel, data transfer, and closing a channel.

When either side wishes to open a new channel, it allocates a local number for the channel and then sends a message of the form:

byte	SSH_MSG_CHANNEL_OPEN
string	channel type
uint32	sender channel
uint32	initial window size
uint32	maximum packet size
....	channel type specific data follows

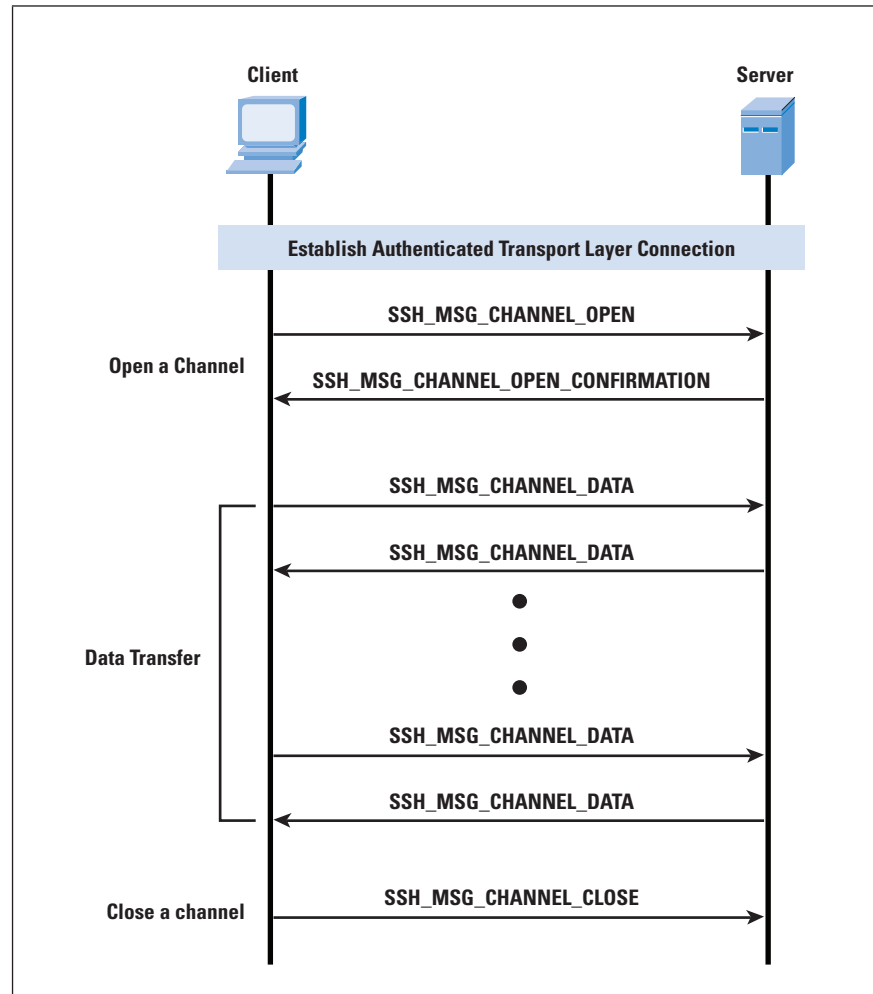
where *uint32* means unsigned 32-bit integer. The *channel type* identifies the application for this channel, as described subsequently. The *sender channel* is the local channel number. The *initial window size* specifies how many bytes of channel data can be sent to the sender of this message without adjusting the window. The *maximum packet size* specifies the maximum size of an individual data packet that can be sent to the sender. For example, one might want to use smaller packets for interactive connections to get better interactive response on slow links.

If the remote side is able to open the channel, it returns a **SSH_MSG_CHANNEL_OPEN_CONFIRMATION** message, which includes the sender channel number, the recipient channel number, and window and packet size values for incoming traffic. Otherwise, the remote side returns a **SSH_MSG_CHANNEL_OPEN_FAILURE** message with a reason code indicating the reason for failure.

After a channel is open, *data transfer* is performed using a **SSH_MSG_CHANNEL_DATA** message, which includes the recipient channel number and a block of data. These messages, in both directions, may continue as long as the channel is open.

When either side wishes to close a channel, it sends a **SSH_MSG_CHANNEL_CLOSE** message, which includes the recipient channel number. Figure 4 provides an example of Connection Protocol Exchange.

Figure 4: Example SSH Connection Protocol Message Exchange



Four channel types are recognized in the SSH Connection Protocol specification:

- *session*: Session refers to the remote execution of a program. The program may be a shell, an application such as file transfer or e-mail, a system command, or some built-in subsystem. When a session channel is opened, subsequent requests are used to start the remote program.
- *x11*: This channel type refers to the X Window System, a computer software system and network protocol that provides a GUI for networked computers. X allows applications to run on a network server but be displayed on a desktop machine.
- *forwarded-tcpip*: This channel type is remote port forwarding, as explained subsequently.
- *direct-tcpip*: This channel type is local port forwarding, as explained subsequently.

One of the most useful features of SSH is *port forwarding*. Port forwarding provides the ability to convert any insecure TCP connection into a secure SSH connection. It is also referred to as *SSH tunneling*. We need to know what a port is in this context. A *port* is an identifier of a user of TCP. So, any application that runs on top of TCP has a port number. Incoming TCP traffic is delivered to the appropriate application on the basis of the port number. An application may employ multiple port numbers. For example, for the *Simple Mail Transfer Protocol* (SMTP), the server side generally listens on port 25, so that an incoming SMTP request uses TCP and addresses the data to destination port 25. TCP recognizes that this address is the SMTP server address and routes the data to the SMTP server application.

Figure 5: SSH Transport Layer Packet Exchanges

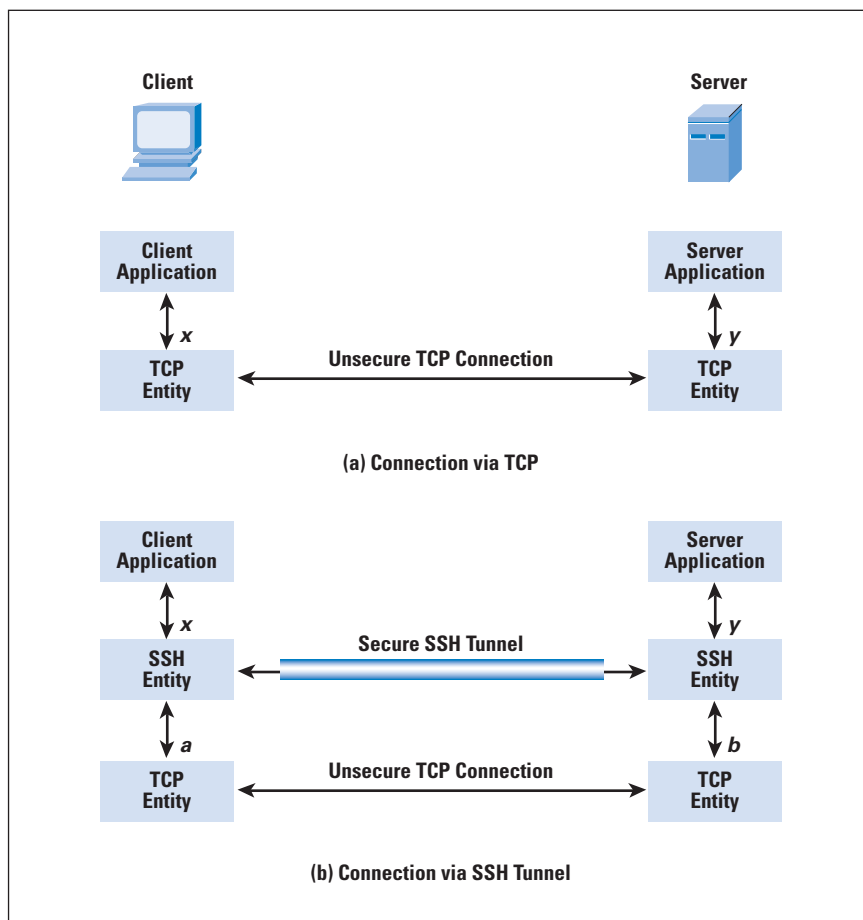


Figure 5 illustrates the basic concept behind port forwarding. We have a client application that is identified by port number x and a server application identified by port number y . At some point, the client application invokes the local TCP entity and requests a connection to the remote server on port y . The local TCP entity negotiates a TCP connection with the remote TCP entity, such that the connection links local port x to remote port y .

To secure this connection, SSH is configured so that the SSH Transport Layer Protocol establishes a TCP connection between the SSH client and server entities with TCP port numbers a and b , respectively. A secure SSH tunnel is established over this TCP connection. Traffic from the client at port x is redirected to the local SSH entity and travels through the tunnel where the remote SSH entity delivers the data to the server application on port y . Traffic in the other direction is similarly redirected.

SSH supports two types of port forwarding: local forwarding and remote forwarding. *Local forwarding* allows the client to set up a “hijacker” process. This process will intercept selected application-level traffic and redirect it from an unsecured TCP connection to a secure SSH tunnel. SSH is configured to listen on selected ports. SSH grabs all traffic using a selected port and sends it through an SSH tunnel. On the other end, the SSH server sends the incoming traffic to the destination port dictated by the client application.

The following example should help clarify local forwarding. Suppose you have an e-mail client on your desktop and use it to get e-mail from your mail server through the *Post Office Protocol* (POP). The assigned port number for POP3 is port 110. We can secure this traffic in the following way:

1. The SSH client sets up a connection to the remote server.
2. Select an unused local port number, say 9999, and configure SSH to accept traffic from this port destined for port 110 on the server.
3. The SSH client informs the SSH server to create a connection to the destination, in this case mailserver port 110.
4. The client takes any bits sent to local port 9999 and sends them to the server inside the encrypted SSH session. The SSH server decrypts the incoming bits and sends the plaintext to port 110.
5. In the other direction, the SSH server takes any bits received on port 110 and sends them inside the SSH session back to the client, which decrypts and sends them to the process connected to port 9999.

With *remote forwarding*, the user’s SSH client acts on the server’s behalf. The client receives traffic with a given destination port number, places the traffic on the correct port, and sends it to the destination the user chooses.

A typical example of remote forwarding follows: You wish to access a server at work from your home computer. Because the work server is behind a firewall, it will not accept an SSH request from your home computer. However, from work you can set up an SSH tunnel using remote forwarding.

This process involves the following steps:

1. From the work computer, set up an SSH connection to your home computer. The firewall will allow this, because it is a protected outgoing connection.
2. Configure the SSH server to listen on a local port, say 22, and to deliver data across the SSH connection addressed to remote port, say 2222.
3. You can now go to your home computer and configure SSH to accept traffic on port 2222.
4. You now have an SSH tunnel that you can use for remote login to the work server.

Summary

SSH is one of the most commonly used cryptographic applications. It provides great flexibility and versatility for a wide variety of tasks, including remote administration, file transfer, web development, and penetration testing.

References

- [1] Cusack, F. and Forssen, M. "Generic Message Exchange Authentication for the Secure Shell Protocol (SSH)," RFC 4256, January 2006.
- [2] Lehtinen, S. and Lonvick, C., "The Secure Shell (SSH) Protocol Assigned Numbers," RFC 4250, January 2006.
- [3] Schlyter, J. and Griffin, W. "Using DNS to Securely Publish Secure Shell (SSH) Key Fingerprints," RFC 4255, January 2006.
- [4] Ylonen, T., "SSH – Secure Login Connections over the Internet," Proceedings, Sixth USENIX UNIX Security Symposium, July 1996.
- [5] Ylonen, T. and Lonvick, C., "The Secure Shell (SSH) Protocol Architecture," RFC 4251, January 2006.
- [6] Ylonen, T. and Lonvick, C., "The Secure Shell (SSH) Authentication Protocol," RFC 4252, January 2006.
- [7] Ylonen, T. and Lonvick, C., "The Secure Shell (SSH) Transport Layer Protocol," RFC 4253, January 2006.
- [8] Ylonen, T. and Lonvick, C., "The Secure Shell (SSH) Connection Protocol," RFC 4254, January 2006.

WILLIAM STALLINGS is a consultant, lecturer, and author of more than a dozen books on data communications and computer networking. His latest book is *Cryptography and Network Security* (Prentice Hall, 2010). He maintains a computer science resource site for computer science students and professionals at WilliamStallings.com/StudentSupport.html and is on the editorial board of Cryptologia. He has a Ph.D. in computer science from M.I.T. He can be reached at ws@shore.net

Book Review

Protocol Politics

Protocol Politics: The Globalization of Internet Governance, by Laura DeNardis, MIT Press, 2009, ISBN 978-0-26204257-4.

In *Protocol Politics*, Dr. Laura DeNardis assembles a variety of stories gleaned from official and unofficial *Internet Engineering Task Force* (IETF) records and firsthand accounts, and supplements them with primer-level descriptions of successive generations of Internet addressing and routing protocols to create a broadly accessible overview of the factors that have shaped the present and evolving state of these most central features of Internet technology.

The author, a former enterprise networking consultant and technology analyst, joined the Yale Law School Information Society Project as a Post-Doctoral Fellow in 2006, and became the Executive Director of the program in late 2008. DeNardis approaches the challenge of organizing these disparate materials by adopting an interpretive framework that highlights the role of power—interpersonal as opposed to electrical—as both the primary input and most important output or consequence of the definition, selection, and implementation of Internet protocols.

The book knits together a wealth of important historical information that has to-date remained largely neglected outside of the technical community. Although DeNardis' choice of framing is perfectly legitimate—and in fact quite common within the academic disciplines that delve into the influence of institutions on industries, economies, and society—in this case it leads her to overreach a bit, and arguably to draw a few prominent conclusions that are not well-supported by the balance of available historical evidence.

Organization of the Book

DeNardis employs this interpretive framework across six densely written chapters, the first four of which directly address the significance of power in a different functional context of relevance to the evolution of Internet addressing and routing. The introductory chapter investigates the significance of scarcity and its effect on protocol resource management and *Internet Governance*. Here she devotes considerable space to detailing the critical importance of IP addresses as the single element among Internet protocols that is both indispensable and nonsubstitutable. DeNardis' insightful overview of the general characteristics of IP addresses is somewhat marred by her mixing together of some basic, intrinsic functional properties of addressing (for example, *identifier* and *locator* functions) with various necessary but extrinsic correlates or consequences of those functional properties (for example, universality, external observability), or with contingent features of current IP address usage conventions (for example, indifference to underlying technologies).

In addition, despite the ostensible focus on scarcity in the chapter, no reference is made to that other, equally essential and quantity-constrained feature of the Internet service landscape—that is, the inherently limited, occasionally overtaxed carrying capacity of Internet routing subsystems, particularly the collectively provisioned inter-domain routing system. Overall, *Protocol Politics* provides almost no exposure to the technical, operational, and economic constraints that define the routing environment, much less to the constraints that those factors impose on number resource distribution arrangements. Chapter One closes with an overview of the priorities that justify and define the sphere of Internet Governance which anticipates many of the concluding observations in the book’s final chapter on “Opening Internet Governance.” Both chapters acknowledge “technical expertise” only as a source of institutional or political legitimacy, without according any special significance to the *content* of such expertise, or why it matters at all. Readers of *Protocol Politics* may thus come away with insufficient appreciation of the fact that before Code can become Law (or anything else), it first must be running code—and *that not every wish is translatable into running code.*^[1]

Piercing the Fog of Protocol War

In the three chapters that follow, DeNardis presents her observations about how power shapes and flows from the definition and selection of Internet protocols. Chapter Two covers the first half of this proposition, focusing on the events that followed the December 1990 IETF meeting where, DeNardis suggests, the twin challenges that would shape the development of Internet addressing intersected with the chief institutional impediment that would ultimately reveal the true political nature of Internet standards development.

The first challenge that she identifies is the foreseeable inadequacy of IPv4 as the exclusive addressing resource pool for a rapidly growing and globalizing Internet. In keeping with the overall theme of the book, the second challenge that DeNardis chooses to highlight is the implicitly political challenge of accommodating greater international participation in the U.S.-centric Internet technical coordination and decision-making bodies. Against this backdrop, DeNardis introduces the other chief protagonist in her story, the *International Organization for Standardization* (ISO), which backed the rival *Open Systems Interconnection* (OSI) family of protocols as an alternative, non-TCP/IP-based foundation for the ongoing, global proliferation of data networking. DeNardis details the convoluted, multidimensional deliberations that followed that 1990 IETF meeting, which eventually culminated in 1994 in the formal recognition of IPv6 as “The Next-Generation Internet Protocol.”

Chapter Three goes on to explore the implications of both IPv4 and IPv6 for important civil liberties—especially privacy—and how such considerations did and did not, *but hypothetically might have*, influenced the choice and form of the most important features of TCP/IP.

Chapter Four rounds out the central thesis of the book by illustrating how various national-level considerations—especially government-directed foreign and domestic economic policies—have resulted in an increasingly diverse global pattern of IPv6 adoption.

DeNardis’ detailed account of the complexities surrounding the *IP Next-Generation* (IPng) debate and its aftermath incorporates a diverse mix of sources, from pointed remarks made on various mailing lists, to conference presentations and official *Internet Architecture Board* (IAB) meeting minutes, and represents a major feat of historical scholarship. That said, her presentation of “relevant historical facts” from the 1990–1994 period is by no means complete, nor is her interpretation of the facts that she does cover or the conclusions that she draws from them immune to criticism. For example, in puzzling over possible hidden forces behind the selection of IPv6, DeNardis states that:

“If anything, there was market pressure to adopt an OSI rather than TCP/IP-based protocol. The ISO alternative had the political backing of most Western European governments (sic) influential technology companies, and users invested in OSI protocols, and was even congruent with OSI directives of the United States. The selection of IPv6...” (p. 61)

Although these facts may be beyond dispute, they do not represent the full picture. To give one illustration, in 1989, almost 2 years before the date that DeNardis marks as the start of the IETF’s lone struggle against the combined forces of Europe, influential carriers and hardware manufacturers, and the U.S. government, an indigenous movement of European network operators emerged and began self-organizing to facilitate the exchange of TCP/IP-based traffic, contact information, and operational tips, and to discuss best practices in areas of networking where individual network-level decisions could have far-reaching effects on internetwork performance.

That organization would go on to become *Réseaux IP Européens Network Coordination Centre* (RIPE NCC), the first independent, transnational registry for Internet Protocol number resources, and the institution that would provide the organizational template for the *Regional Internet Registries* (RIRs) that subsequently sprang up in Asia (APNIC, 1993), North America (ARIN, 1997), Latin America (LACNIC, 2002), and Africa (AFRINIC, 2004). These facts point to a level of active indigenous European support for TCP/IP-based networking that would seem to be at odds with any suggestion of a continent united in support of OSI against a less-attractive standard being pushed by an insular foreign organization.

Thus, regardless of whether DeNardis’ concerns about institutions and power relations are well-founded, her intuitions about the division of contestants in the great protocol power struggle clearly are not.^[2]

Market Contrast

Another question that DeNardis raises, obliquely but repeatedly, relates to the possibility of “free markets” as an alternative mechanism for defining, selecting, and distributing Internet protocols and the virtual resources that they create.

In no less than a dozen separate passages scattered across each of the chapters in the book, DeNardis sharply contrasts a range of IETF and RIR institutional processes to the workings of the “free market.” For example, she observes that the value of IP addresses is unknown because they have never been exchanged in free markets (p. 16); that Internet addresses have never been exchanged in free markets (pp. 23, 190); that the privacy potential of Internet technologies is enhanced by selection pressures from free markets (p. 74); that the IETF refused to countenance an IPng protocol selection made by free markets (p. 51); that the selection of IPv6 happened outside the realm of free markets (p. 69); that widespread adoption of IPv6 is impeded by the absence of a free market for protocols (p. 137); that IETF philosophy holds that it would be inappropriate to exchange protocol resources in free markets (pp. 163, 183–184); that the *Internet Assigned Numbers Authority* (IANA) refused to relinquish IP addresses to free markets (pp. 163, 164); that traditional opposition to the exchange of protocol resources in free markets fortified and centralized the IETF’s institutional control (p. 184); and that exchanging IPv4 in free markets has pragmatic appeal, if only as a temporary stopgap (p. 228), although such exchanges might have unintended consequences (p. 229).

Given this frequency of repetition, it is impossible to avoid forming a strong impression of DeNardis’ underlying opinion about the intrinsic merits of “free markets” as compared to the seemingly market-antithetical goals and practices of the IETF and the other TCP/IP-centric standards-setting and technical coordination bodies. However, even if one stipulates that “free markets” would by definition represent a superior alternative to the enumerated protocol design and distribution mechanisms, DeNardis never provides any clear indication of where a model for such “free markets” might be found—whether in Europe, the United States, or anywhere else, now or anytime in the past.

Even her own description of that fateful moment in networking history when IPv6 was selected clearly suggests that the alternative to the IETF process that ultimately prevailed was itself neither “free” nor especially market-like:

“... congruent with OSI directives of the United States. The selection of IPv6, an expansion of the prevailing IPv4 protocol *over such a politically sanctioned OSI alternative* solidified and extended the position of the Internet’s traditional standards-setting establishment as the entity responsible for the Internet’s architectural direction.” (p. 61, emphasis added).

Arguably, the non-inclusion of a pure “free market” example is not merely a coincidence, but rather reflects a more fundamental problem inherent in the concept itself. Further, if one grants that the market mechanism that is *most free* is the one that fosters the broadest participation in those activities that make markets attractive—including openness to participation, exercise of individual choice, competition, accelerated innovation, and wealth creation—then one might interpret the two-plus orders-of-magnitude growth in the number of independent network services providers operating on both sides of the Atlantic since that time as a solid indicator that markets have not suffered too badly from the 1994 decision to extend the lifetime of TCP/IP through IPv6.

Clearly the looming inflection point in IP addressing will provide many irresistible opportunities to revisit that choice in the days ahead. Meanwhile, the question of whether the embrace of an OSI-friendlier IPng by the IETF would have been sufficient to offset the varied negative externalities that might have accompanied such a choice must forever remain unanswered. Would an IETF endorsement have trumped the as-yet incomplete state of OSI standards, as well as OSI’s tighter associations with non-standards-based operating systems, proprietary hardware platforms, and the connection-oriented networking technologies favored by then Internet-averse incumbent *Public Switched Telephone Network* (PSTN) operators? Would that choice alone have created or been likely to foster a freer market, or to have led to a more enthusiastic, widespread embrace of a different post-IPv4 addressing format—or alternately would it have led to the appearance of books like *Protocol Politics*, albeit written from the opposite perspective, and possibly a decade sooner? Contrary to the popular adage, hindsight is not 20/20, any more than is our vision of where to go from here.^[3]

Beyond the Clash of Idealizations

Writing a book review is an inherently risky undertaking, one that is vulnerable to many of the same human biases and errors that have unquestionably informed both the selection and development of various technical standards, just as they have influenced the embrace, rejection, or modification of various market arrangements throughout history.

Even when people (book reviewers, for example) recognize that real-world decisions and their consequences tend to be irreducibly complex—or perhaps precisely because they recognize that complexity—they nevertheless tend to gravitate toward explanatory frameworks and cognitive models that promise to invest their perceptions and choices with the kind of absolute certitude that is very rarely found outside of the physical world (and only infrequently found there).

The problem, of course, is that many such explanatory frameworks can be found to fit quite nicely with the same set of human experiences, even though some of those models may be mutually orthogonal, and some may be quite mutually and actively antagonistic. In this sense, the juxtaposition of pure, frictionless “free markets” alongside the idea of absolutely pure scientific or technical decision making divorced from all other human considerations, while well-calibrated to inflame passions, represents less a contrast of opposites than a rather less illuminating pairing of two deeply unrealistic ideal types. Distilling a book as rich and informative as *Protocol Politics* down to one possible review-sized essence is much easier to accomplish from just such a privileged vantage point, and no doubt this particular review suffers from the all-too-predictable effects described herein. However, with that caveat firmly established, a few more things about *Protocol Politics* deserve to be mentioned here.

First, *Protocol Politics* is an important book. It is the well-written and informative, and is the first to be written for a general audience that draws on the right historical sources (or at least most of the right ones that remain accessible) to cover this critical period in the development of the Internet’s core addressing and routing protocols. Even those who are least likely to be sympathetic to its findings are likely to find *Protocol Politics* to be a thoughtful and engaging read.

Second, IPJ readers and other technologists should not dismiss the inherently political, power-oriented framework that DeNardis employs in *Protocol Politics*. In general, the most honest and effective response to an assertion of *systemic* political or institutional bias is not to claim an equally absolute, otherworldly detachment from the affairs of man, but rather to remind the critic that in a world where all institutions are regarded as manifestations of somebody’s will to power, specific targeted criticisms based *solely* on that fact lose all coherence. Would-be institutional critics who espouse such views thus have no choice but to make a positive argument as to which arrangement, among all of the equally power-tainted institutional arrangements that are possible, should be regarded as the preferable outcome, for whom, and why. Judged in this light, this reviewer feels that “the IETF way” still stands up pretty well, foibles and all. There is always room for improvement, but just as in matters of code, a concrete proposal for improvement is worth a thousand critiques of the past.

Finally, the careful reader may notice a pattern within this review, one composed of points highlighted here even though they may not be equally central to the story presented in *Protocol Politics* (for example, about the role of technical expertise in Internet governance, the dynamic limitations of routing system carrying capacity, the possibility of free market alternatives to current Internet address distribution arrangements, and so on).

Each of these points merits special attention because taken together they help to illuminate the existence of an identical set of critiques that have reappeared periodically in the course of another, much older (actually, centuries-old) debate that parallels the as-yet unresolved debates outlined by DeNardis in *Protocol Politics*.

In both instances, the question at issue involves the relative merits of nonmarket, technical expert-based systems as a means of managing resources that are uniquely central to economic growth, and for mitigating the systemic risks that can threaten that growth. In that other debate, arguments in favor of pure free market solutions have generally been dismissed as extreme and unrealistic for more than a century, ever since the last real-world implementation of such a system finally succumbed to its own chronic instabilities and was replaced by a nonmarket coordination arrangement. More recently, however, a resurgence of extreme turmoil in that parallel industry has undermined belief in expert management, if not in the underlying “hard realities” that were supposed to constitute the managers’ technical domain of expertise. In turn this turmoil has sparked renewed interest in the long-marginalized pure free market proposals, as well as in alternative remedies involving much tighter industry control by nonmarket authorities.

How the current chapter in either of these parallel stories will play out remains to be written. However, those who are eager to anticipate the kind of language that is likely to play a central role in both outcomes will find that a close reading of *Protocol Politics* provides a wealth of possibilities to consider, and more than a few to keep one up at night.

—Tom Vest, Consultant
tvest@eyeconomics.com

References

- [1] DeNardis makes several references to the idea that *Code is Law*, which was first articulated by Larry Lessig in *Code and Other Laws of Cyberspace* (1999) [Editor’s note: *Code* was reviewed in IPJ Volume 11, No. 3]. Here the phrase is juxtaposed with David Clark’s famous paean to “rough consensus and running code,” which DeNardis describes as an “articulation of the IETF’s core philosophy” (p. 47), and amended with a paraphrasing of an early (c. 1992) observation made by Marshall Rose about a common problem encountered when attempting to implement code to satisfy a non-operationally developed standard. The original staying was, “The problems of the real world are remarkably resilient to administrative fiat.”

- [2] Several formerly obscure insights on the events of this period were recently illuminated by RIPE co-founders Rob Blokzijl and Daniel Karrenberg, during RIPE's 20th Anniversary Commemoration at the RIPE 58 meeting in Amsterdam (May 2009). Some of these are available at:

<http://www.ripe.net/ripe/meetings/ripe-58/content/presentations/Blokzijl-RIPE-20-years.pdf>

and

<http://www.ripe.net/ripe/meetings/ripe-58/content/presentations/the-origins-of-ripe.pdf>

- [3] Those wishing to investigate these questions further may benefit substantially from yet another unique historical resource that has recently been made available online. Thanks to the Charles Babbage Institute and the Institute of Technology at the University of Minnesota, the entire ten-year archive of *ConneXions—The Interoperability Report* (1987–1996) is now available online at: <http://www.cbi.umn.edu/hostedpublications/Connexions/index.html>

In keeping with its mandate to track the interoperability of emerging network technologies, *ConneXions* published more than sixty substantive articles on OSI and GOSIP during the period leading up to and following the IPng debates recounted in *Protocol Politics*.

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. Contact us at ipj@cisco.com for more information.



Lorenzo Colitti (L) and Erik Kline
Photo: Matsuzaki Yoshinobu

Colitti and Kline Receive First Itojun Service Award

The first *Itojun Service Award* was presented at the recent IETF meeting in Hiroshima, Japan to Lorenzo Colitti and Erik Kline of Google for their outstanding contributions to the development and deployment of IPv6.

The award honours the memory of Dr. Jun-ichiro “Itojun” Hagino, who passed away in 2007, aged just 37. Established by the friends of Itojun and administered by the *Internet Society* (ISOC), the award recognises and commemorates the extraordinary dedication exercised by itojun over the course of IPv6 development.

“The sustained efforts of Lorenzo and Erik have tangibly increased the availability of Web-based services that use IPv6, reflecting the Itojun Service Award’s focus on pragmatic contributions in the spirit of serving the global Internet’s continued evolution,” said Jun Murai of the Itojun Service Award committee and Director of the WIDE Project. “The award aims to recognize how important both the development of IPv6 and related protocols and efforts to advance their deployment are to ensuring the Internet continues to serve as a platform for innovation around the world.”

The award, expected to be presented annually, includes a presentation crystal, a US\$3,000 honorarium and a travel grant.

Lorenzo Colitti, Network Engineer at Google said, “This is a great honour. Itojun is a legend in the IPv6 community, and the Internet is indebted to him. Without his foundational work, none of what we achieved with IPv6 would be possible—we stand on the shoulders of giants. Itojun has been a source of inspiration, and I regret never being able to meet him, to show him our work, and show him that we too shared his vision of bringing IPv6 to the users of the Internet.”

Erik Kline, IPv6 Software Engineer at Google said, “It’s humbling to be sharing the Itojun Service Award, having achieved by comparison only a small fraction of the impact of his widely influential body of work. For me personally, Google’s IPv6 efforts are not just for the Internet and its future, but also a way to honor his vision, dedication, and passion.”

More information on the Itojun Service Award is available at: <http://www.isoc.org/itojun>

ISOC Donation to Support Evolution of W3C Organization

ISOC and the *World Wide Web Consortium* (W3C) recently announced a donation from ISOC for the purpose of advancing the evolution of W3C as an organization that creates open Web standards. Citing strongly aligned views on the value of an open global Internet and support for the current Internet governance and management model, ISOC pledged to support W3C efforts to implement a more agile, inclusive, and flexible organizational structure.

“ISOC and W3C have worked together for years in a number of areas, and have deeply shared values about the Internet’s development,” said Lynn St. Amour, President and CEO of ISOC. “Our support to the W3C in their transition efforts demonstrates our commitment to ensuring the Internet continues to be a global platform for innovation. What’s at stake is the Internet’s openness, which is a critical enabler of new products and services to billions of users worldwide.”

“ISOC and W3C have a long history of cooperation and the Internet ecosystem has benefited from our shared yet independent voices,” said Tim Berners-Lee, W3C Director. “The W3C staff, Members, and community continue to work on making W3C more relevant and valuable to the Web and Internet communities. ISOC support will allow W3C to evolve its structure to ensure we continue to forge solid working relationships with the increasing numbers of developers and users, worldwide.”

The two organizations will continue to operate independently, and will maintain their long-standing, informal collaboration. ISOC’s pledge of support is for three years, with both organizations working to ensure progress. A FAQ with additional information is available on both the ISOC site and the W3C site, see <http://www.isoc.org> and <http://www.w3.org>

DNSSEC Deployment in the Root Zone

In December 2009, ICANN and VeriSign began to deploy DNSSEC across the root server system and launched a website that provides information about DNSSEC for the root zone. The website is a repository for the documentation relating to the deployment of DNSSEC, and it includes information such as technical status updates and the full timetable for the deployment of DNSSEC.

See: <http://www.root-dnssec.org/>

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2009 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

March 2010

Volume 13, Number 1

A Quarterly Technical Publication for
Internet and Intranet Professionals

In This Issue

From the Editor	1
Rolling Over DNSSEC Keys	2
Virtual Aggregation	17
RFC Editor	26
Fragments	33
Call for Papers	35

FROM THE EDITOR

Previous articles in IPJ have described *Domain Name System Security Extensions* (DNSSEC), the security system for the *Domain Name System* (DNS). DNSSEC introduces security into the DNS through the use of cryptographic keys and digital signatures. Interest in DNSSEC has grown in recent months, as the *Internet Corporation for Assigned Names and Numbers* (ICANN) and VeriSign have undertaken a phased program to deploy DNSSEC across the root server system in the first half of 2010. In an article by four DNS practitioners, we will explore some side effects of DNSSEC, and examine what happens in two widely used DNS resolver implementations when DNS clients lag behind in synchronizing their local copy of trust keys with the master keys used by the zone administrators to sign their DNS data.

Several articles in IPJ have dealt with various concerns related to scaling of the Internet. In this issue, Paul Francis and Xiaohu Xu describe *Virtual Aggregation*, a new routing technology being developed by the GROW working group of the IETF to reduce the size of the *Forwarding Information Base* (FIB) held in memory by routers.

The *Request For Comments* (RFC) Series has been the main publication channel for Internet standards and related documents for more than 40 years. The RFC Editor function is in the process of being restructured and moved from its original home at the *University of Southern California Information Sciences Institute* (USC/ISI). Leslie Daigle describes the history and future of the RFC Editor mechanism.

If you are reading this online and did not receive the March 2010 edition of IPJ, it may be because your subscription has expired. You can still renew your subscription by visiting the "Subscriber Services" section of our webpage at www.cisco.com/ipj. Enter your subscription ID and e-mail address to gain access to your subscription record. If you don't know your subscription ID or have changed e-mail address recently, just send a message to ipj@cisco.com and we will take care of the renewal and update for you.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

Rolling Over DNSSEC Keys

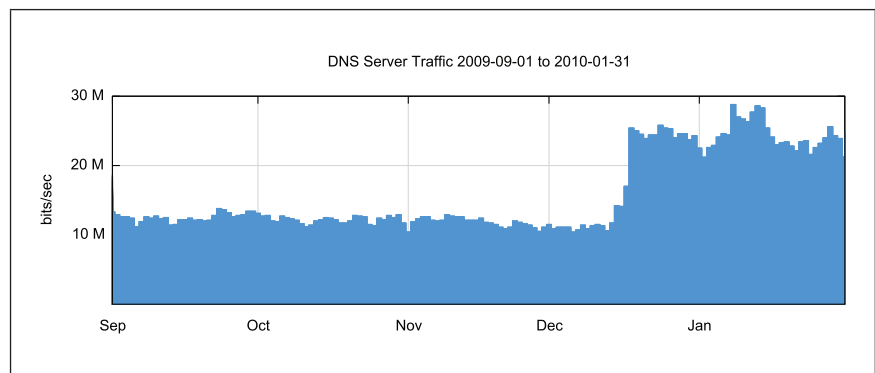
by George Michaelson, APNIC, Patrick Wallström, .SE, Roy Arends, Nominet, Geoff Huston, APNIC

As we are constantly reminded, the Internet can be a very hostile place, and public services are placed under constant pressure from a stream of probe traffic, attempting to exploit any one of numerous vulnerabilities that may be present at the server. In addition, there is the threat of *Denial of Service* (DoS)^[1] attacks, where a service is subjected to an abnormally high traffic load that attempts to saturate and take it down. This story starts with the detection of a possible hostile DoS attack on *Domain Name System* (DNS) servers, and narrates the investigation as to the cause of the incident, and the wider implications of what was found in this investigation.

Detecting the Problem

The traffic signature in Figure 1 is a typical signature of an attempted DoS attack on a server, where the server is subjected to a sudden surge in queries. In this case the traffic log is from a secondary DNS Name Server that is authoritative for a number of subdomains of the `in-addr.arpa` zone; the traffic surge shown here commenced on December 16, 2009. The traffic pattern shifted from a steady state of some 12 Mbps to a new steady state of more than 20 Mbps, peaking at 30 Mbps.

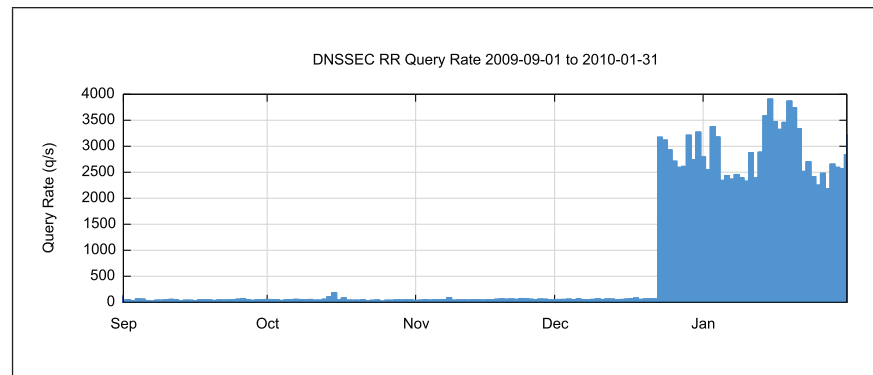
Figure 1: Traffic Load for `in-addr.arpa` Server (provided by George Michaelson)



Because the traffic shown in Figure 1 is traffic passed to and from a Name Server, the next step is to examine the DNS traffic on the Name Server, and in particular look at the rate of DNS queries that are being sent to the Name Server (Figure 2). The bulk of the additional query load is for DNSKEY *Resource Records* (RRs), which are queried as part of the operation of *Domain Name System Security Extensions* (DNSSEC)^[2].

Because this zone is a DNSSEC signed zone, DNSKEY queries will cause the server to respond with a DNSKEY RR and the related RRSIG RR in response to each query. This pair of RRs generates a response that is 1,188 bytes in this case. At a peak query rate of some 3,000 DNS queries per second, a traffic response from the server in excess of 35 Mbps will be generated.

Figure 2: Query Rate for
in-addr.arpa Server
(provided by George Michaelson)



There are many possibilities as to what is going on here:

- This problem could be caused by a DoS attack directed at the server, with the attacker attempting to saturate the server by flooding it with short queries that generate a large response.
- This problem could be caused by a DNS reflection DoS attack, where the attacker is placing the address of the intended victim or victims in the source address of the DNS queries and attempting to overwhelm the victim with this DNS response traffic.

Although it is good to be suspicious, it is also useful to remember the old adage that we should be careful not to ascribe to malice what could equally be explained by incompetence, so numerous other explanations should also be considered, including:

- This problem could be a DNS resolver problem, where the resolver is not correctly caching the response, and some local event is triggering repeated queries.
- This problem could be a bug in an application where the application has managed to wedge itself in a state of rapid-fire queries for DNSKEY RRs.

The next step is to examine some of these queries more closely, and, in particular, look at the distribution of query source addresses to see if this load can be attributed to a small number of resolvers that are making a large number of queries, or if the load is spread across a much larger set of resolvers. The server in question typically sees on the order of 500,000 to 1,000,000 distinct query sources per day.

Closer inspection of the query logs indicates that the additional load is coming from a relatively small subset of resolvers, on the order of 1,000 distinct source addresses, with around 100 “heavy hitters.” In other words, all this DNS traffic is being generated by some 0.01% of the DNS clients. The sequence of queries from one such resolver that is typical of the load being imposed on the server is shown in Figure 3.

Figure 3: DNS Query Sequence
Packet Capture

1. Client requests the Delegation Signer (DS) records for the `211.89.in-addr.arpa` zone.
2. Reply says “no such delegation,” and sends *DNSSEC Signature* (RRSIG) and *Next-Secure record* (NSEC) from the parent zone, and surrounding records.
3. Client requests DNSKEY from the parent zone.
4. Server sends DNSKEY and RRSIG set for the parent zone.

Having established an initial query state and the DNSKEY and signature set over the original request, the client then paradoxically repeatedly re-queries the parent-zone DNSKEY state. This process is elided as follows because the query and response do not differ during this exchange:

5. Client repeats the DNSKEY request.
6. Server repeats the DNSKEY and RRSIG response.
7. Client repeats the DNSKEY request.
8. Server repeats the DNSKEY and RRSIG response.

This exchange of DNSKEY request and DNSKEY and RRSIG response is repeated a further 6 times.

If this additional query load had appeared at the server over an extended period of time, it would be possible to ascribe this problem to a faulty implementation of a DNS resolver, or a faulty client application. However, the sudden onset of the additional load tends to suggest that something else is happening. The most likely explanation is that some external “trigger” event exacerbated a latent behavioral bug in a set of DNS resolver clients. And the most likely external trigger event is a change of the contents of the zones being served.

So we can now refine our set of possible causes to concentrate consideration on the possibility that:

- Something changed in the zones being served by this secondary server that triggered a pathological query response from a set of resolvers.

And indeed the contents of the zones did change on the day when the traffic profile changed, with a key change being implemented on that day.

DNSSEC Key Management

It is considered good operational practice to treat cryptographic keys with a healthy level of respect. As RFC 4641^[3] states: “The longer a key is in use, the greater the probability that it will have been compromised through carelessness, accident, espionage, or cryptanalysis.” Even though the risk is considered slight if you have chosen to use a decent key length, RFC 4641 recommends, as good operational practice, that you “roll” your key at regular intervals. Evidently it is a popular view that fresh keys are better keys.

The standard practice for a “staged” key rollover is to generate a new key pair, and then have the two public keys coexist at the publication point for a period of time. This practice allows relying parties, or clients, some period of time to pick up the new public key. Where possible during this period, signing is performed twice, once with each key, so that the validation test can be performed using either key. After an appropriate interval of parallel operation, the old key pair can be deprecated and the new key can be used exclusively for signing.

This key rollover process should be a routine procedure, without any intended side effects. Resolvers that are using DNSSEC should refresh their local cache of zone keys in synchronization with a published schedule of key rollover, and ensure that they load a copy of the new key within the period when the two keys coexist. In this way when the old key is deprecated, responses from the zone servers can be locally validated using the new key.

The question here is why did this particular key rollover for the signed zone cause the traffic load at the server to spike? And why is the elevated query rate sustained for weeks after the key rollover event? The key had changed 6 months earlier and yet the query load prior to this most recent key change was extremely low.

DNSSEC DNS Resolver Behavior with Outdated Trust Keys

It is possible to formulate a theory as to what is going on from this collection of information. It could be that one or more DNS resolver clients has been using a local *Trust Anchor* that has been manually downloaded from the zone administrator prior to the most recent key rollover, but has not been updated since. When the key rollover occurred in December 2009, these clients could no longer validate the response with their locally stored Trust Anchors.

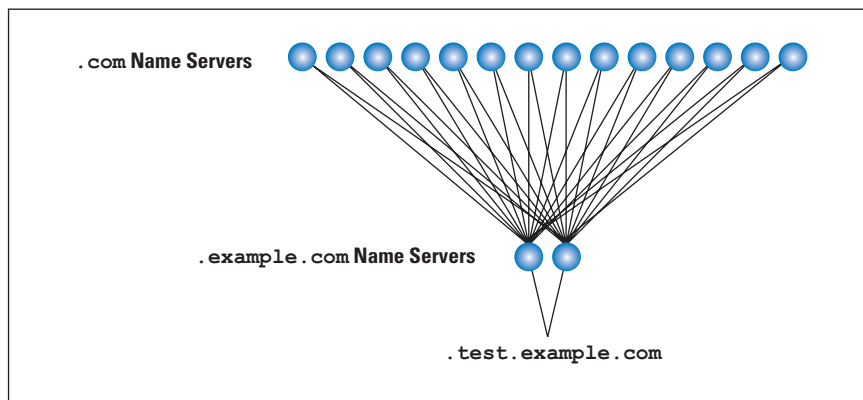
Upon detecting an invalid signature in the response, the client appears to have reacted as if there were a “man-in-the-middle” injection attempt, and immediately repeated the request in an effort to circumvent the supposed attack by rapidly repeating the query. If this instance were really a man-in-the-middle injection attack, this response would be plausible, because there is the hope that the query will still reach the authoritative server and the client will receive a genuine response that can be locally validated.

Why does the client really perform this repeated query pattern? In this case the contributory factor is the use of multiple name servers in the DNS. When the DNS client performs a key validation, it performs a bottom-up search to establish the trust chain from the initial received query to a configured Trust Anchor.

Example DNSSEC Validation

As a hypothetical example, assume a TXT RRset for `test.example.com` in a signed `example.com` zone. The zone `example.com` resides on two Name Server addresses. The `example.com` zone has a *Key Signing Key* (KSK), which is referred to by the DS record in the `.com` zone. The `.com` zone is signed, and it resides on 14 addresses (11 IPv4 and 3 IPv6). The `.com` zone has a KSK, which is referred to by a Trust Anchor in the local configuration of the resolver (Figure 4).

Figure 4: Example Configuration



Assume that the locally held Trust Anchor for `.com` in the resolver has become stale. That is, the DS record for `.com` in the root zone validates, but there are no DNSKEYs in `.com` that match the DS record in the root zone.

When a client is resolving a query relating to `test.example.com`, the following search occurs:

- *Berkeley Internet Name Domain* (BIND)^[9] resolves the `test.example.com` RRset. It attempts to validate it. To do so, it needs the `example.com` DNSKEY RRset.
- It resolves the DNSKEY RRset for `example.com` from a Name Server of `example.com`. It attempts to validate it. To do so, it needs the `example.com` DS RRset.
- It resolves the DS RRset for `example.com` from a Name Server of `.com`. It attempts to validate it. To do so, it needs the `.com` DNSKEY RRset.
- It resolves the DNSKEY RRset for `.com` from a Name Server of `.com`. It attempts to validate it with the locally configured Trust Anchor.

However, the resolver cannot validate the `.com` DNSKEY RRset because it does not have the proper Trust Anchor for it. It queries all remaining 13 `.com` servers for the DNSKEY RRset for `.com`. Then the resolver still does not have the proper `.com` DNSKEY, and tracks back one level:

- It resolves the DS RRset for **example.com** from the next authoritative Name Server. It attempts to validate it. To do so, it needs the **.com** DNSKEY RRset. The search goes forward again.
- It resolves the DNSKEY RRset for **.com**. It attempts to validate it with the locally configured Trust Anchor.

Because the DNSKEY RRset for **.com** has not changed, this attempt will fail as well.

The complete in-depth first search consists of:

- TXT records on 2 **example.com** servers, signed by:
- DNSKEY records on 2 **example.com** servers, referred to by:
- DS records on 14 **.com** servers, signed by:
- DNSKEY records on 14 **.com** servers.

When all possible paths are exhausted, the client will have sent the following:

- 784 ($2 \times 2 \times 14 \times 14$) **.com** DNSKEY requests to 14 **.com** Name Servers
- 56 ($2 \times 2 \times 14$) **example.com** DS requests to 14 **.com** Name Servers
- 4 (2×2) **example.com** DNSKEY requests to 2 **example.com** Name Servers

In other words, in this example scenario with stale Trust Anchor keys in a local client's resolver, a single attempt to validate a single DNS response will cause the client to send a further 844 queries, and each **.com** Name Server to receive 56 DNSKEY RR queries and 4 DS RR queries.

The breadth and level of the search is important here, because the longer the validation chain and the more the number of authoritative Name Servers for those zones that lie on the validation chain path, the more queries that will be sent in an effort to validate a single initial response. In this example, the level of search is three deep, and terminates at **.com**. If the **.com** zone were signed by the root Name Servers and the client were using a stale root zone key, then the 20 distinct root zone server addresses (13 in IPv4 and 7 IPv6 addresses) would also be queried:

- 313,600 ($2 \times 2 \times 14 \times 14 \times 20 \times 20$) root DNSKEY requests to 20 root Name Servers
- 15,680 ($2 \times 2 \times 14 \times 14 \times 20$) **.com** DS requests to 20 root Name Servers

It is worthwhile noting in this context that reverse trees and enum trees in the **.arpa** zone are longer on average. Though delegations in those subtrees might span several labels, it is not uncommon to delegate per label. Note also that the entire effort is done per incoming query—the entire search is repeated for each query.

Though this example shows an enormous query load, there are a few ceilings. In commonly used validating resolvers, such as BIND 9.7rc2, every search is performed in serial, and each search is halted after 30 seconds.

The *Unbound* client^[4] also appears to have a similar request behavior, although it is not as intense because of the cache management in this implementation. Unbound will “remember” the query outcome for a further 60 seconds, so repeated queries for the same name will revert to the cache. But the DNSSEC key validation failure is per zone, and further queries for other names in the same zone will still exercise this re-query behavior. In effect, for a zone that has sufficient “traffic” of DNS load in subzones or instances inside that zone, the chain of repeated queries is constantly renewed and kept alive.

If one such client failed to update its local trusted key set, then the imposed server load on DNSSEC key rollover would be slight. However, if a larger number of clients were to be caught out in this manner, then the load signature of the server would look a lot like Figure 2. The additional load imposed on the server comes from the size of the DNSKEY and RRSIG responses, which are 1,188 bytes per response in the specific failure case that triggered this investigation.

So far we’ve been concentrating attention on the **in-addr.arpa** zone, where the operational data was originally gathered. However, it appears that this problem could happen to any DNSSEC signed domain where the zone keys are published so as to allow clients to manually load them as trust points, and where the keys are rolled on a regular basis.

It is likely that one possible cause for this situation is in the way in which some DNSSEC distributions are packaged with operating systems. For example, the *Fedora*^[5] Linux distribution has bundled numerous trust keys with its packaging of a DNS resolver client and local Trust Anchor key set. When the keys associated with sub zones of **in-addr.arpa** rolled over in December 2009, users of this version of the Fedora Linux distribution would have been caught with stale trust keys.

So there appears to be a combination of three factors that are causing this situation:

- The use of prepackaged DNSSEC distributions that included pre-loaded keys in the distribution
- The use of regular key rollover procedures by the zone administrator
- Some implementations of DNS resolvers that react aggressively when there is a key validation failure by performing a rapid sequence of repeat queries, with either a very slow, or in some cases no apparent back-off in query load

This combination of circumstances makes the next scheduled key rollover for **in-addr.arpa**, scheduled for June 2010, appear to be quite an “interesting” event. If there is the same level of increase in use of DNSSEC with manually managed trust keys over this current 6-month interval as we’ve seen in the previous 6 months, and if the same proportion of clients fails to perform a manual update prior to the next scheduled key rollover event, then the increase in the query load imposed on **in-addr.arpa** servers at the time of key rollover promises to be truly biblical in volume.

Signing the DNS Root

There is an end in sight for this situation for the subzones of **in-addr.arpa**, and for all other such subzones that currently have to resort to various forms of distribution of their zone keys. The *Internet Corporation for Assigned Names and Numbers* (ICANN) has announced that on July 1, 2010, a signed root zone for the DNS will be fully deployed^[6]. Assuming that the **.arpa** and **in-addr.arpa** zones will be DNSSEC-signed in a similar time frame, the situation of escalating loads being imposed on the servers for delegated subdomains of **in-addr.arpa** at each successive key rollover event will be curtailed. It would then be possible to configure the client with a single trust key, the public key signing key for the root zone, and allow the client to perform all signature validation without the need to manually manage other local trust keys.

There are two potential problems with this scenario.

The first is that for those clients that fail to remove the local Trust Anchor key set, these repeated queries may not go away. When there are multiple possible chains of trust, the resolver will attempt to validate using the shortest validation chain. As an example, if a client has configured the DNSKEY for, say, **test.example.com** into its local Trust Anchor key set, and it then subsequently adds the DNSKEY for **example.com**, the resolver client will attempt to validate all queries in **test.example.com** and its subzones using the **test.example.com** DNSKEY.

A more likely scenario is where an operator has already added local Trust Anchor keys for, say, **.org** or **.se**. When the root of the DNS is signed, the operator may also add the keys for the root to the local Trust Anchor set. If the operator fails to remove the local copies of the **.org** and **.se** Trust Anchor keys, in the belief that this root key value will override the **.org** and **.se** local keys, then the same validation failure behavior will occur. In such a case, when the local keys for these second-level domains become stale, their resolver will exhibit the same re-query behavior, even when they maintain a valid local root Trust Anchor key.

As a side note, the same behavior may occur when *DNSSEC Lookaside Validation* (DLV) is used. If the zone key management procedures fall out of tight synchronization with the DLV repository, it is possible to open a window where the old key remains in the DLV repository, but is no longer in the zone file. This situation can lead to a window of vulnerability where the keys in the DLV repository are unable to validate the signed information in the zone file, a situation that, in turn, introduces the same problem with re-query.

The second potential problem lies with the phase-in approach of signing the root. The staged rollout of DNSSEC for the root zone envisages a sequenced deployment of DNSSEC across the root server clusters, and through this sequence the root will be signed with a key that has no valid published public part, creating a *Deliberately Unvalidatable Root Zone* (DURZ).

What happens when a client installs this key in its local Trust Anchor set and performs a query into the root zone?

As an experiment, this DURZ key was installed into an instance of BIND 9.7.0rc2, with a single upstream root, pointing at the “L” root, the only instance of the 13 authoritative root servers enabled with DNSSEC signed data in February 2010. On startup the client made 13 consecutive DNSKEY requests, one to each of the root zone server addresses. When the client started its first query in a subzone, the client issued a further 156 DNSKEY queries in a period of 19 seconds, making 12 queries to each of the 13 root zone server addresses.

This scenario should sound familiar, because it is precisely the same query pattern as happened with the `in-addr.arpa` servers and the `.se` servers, although the volume of repeated DNSKEY queries is somewhat alarming. When the client receives a response from a subdomain that needs to be validated against the root, and when the queries to the root are not validatable against the local trust key, the client goes into a sequence of repeated queries that explore each potential validation path. Anchoring the local resolver with a key state that invalidates the signatures of all authoritative servers of the zone—but authoritatively (absent DNSSEC) confirms them as valid servers of the zone—places the client instance in an unresolvable situation: no authoritative Name Server that it can query has a signature that the client can validate, but the root zone informs it that only these Name Servers can be used.

Further tests of this behavior show that the client does not cache the outcome that the DNSKEY cannot be validated for a zone, and the client reinitiates this spray of repeated queries against the zone Name Servers when a subsequent DNSSEC query is made in a subzone. Therefore the behavior is promiscuous in two distinct ways. First it is evident that any Name Server so queried is repeatedly queried. Second, it is evident that all Name Servers of a zone are queried. The other part of the client response is not to cache validation failure for the zone in case this repeated query phase does not provide the client with a locally validated key.

After all, the data is provably false, so caching it would be to retain something that has been “proven” to be wrong.

The emerging picture is that misconfigured local trust keys in a DNS resolver for a zone can cause large increases in the DNS query load to the authoritative Name Servers of that zone, where the responses to these additional queries are themselves large, of the order of 1,000 bytes in every response. This situation can occur for any DNSSEC signed zone.

The conditions for the client to revert to a rapid re-query behavior follow:

- The *DNSSEC OK* (DO) bit is honoured by the server.
- The DNS data appears to be signed.
- The signature check fails.
- The client does not cache the validation failure for this zone.

The conditions being set up for the DURZ approach for signing the root follow:

- The DO bit is honoured by the server.
- The DNS data appears to be signed.
- The signature check fails.
- The client does not cache the validation failure for this zone.

What is to stop the DNS root servers from being subjected to the same spike in the query load?

The appropriate client behavior for this period of DNSSEC deployment at the root is not to enable DNSSEC validation in the resolver. Although this advice is sound, it is also true that many resolvers have already enabled validation in their resolvers, and are probably not going to turn off for the next 6 months while the root servers gradually deploy DNSSEC using DURZ.

But what load will appear at the root servers if a subset of the client resolvers starts to believe that these unvalidatable root keys should be validated?

What If...?

The problem with key rollover and local management of trust keys appears to be found in around 1 in every 1,500 resolvers in the **in-addr.arpa** zones. With a current client population of some 1.5 million distinct resolver client addresses each day for these **in-addr.arpa** zones, there are some 1,000 resolvers who have lapsed into this repeated query mode following the most recent key rollover of December 2009. Each subzone of **in-addr.arpa** has six Name Server records, and all servers see this pathological re-query behavior following key rollover.

The root servers see a set of some 5 million distinct resolver addresses each day, and a comparable population of nonupdated resolvers would be on the order of some 3,000 resolvers querying 13 zone servers, where each zone server would see an incremental load of some 75 Mbps.

Because the re-query behavior is caused by the client's being forced to reject the supposedly authoritative response because of an invalid key, and because DURZ is by definition an invalid key, the risk window for this increased load is the period during which DURZ is enabled, which for the current state of the root signing deployment is from the present date until July 2010. Because not all root servers have DNSSEC content or respond to the DO bit—and therefore do not return the unvalidatable signatures—the risk is limited to the set of DNSSEC-enabled roots, which is increasing on a planned, staged rollout. It has been reported that a decision to delay deployment of the DNSSEC/DURZ sign state to the “A” root server instance was made because this root server receives a noted higher query load for the so-called “priming” queries, made when a resolver is reinitialized and uses the offline root “hints” file to bootstrap more current knowledge. It is therefore likely that the “A” root server would also see increased instances of this particular query model, if the priming query is implicated in this form of traffic.

Arguably, this situation is unlikely. For most patterns of DNS query, failure to validate is immediately apparent. After all, where previously you receive an answer, you now see your DNS queries time out and fail.

However, because the typical situation for a client host (including *Dynamic Host Configuration Protocol* [DHCP] initialized hosts in the customer network space, the back office, etc.) is to have more than one listed resolver, there is the possibility of a misconfiguration being unnoticed during the period of a rolling deployment of DNSSEC-enabled services. In this situation if only one of the resolver's “nserver” entries is DNSSEC-enabled, either it is not queried or it is queried, but then passed over by the resolver timeout setting. Users see slower DNS resolution, but can attribute it to network delay or other local problems.

A second argument is that installation of hand-trust material is not normal, so the servers in question will be immediately known because a nonstandard process has to be invoked. Unfortunately, this situation is demonstrably not true. For example, the *Fedora*^[5] release of Linux has included a simple DNSSEC-enabling process including a preconfigured trust file covering the reverse-DNS ranges. Because a previous release of this software included now stale keys (which have since been withdrawn in subsequent releases), any instance of *Fedora* for this release state being enabled will not only be unable to process reverse-DNS, it may also invoke this re-query mode of operation that places the server under repeated load of DNSKEY requests.

Because reverse-DNS is the “infrastructure” DNS query that is typically logged, but not otherwise used, unless the server in question is configured to block service on failing reverse (unlikely, given that more than 40 percent of reverse-DNS delegations are not made for the currently allocated IP address ranges), the end user simply might never notice this behavior. The use of so-called “Live CDs” can exacerbate this problem of pre-primed software releases that include key material that falls out-of-date. Even when the primary release is patched, the continued use of older releases in the field is inevitable. So perhaps this second argument is not quite as robust as originally thought.

Lastly, distinct from hand-installed local trust is the use of DNSSEC look-aside validation, which is known as DLV. This DNS namespace is privately managed and has been using the ICANN-maintained *Interim Trust Anchor Repository*, or ITAR. The DLV service is configured to permit resolvers to query it, in place of the root, to establish trust over subzones that exist in a signed state, but cannot be seen as signed from the root downward before the deployment of a signed root. There is now evidence that part of this query space exists, covering zones of interest to this situation. The `.se` zone key, for instance, is in the ITAR, as are the `in-addr.arpa` spaces signed by the RIPE NCC. Evidence suggests that if the DLV chain is being used and a key rollover takes place, some variants of BIND resolver clients fail to reestablish trust over the new keys until the client is rebooted with a clean cache state. This theory is difficult to confirm because as each resolver is restarted, the stale trust state is wiped out and the local failure is immediately resolved.

Post DURZ

Of course this phase is transitory, and even if there are concerns in terms of DURZ and queries to the root servers, all will be resolved when the root key is rolled to a validatable key on July 1, 2010.

Yes? Maybe not.

The current plan is to roll the root zone Key Signing Key every 2 to 5 years. The implication is that sometime every 2 to 5 years all DNS resolvers will need to ensure that they have fetched a new root trust key and loaded it into their resolver’s local trust key cache.

If this local update of the root trust key does not occur, then the priming query for such DNSSEC-enabled resolvers will encounter this problem of an invalid DNSKEY when attempting to validate the priming response from the root servers. The fail-safe option here for the resolver client is to enter a failure mode and shut down, but there is a strong likelihood that the resolver client will try as hard as it can to fetch a validatable DNSKEY for the root before taking the last resort of a shutdown, and in so doing will subject the root servers to this intense repeated query load that we are seeing on the `in-addr.arpa` zone.

A reasonable question to ask follows: “Are there any procedural methods to help prevent stale keys from being retained during key rollover?” Reassuringly, the answer is “Yes.” There is a relatively recent RFC, “Automated Updates of DNS Security (DNSSEC) Trust Anchors,” RFC 5011^[7], which addresses this problem.

RFC 5011 provides a mechanism for both signaling that a key rollover needs to take place and forward declaring the use of keys to sign over the new trust set to permit in-band distribution of the new keys. Resolvers are required to be configured with additional keying, and a level of trust is placed on this mechanism to deal with normal key rollover. RFC 5011 does not solve initial key distribution problems, which of course must be made out of band, nor does it attempt to address multiple key failures. Cold standby equipment, or decisions to return to significantly older releases of systems (for example, if a major security compromise to an operating system release demands a rollback) could still potentially deploy resolvers with invalid, outdated keys. However, RFC 5011 will prevent the more usual process failures, and it provides an elegant in-band rekeying method that obviates a manual process of key management that all too often fails through neglect or ignorance of the appropriate maintenance procedures to follow.

It is unfortunate that RFC 5011-compliant systems are not widely deployed during the lifetime of the DURZ deployment of the root, because we are definitely going to see at least one key rollover at the end of the DURZ deployment, and we can expect a follow-up key rollover within a normal operations window. The alternative is that no significant testing of root trust rollover takes place until we are committed to validation as a normal operational activity—a situation that invites the prospect of production deployment across the entire root set while many production operational processes associated with key rollover remain untested. The evidence from past concerns in resolver behavior is that older deployments have a very long lifetime for any feature under consideration, and because BIND 9.5 and older prerelease BIND 9.7 systems can be expected to persist in the field in significant numbers for some years to come, it is likely a significant level of pathological resolver behavior in re-querying the root services by active resolvers will have to be tolerated for some time.

It is also concerning that aspects of the packet traces for the **in-addr.arpa** zone suggest that for all key rollovers, albeit at very low levels of query load, some of the resolvers have simply failed to account for the new keys—and may never do so. Therefore, with increasing deployment of key validation, it is possible that a substantial new traffic class that grows, peaks, and then declines, but always declines to a slightly higher value than before, has to be borne, and factored into deployment scaling and planning.

Because this traffic is large—generating a kilobyte of response per query and potentially generally prevalent—it has the capability to exceed the normal response requirements for “normal” DNS query loads by at least one, if not two orders of magnitude. This multiplication factor of load is defined by the size of the resolver space and the number of listed Name Servers for the affected zone.

Mitigation at the server side is possible if this problem becomes a major one. The pattern of re-query here (the sequence of repeated queries for DNSKEY RRs) appears a potential signature for this kind of problem. Given that for any individual server the client times its repeat queries on the reception of the response from the previous query, delaying the response of the server to the repeated query will further delay the client’s making its repeated query to this server. If the server were in a position to delay such repeated responses, using a form of exponential increase in the delay timer or similar form of time penalty, then the worst effects of this form of client behavior in terms of threats to the integrity of the ability of the server to service the “legitimate” client load could be mitigated.

Conclusion

It is an inherent quality of the DNSSEC deployment that in seeking to prevent lies, an aspect of the stability of the DNS has been weakened. When a client falls out of synchronization with the current key state of DNSSEC, it will mistake the current truth for an attempt to insert a lie. The subsequent efforts of the client to perform a rapid search for what it believes to be a truthful response could reasonably be construed as a legitimate response, if indeed this instance was an attack on that particular client. Indeed, to do otherwise would be to permit the DNS to remain an untrustable source of information. However, in this situation of slippage of synchronized key state between client and server, the effect is both local failure and the generation of excess load on external servers—and if this situation is allowed to become a common state, it has the potential to broaden the failure state to a more general DNS service failure through load saturation of critical DNS servers.

This aspect of a qualitative change of the DNS is unavoidable, and it places a strong imperative on DNS operations and the community of the 5 million current and uncountable future DNS resolvers to understand that “set and forget” is not the intended mode of operation of DNSSEC-equipped clients.

For Further Reading

- [0] A longer version of this article can be found in our online companion publication, *The Internet Protocol Forum*,
<http://www.ipjforum.org/?p=226#more-226>
- [1] Charalampos Patrikakis, Michalis Masikos, and Olga Zouraraki, “Distributed Denial of Service Attacks,” *The Internet Protocol Journal*, Volume 7, No. 4, December 2004.

- [2] Miek Gieben, “DNSSEC: The Protocol, Deployment, and a Bit of Development,” *The Internet Protocol Journal*, Volume 7, No. 2, June 2004.
- [3] O. Kolkman and R. Gieben, “DNSSEC Operational Practices,” RFC 4641, September 2006.
- [4] <http://www.unbound.net>
- [5] <http://fedoraproject.org>
- [6] <http://www.root-dnssec.org>
- [7] M. St. Johns, “Automated Updates of DNS Security (DNSSEC) Trust Anchors,” RFC 5011, September 2007.
- [8] Geoff Huston, “Resource Certification,” *The Internet Protocol Journal*, Volume 12, No. 1, March 2009.
- [9] <https://www.isc.org/software/bind>
- [10] <https://www.isc.org/community/blog/201002/signed-root-coming-and-what-means-you>

GEORGE MICHAELSON has a B.Sc. from the University of York and is a research scientist at the *Asia Pacific Network Information Centre* (APNIC), the Regional Internet Registry serving the Asia Pacific region. George explores problems in Internet Number Resource management, Internet standards, and network measurement by collaborative research. George has more than 28 years experience in computer science, networking, ICT administration, and research conducted in Australia and the UK. He participates in standards development in the IETF and has been a working group chair as well as an RFC author. He is a member of the *British Computer Society*. E-Mail: ggm@apnic.net

PATRIK WALLSTRÖM is a senior researcher at .SE, the Internet registry for SE domain names, and has been with the registry for eight years developing registry systems and working with the deployment of DNSSEC. Patrik is currently working on the *OpenDNSSEC* project, producing tools for a wider deployment of the technology. At .SE he is also managing the *Healthcheck* project, a new open source platform for measuring the quality of DNS, E-mail, Web and IP within Sweden. Patrik is also a board member of the *Swedish Network Users' Society* (SNUS). E-mail: pawal@iis.se

ROY ARENDS is a senior researcher at Nominet UK, the Internet registry for UK domain names. He co-authored several IETF standards on DNSSEC, resides on the board of DNS-OARC, is a member of ICANN's *Security and Stability Advisory Committee*, and is part of IETF's DNS-Directorate. As an expert on DNS and DNSSEC, Roy has co-initiated several DNS-related open source projects, such as *Unbound* and *OpenDNSSEC*. In the past, Roy was a member of, and chaired, CERT-NL. E-mail: roy@nominet.org.uk

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. The author of numerous Internet-related books, he is currently the Chief Scientist at APNIC. He was a member of the *Internet Architecture Board* (IAB) from 1999 until 2005, and served on the Board of the Internet Society from 1992 until 2001. E-mail: gih@apnic.net

Extending Router Lifetime with Virtual Aggregation

by Paul Francis, Max Planck Institute for Software Systems, and Xiaohu Xu, Huawei Technologies

Biologists believe that human life is limited by the number of times cells can replicate; noncancerous cells have a kind of internal counter that prevents them from replicating forever. Even if humans are kept healthy in every respect, they will eventually die simply because their cells will cease to replicate. Internet routers also have a finite lifetime. They are built with a fixed amount of hardware memory for storing the forwarding table (the memory structure that tells the router where to forward any IP packet, also called the *Forwarding Information Base* [FIB]). As the Internet global routing table grows, it eventually overflows the FIB, and the router ceases to be able to hold the full routing table. Even if the router is healthy in every respect (all of its hardware components still operate), it can no longer function as a router in the Internet *Default-Free Zone* (DFZ), where no default routes can be used.

In the past, router vendors have been reasonably good at predicting how long FIBs will last because the growth of the global DFZ routing table has stayed fairly predictable. As a result, *Internet Service Providers* (ISPs) can plan their capital budgets, and where necessary use a set of tricks (discussed in the next section) to squeeze additional life out of routers even after their “FIB death.” But there are two problems.

First, these tricks work only in limited situations, they require extra configuration, and they can lead to increased traffic loads. Second, and potentially much more serious, the rate of routing table growth may dramatically accelerate in the near future, thus shrinking the lifetime of the installed router base. This expected acceleration is due to the imminent exhaustion of IPv4 addresses. In the past, address authorities such as the *American Registry for Internet Numbers* (ARIN) could assign large contiguous blocks of addresses to ISPs, which in turn assigned smaller blocks to their customers. Therefore, routers in other ISPs’ networks need only a single routing table entry—that of the large block—to route to destinations in the ISP. This approach to scaling is called *address aggregation*. There is a fear that, as IPv4 addresses become increasingly unavailable, ISPs will start buying and selling smaller and smaller blocks of IP addresses from each other in an effort to squeeze out as many addresses as possible. These small blocks will appear all over the Internet thus significantly increasing the size of the routing table.

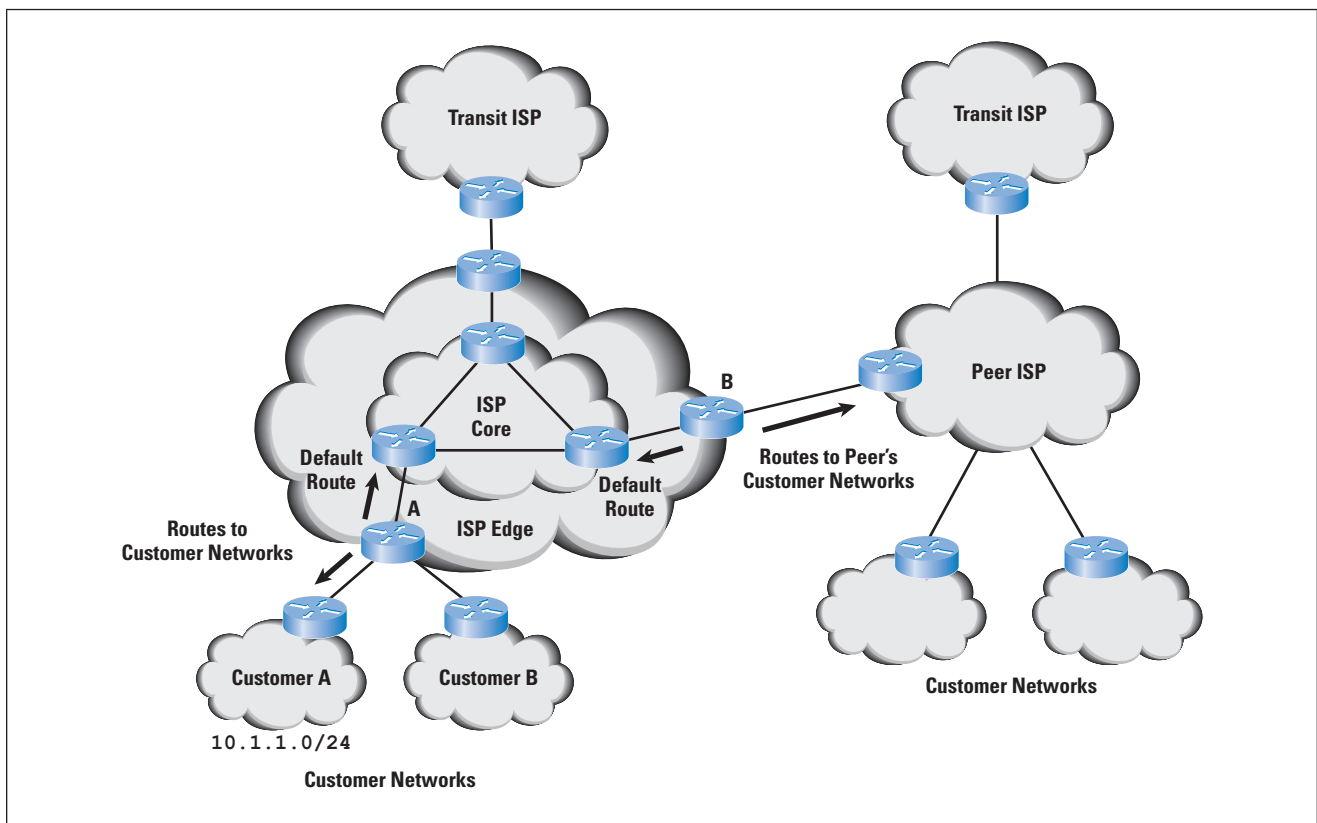
This article describes a new routing technology, called *Virtual Aggregation* (VA), which mitigates these problems. It makes extending the lifetime of old routers much easier, and makes it possible for existing routers to absorb a surge in the routing table size. Virtual Aggregation is a working item in the *Global Routing Operations Working Group* (GROW) working group of the IETF^[7], and is documented in `draft-ietf-grow-va`^[6] and related drafts.

Tricks for Keeping Old Routers Deployed

ISPs frequently want to extend the usefulness of a router beyond its “FIB death,” and there are many tricks for doing just this. The most common is to structure the ISP in a core-edge arrangement. In this setup, a core of routers forms the backbone of the network. Edge routers connect to other networks and feed into the core. In many cases these edge routers do not need to know how to route to everything in the Internet. Rather, they often need to know only what addresses are reachable in their directly connected networks.

For instance, Figure 1 shows an ISP whose edge routers connect to three types of other networks: customer networks, peer ISP networks, and transit ISP networks. Each customer network has only one or a small number of address prefixes. The edge routers connecting customer networks must know what addresses are reachable in the customer networks, but everything else can be “default routed” to the core. Likewise, the routers connected to peer ISPs need to know how to route to the peers’ customer addresses. Everything else can be defaulted to the core. The core routers and the edge routers that connect to transit ISPs, however, need to know how to route to everything.

Figure 1: With a core-edge style of deployment, some routers need to keep full routing tables, while others can keep partial routing tables and default route everything else to the ISP core.



A common practice is for ISPs to delegate FIB-dead routers to the customer or peer edges, and to have the core routers filter the routing information given to the edge routers. For instance, router A in Figure 1 learns the addresses reachable in customer network A (say, 20.1.1.0/24) and conveys them to the ISP core, but the core tells router A only that “everything else” is reachable through it (0.0.0.0/0). But what if customer A itself wants the full DFZ routing table? For instance, customer A might be multihomed to some other ISP, and might want to know which Internet destinations are best reachable through each ISP. To do this, it needs to receive the whole routing table from each ISP, a situation that, of course, cannot happen if the core withholds routes from router A.

As another example, what if two peer ISPs later decide that they want to offer transit service to each other? Now additional routes need to be conveyed to the peer-connected edge routers (router B), and this process may not be possible with limited FIB.

Another way an ISP can shrink its routing table is to default route to its transit ISPs. For instance, routers keep track only of how to route to customers and peers, and everything else is defaulted to the transit ISPs. When this default routing is done, even an ISP’s core routers do not need the full routing table. A simple approach is for an ISP to send all defaulted packets to the nearest transit ISP. This process, however, may result in many packets taking a longer Internet path than necessary. Reference [1] describes a more complicated approach where the ISP maintains “semidefaults” for different transit networks in order to improve its global routing while reducing routing table size by about half. This approach, however, can be hard to manage.

In addition, any form of ISP-level default (simple or complex) results in sending extra traffic to the transit ISPs. A substantial amount of Internet traffic is targeted to nonroutable prefixes. When an ISP has the full routing table, it can identify this traffic and drop it before sending it to its transit ISPs. When an ISP defaults, it sends this traffic to its transit ISPs, and pays for it.

To summarize, dealing with FIB-dead routers leads to more complex management, limitations in business arrangements with peers and customers, poor paths over the Internet, and increased traffic load.

The Idea of Virtual Aggregation

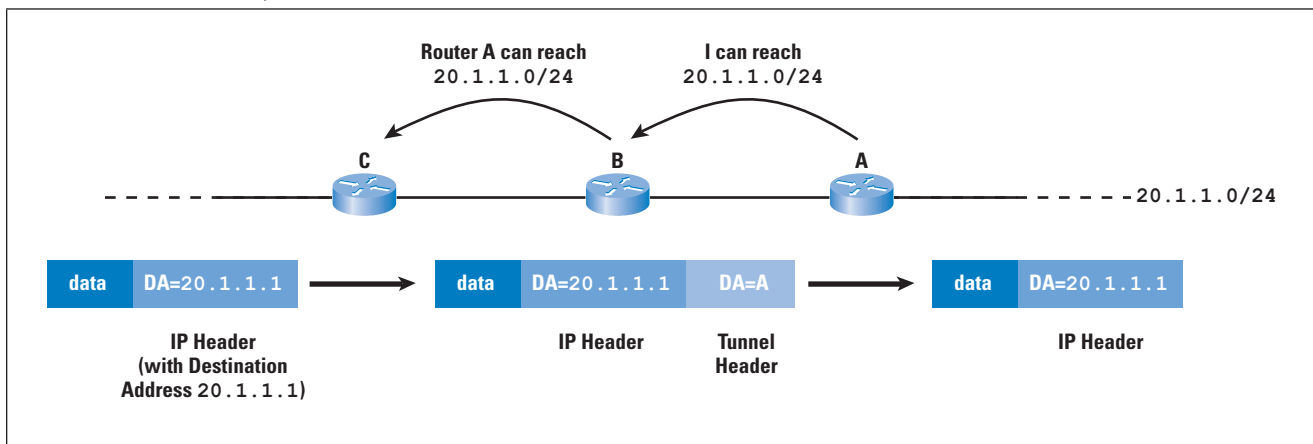
In its simplest form, Virtual Aggregation allows an ISP to use FIB-dead routers as edge routers, in any edge router position (neighbor is a transit provider, a peer, or a customer) without limiting what routing information is exchanged. Configuration requirements are minimal. In a more complex form, Virtual Aggregation allows all ISP routers (not just edge routers) to be FIB-dead routers, without requiring ISP-level default routing.

Virtual Aggregation uses two basic mechanisms, FIB suppression and tunneling. Before discussing FIB suppression, a small amount of background is needed. Internet routers have a “data plane” and a “control plane.” The data plane is what forwards packets, and includes such functions as header parsing, FIB lookup, queuing, and packet transmission. The control plane operates the background protocols that gather much of the information needed by the data plane. Examples include routing protocols such as the *Border Gateway Protocol* (BGP) and *Open Shortest Path First* (OSPF), and tunnel establishment protocols such as the *Label Distribution Protocol* (LDP).

The idea of FIB suppression is that the control plane operates as normal, but that certain routing table entries are not loaded into the FIB. This idea exploits the fact that it is (data plane) FIB memory, not control plane routing table memory that is the more severe bottleneck. By allowing the control plane to operate as normal, no changes are required to routing protocols or, for the most part, the management of routing protocols.

Tunneling is used to pass packets through routers that have suppressed FIB entries. The principle is illustrated in Figure 2. Here router A tells router B that it can reach `20.1.1.0/24`. Router B in turn tells router C that router A can reach `20.1.1.0/24`. As a result, router C tunnels packets destined for `20.1.1.0/24` to router A through router B. In other words, it wraps the IP header in another IP or a *Multiprotocol Label Switching* (MPLS) header that first gets the packet to router A. Router A strips that header, and sends the packet toward the destination. Notice that router B can suppress the route to `20.1.1.0/24` from the FIB—it only needs to know how to route the packet to router A. In other words, even though router B fully participates in the control plane, it is able to shrink its FIB through FIB suppression and tunneling.

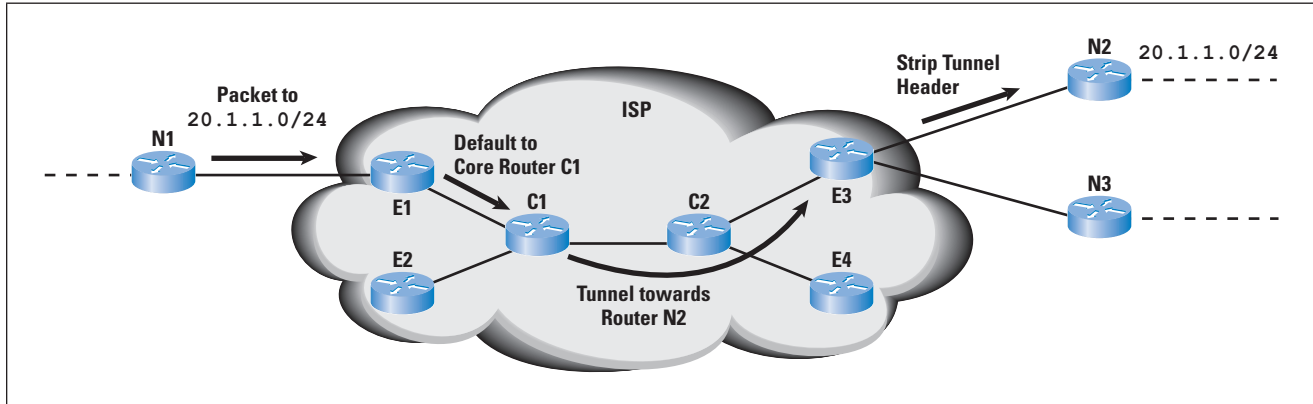
Figure 2: Because router C tunnels the packet to router A, router B does not need to know how to forward packets with addresses in `20.1.1.0/24`.



Virtual Aggregation in Practice, Simple Version

In the simplest version of Virtual Aggregation, a core-edge configuration is used. The core routers maintain full FIB tables. The edge routers FIB-install at least a default route to the core, and potentially additional routes if there is space in the FIB. This process is illustrated in Figure 3. Here there are two core routers, C1 and C2, and four edge routers, E1, E2, E3, and E4. The edge routers have external neighbors, N1, N2, and N3, as shown.

Figure 3: Packets can be delivered to `20.1.1.0/24` even if none of the edge routers has a FIB entry for `20.1.1.0/24`.



The operation is best explained by example. Suppose that N2 advertises a route to destination `20.1.1.0/24` to E3 using *External BGP* (eBGP) and giving itself as the next hop. E3 in turn advertises this route to the other internal routers using *Internal BGP* (iBGP), with the next hop still as N2. The core routers install this route in their FIBs, with an indication that packets matching the route should be tunneled to the next hop, N2. Assume for now that all edge routers FIB-suppress the entry. When a packet for say `20.1.1.1` arrives at E1 from N1, E1 does not find an entry for `20.1.1.0/24`, but does find the default route `0/0` telling it to forward the packet to its core router C1. C1 looks into its FIB and indeed finds an entry for `20.1.1.0/24` telling it to tunnel the packet to N2. C1 wraps the packet in another header, typically IP or MPLS, addressed to N2. When the packet reaches E3, however, E3 notes that the header directs it to send the packet to N2, strips off the outer header, and sends the packet to N2. E3 can do this without a FIB entry for `20.1.1.0/24`.

MPLS already has all the mechanisms needed to perform this packet forwarding. E3 can use LDP to signal a *Label Switched Path* (LSP) to N2, and *Penultimate Hop Popping* can be used to strip off the MPLS header before forwarding the packet to the external neighbor N2 (as described in section 4.1.4 of [4]).

Alternatively, stacked MPLS label technology can be used; for example, the inner label is signaled with BGP (see “Carrying Label Information in BGP-4”^[3]) while the outer label is signaled with LDP. Here E3 sets itself as the next hop for all the routes learned from external neighbors (for example, **20.1.1.0/24**) when advertising them to its iBGP peers, and uses the inner label to identify the external neighbor (see section 4.3, “Label Stacks and Implicit Peering” of [4]). IP-in-IP tunneling can also be used, in this case signaled with softwires BGP attributes^[5].

Now let’s see what happens if a packet to **20.1.1.1** is received by E3 from external neighbor N3. If E3 has not FIB-installed the route for **20.1.1.0/24**, it uses its default entry and forwards the packet to C2. C2 finds its entry for **20.1.1.0/24**, which instructs it to tunnel the packet to N2. The packet is sent back to E3, which strips off the outer header and delivers the packet to E2. In this case, the packet has traveled an extra hop and back, a process that is not acceptable if done too much. As long as there is space in the FIB, however, routers are free to FIB-install additional routes. A good policy is to always install routes when external neighbors are the next hop. This policy avoids the longer path. In some cases, such as edge routers that connect to transit networks, there may not be enough FIB space to hold all routes from all external neighbors. In this case, the router may FIB-install the routes for which the most traffic is forwarded. Studies have shown that a small number of routes account for majority of the traffic, making Virtual Aggregation a very efficient solution^[2].

Note that this simple form of Virtual Aggregation is very easy to configure. Essentially all that is needed is to tell the routers that they are using simple Virtual Aggregation, and to tell them if they are a core or an edge router. The routers can automatically configure everything else. Virtual Aggregation requires configuration of tunnels from every router to every other router, but these configurations also can be automatic. In any event, increasingly these tunnels are created anyway for the purpose of traffic engineering.

Simple Virtual Aggregation solves most of the problems described earlier. It can save FIB on any edge router without having to compromise BGP service to customers or flexibility in using peer networks for some transit. It also allows FIB-dead routers to be used as edge routers with transit ISPs. Finally, it prevents the need for ISP-level default routing to transits, thus avoiding unnecessarily sending unroutable traffic to the transit. And it does all this with much less configuration than is required to operate with FIB-dead routers today.

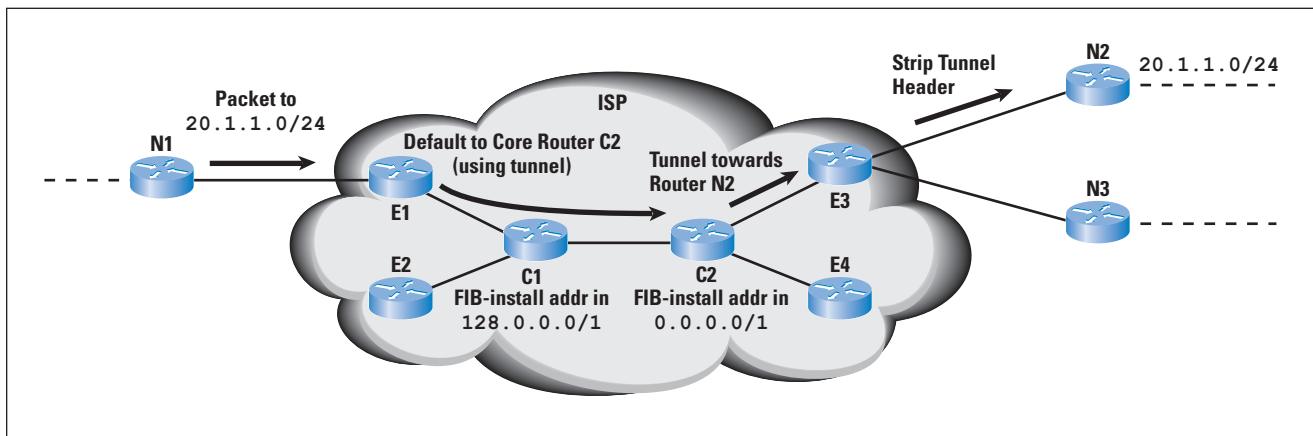
Virtual Aggregation in Practice, Complex Version

The simple version of Virtual Aggregation is satisfactory for edge routers, but it does nothing to reduce FIB size on core routers. What if an ISP wishes to also extend the lifetime of its core routers? Or wants to move away from a core-edge model, and rather connect all edge routers directly through a Layer 2 substrate like MPLS?

What if indeed there is a surge in routing table growth, thus causing ISPs all over the world to suddenly find themselves FIB-starved? There is a version of Virtual Aggregation that allows for FIB reduction in any and all routers in an ISP network.

The basic idea is to divide the address space so that different routers maintain full routes within different parts of the address space. So for instance, rather than have all core routers responsible for all of the address space, you could have half of the core routers responsible for the lower half of the address space, and the other half of the core routers responsible for the upper half of the address space. Figure 4 shows how this setup would look for the simple topology of Figure 3, keeping in mind that this example is rather simplistic.

Figure 4: In a complex version of Virtual Aggregation, even core routers do not need to hold the full routing table.



Assume that C2 FIB-installs only the lower half of the address space (0.0.0.0/1) and C1 FIB-installs the upper half (128.0.0.0/1). With this arrangement, the edge routers have two defaults instead of one. Packets to addresses in 0.0.0.0/1 are defaulted, through a tunnel, to C2, and packets to addresses in 128.0.0.0/1 are defaulted to C1. These defaults are learned simply by having C1 and C2 advertize their respective default routes with themselves as the next hop in iBGP.

As with the previous example, assume that router N2 advertises a route to 20.1.1.0/24, with itself as the next hop, to E3. E3 advertises this route to all other routers using iBGP. Only C2, however, FIB-installs this route—C1 suppresses it. When a packet to 20.1.1.1 arrives at E1, it looks in its FIB, finds a matching route to 0.0.0.0/1, and so tunnels the packet to C2. C2 terminates the tunnel, finds its FIB entry for 20.1.1.0/24, and tunnels the packet toward N2. E3 uses the tunnel information to know to forward the packet to N2, strips away the tunnel header, and forwards the packet to N2.

Now suppose that a packet for 20.1.1.1 arrives at E3 from N3. Ideally E3 has already automatically FIB-installed the route for 24.1.1.0/24 either because its external neighbor provides the next hop, or because the route is a high-volume destination. In this case, of course, the packet is directly forwarded to N2. However, if E3 has not FIB-installed the route, then its best match is the default to 0.0.0.0/1, and it tunnels the packet to C2. C2 in turn tunnels the packet back toward N2 through E3 as described before. Worse, if C1 rather than C2 FIB-installed the lower half of the address space, the packet would have detoured all the way to C1. Clearly these routes are not optimal, and so we must ask how nonoptimal would the complex version of Virtual Aggregation be in real ISPs.

The USENIX NSDI paper^[2] answers this question for one large transit ISP. In this study, both the topology and the traffic matrix of the ISP are considered. The deployment strategy is substantially more complex than the simplistic example given previously. An upper limit is placed on the maximum increase in latency (5 ms) for any path through the ISP. There is a requirement that within a *Point of Presence* (PoP) at least two routers must cover the same address space. The number and size of address partitions are engineered to spread FIB load evenly. The “additional” routes installed in the FIB are designed to cover high-traffic destinations to the extent possible.

With these requirements in mind, this study found that FIB size could be reduced in all routers by at least an order of magnitude with a negligible increase (1–2%) in overall traffic load due to the occasional extra hops from the detours. This result ultimately translated into an increased router lifetime of easily 10 years.

The management requirements for the complex version are substantially greater than those for the simple version. The address partitions must be chosen, the routers assigned to address partitions must be chosen, and possibly some strategy for deciding what “additional” routes should be FIB-installed is needed. Whether this added configuration and the associated difficulties due to, for instance, misconfiguration are worth the cost savings for extending router lifetime is up to each ISP. Virtual Aggregation at least provides an option that was not previously available.

Status

Virtual Aggregation is a working-group item in the *Global Routing Operations Working Group* (GROW) in the IETF. The primary draft is **draft-ietf-grow-va**^[6]. This draft has gone through several revisions, and is very close to its final form. Huawei is currently implementing Virtual Aggregation. A second open-source implementation has been built by Paul Francis’ research group for the *Quagga* open-source routing platform, and is still being enhanced.

Acknowledgements

The authors would like to thank the co-authors of the Virtual Aggregation drafts, Hitesh Ballani, Dan Jen, Robert Raszuk, and Lixia Zhang. In particular, it was Robert who suggested the simple version of Virtual Aggregation.

References

- [1] Andre Chapuis, “BGP Filtering,” Presentation from SWINOG7, www.swinog.ch/meetings/swinog7/BGP_filtering-swinog.ppt
- [2] Hitesh Ballani, Paul Francis, Tuan Cao, and Jia Wang, “Making Routers Last Longer with ViAggre,” USENIX NSDI 2009, April 2009.
- [3] Y. Rekhter and E. Rosen, “Carrying Label Information in BGP-4,” RFC 3107, May 2001.
- [4] E. Rosen, A. Viswanathan, and R. Callon, “Multiprotocol Label Switching Architecture,” RFC 3031, January 2001.
- [5] P. Mohapatra and E. Rosen, “BGP Encapsulation SAFI and BGP Tunnel Encapsulation Attribute,” RFC 5512, April 2009.
- [6] P. Francis, X. Xu, H. Ballani, D. Jen, R. Raszuk, and L. Zhang, “FIB Suppression with Virtual Aggregation,” October 2009, [draft-ietf-grow-va-01](#).
- [7] IETF Global Routing Operations Working Group (GROW), <http://www.ietf.org/dyn/wg/charter/grow-charter.html>

PAUL FRANCIS is a faculty member at the Max Planck Institute for Software Systems in Germany. He has been active periodically in the IETF for nearly 20 years. Dr. Francis has held research positions at Cornell University, ACIRI, NTT Labs, and Bellcore, and was Chief Scientist at Fast Forward Networks and Tahoe Networks. E-mail: francis@mpi-sws.org

XIAOHU XU graduated from Beijing University of Posts and Telecoms in 2000. He has been working in the telecom industry for about 10 years and now is a research engineer with IP Advanced Technology Research Department of Huawei Technologies. Before joining Huawei at the end of 2004, he was the chief engineer of the Technical Support Department for Harbour Networks. E-mail: xuxh@huawei.com

RFC Editor in Transition: Past, Present, and Future

by Leslie Daigle, ISOC

In April 2009, the *Request For Comments* (RFC) Editor published RFC 5540^[1], “40 Years of RFCs,” which summarized the publication history of the RFC Series. The series has been the technical publication series for Internet technology since long before there was an *Internet Engineering Task Force* (IETF). Although the RFC Series is the publication vehicle for the IETF, it has been, and remains, scoped more broadly than that (refer to RFC 4844^[2], “The RFC Series and RFC Editor”). The RFC Series is the archival series dedicated to documenting Internet technical specifications, including general contributions from the Internet research and engineering community as well as standards documents.

For the past three of the four decades of the history of the series, the RFC Editor work has been carried out at the *University of Southern California Information Sciences Institute* (USC/ISI). The RFC Editor role now faces another evolutionary step: The work involved in managing the overall series is being split up to recognize the different components of the editing, production, and archiving activities and to lay the groundwork to ensure its continued success, as outlined in RFC 5620^[3], “RFC Editor Model (Version 1).”

At the IETF 76 plenary in Hiroshima, Japan, in November 2009, USC/ISI and the role it has played in supporting the RFC Editor over the past 30 years were given special recognition. Some members of the team will move from USC/ISI to the RFC Editor’s new home, where they will continue their work. We took the opportunity to talk with current and future RFC Editor staff and advisory board members, including current RFC Editor staff members Bob Braden, Sandy Ginoza, and Alice Hagens, as well as Bob Hinden, who is a member of the RFC Editor advisory board.

The People Behind the RFC Editor

Jon Postel was the first RFC Editor, starting the position in 1969 as an activity to keep track of RFC Series documents. Bob Braden, who was then part of the *Advanced Research Project Agency Network* (ARPANET) research program, told how he got started with the RFC Series: “I wrote my first RFC in the early 1970s, when it was somewhere around RFC 100. I was at that point manager of programming for the Computing Center at the *University of California, Los Angeles* (UCLA), and *Advanced Research Projects Agency* (ARPA) wanted to connect it to ARPANET as a resource.” This was all pre-TCP/IP, and Bob’s staff had to implement file transfer and Telnet. At the same time, Jon was a graduate student at UCLA, and Bob worked with him as a colleague. It was before Jon got his Ph.D. and moved to SRI in 1973–1974. In 1980, Jon moved to USC/ISI, taking the RFC editorship with him. Joyce Reynolds went to work for Jon at USC/ISI. She did much of the actual editing and became an important part of making the RFC Editor activity viable.

Jon was responsible for quality control, running the operation, and generally being the series editor. When Jon died suddenly in 1998, Bob, who joined USC/ISI in 1986, and Joyce both felt a keen sense of loss. “Jon was a very remarkable guy in many ways,” Bob said. “We knew how much the RFC Series meant to Jon, and we volunteered to carry it on.”

Sandy Ginoza joined USC/ISI to work on the RFC Editor activity in 1999, just after Jon passed away. Alice Hagens came onboard in 2005, taking on more of the computer-oriented aspects of the work.

RFC Series

Although we tend to reference and read individual RFC documents, it is important to understand that there is significant value in the collection of published RFCs as a *series*. On the importance of the RFC Series, Bob Hinden said, “This community is IETF-focused, but to the larger world not centered around the IETF; it’s really the RFCs that are how you build the Internet. One of the things that made the Internet possible was the RFC Series: that you could build things and deploy things without coming to IETF meetings was valuable.” Bob went on to outline his own experiences, such as meeting engineers in Taipei, for whom it was the first time they had ever met anyone who had written an RFC. Even the notion of going to an IETF meeting was in another dimension. “The RFC Series is what enables people to build products, networks, and the Internet,” he said.

And it is quite an active series. Currently, some 300 documents (10,000 pages) are published every year, and although it might be interesting to review the material to detect trends or arcs of work in the Internet technical community, that type of activity is beyond the current scope of the RFC Editor. Focusing on consistency of the series, Bob Braden wondered, “Will we eventually have good enough statistics from the errata system to gauge our error rate?”

The intent of the RFC Series is to serve the broader Internet community; it is not just for or by the IETF. Sandy’s perspective on the value of the *Independent Stream* of RFCs is that “it offers an alternate view than what happens in the IETF and what working groups have decided to take on as part of their chartered activities. It’s good to document that work was done, results were generated, lessons learned, etc. ‘We tried it; don’t do it this way.’ We often get asked why it’s called RFC when we’re not really requesting comments anymore, but that is the genesis, and the Independent Stream keeps some of that alive.”

Bob Braden offered his own perspective on the Independent Stream.

“Historically, the RFC Series is supposed to be larger than the IETF, and while Jon was alive, the editor did whatever he thought he ought to do; the community didn’t question it much.”

However, in the absence of Jon as an authority figure, the community began to ask questions and build its own set of beliefs, eventually coming to believe that RFCs were only for the IETF. That matter was resolved with RFC 4846^[4], which explained that there is a separate set of independent submissions that do not come through the IETF.

“It’s not a big stream, not a lot of documents, but it is important philosophically,” Bob added. “The Internet community is bigger than the IETF.”

The RFC Series is, nevertheless, entwined with the IETF and its activities. For instance, the discussion of (IETF) *Intellectual Property Rights* (IPR) has led to an impasse in assigning boilerplate to RFCs that allow the continued publication of the Independent Stream documents. That subject is being worked on and resolved, but it offers an example of some of the complexities—and frustrations—that can arise as part of the RFC Editor process. “The current situation—that the independent submissions cannot be published because we don’t know what the boilerplate is—is just terrible,” said Bob Braden.

Bob Hinden, who has been tracking the IPR work from the IETF side, agreed and elaborated on some important lessons learned: “The IETF created a process in the IPR working group that focused on trying to provide a solution to what they perceived as a problem. But they lost sight of the complexity and cost of implementing that solution compared with the actual risk of something bad happening. We have learned a lot about doing this in the future. This isn’t like a protocol spec where you fix a bug in the finite state machine. This has a real effect on people doing stuff. When you ask for legal opinions you get the answer about how to solve the problem, but that’s not the end of the process. You need to balance the cost of solving the problem with the risk of what you’re trying to avoid. Lawyers are supposed to give you the lowest-risk answer. You need to follow through with questions about likelihood and consequences. This is all great hindsight, and I hope we can apply it in the future.” Hinden also said he believes the current impasse could have been avoided if the new procedure had specified that it go into effect when appropriate supporting conditions were met, instead of on a specific flag day, such as the date of publication of the RFC.

The effects of entwining the RFC Series and the IETF go both ways. For example, the RFC Series recognizes three levels of standards documents: *Proposed*, *Draft*, and *Full*. The expectation, documented in the IETF standards process, is that standards-track specifications should be published as Proposed and then advanced to Draft and Full as the specification gets tested commercially and acknowledged as appropriately mature to move to the next stage.

In reality, as observed at the IETF 76 plenary, many of the important specifications that form the basis of the operating Internet are still published only as Proposed Standard. Bob Braden explained the history of the standards-track RFC maturity system this way: “Labels were invented whole cloth by the original *Internet Architecture Board* (IAB), who were a bunch of academics. At that point the Internet had not been commercialized—there were no commercial pressures—so we imagined that it made sense to step through progressions in a theoretical world. In the real world, companies are putting out products. There is no financial incentive for people to spend time advancing documents. Plus, the IETF is so large and there are so many working groups that we try to dispatch them as fast as we can; there is no one around to advance a document.” There have been, and will continue to be, proposals for moving important, current standards (such as the *Border Gateway Protocol* [BGP]) forward in maturity or for collapsing the maturity scale and labeling system.

On the fun side of the RFC Series, there remains a tradition of “April 1st” RFCs. “That people want to participate in that is cool,” said Sandy. “And we get to see the runners-up and the really-not-so-good ideas!”

Alice agreed, adding that “there are high standards for straight-faced satire.”

RFC Editor

Traditionally, the RFC Editor has not only populated the series with new (approved) documents but also kept all the threads together in the RFC Series. Describing the origins of the role, Bob Braden pointed out that “Originally, Jon was prince of his kingdom. As RFC Editor, he was an honorary member of the IAB informally called the *Protocol Czar*. He used the RFC Editor position to actively prevent bad ideas from getting pushed. Jon imposed a consistency of style on the document series. You pick up RFC 1001 and compare it with 2001, and they look very similar.” Jon believed, and the RFC Editor continues to believe today, that consistency was a worthwhile attribute, promoting stability in the series.

Reflecting back, Bob Braden said, “In discussions over the last five years, people have expressed the view that we don’t need an RFC Editor—just take an Internet Draft and publish it. That notion drives me crazy. The implication is that it doesn’t matter whether it is good English, correctly referenced, consistent, etc. I can’t stand that view.” One of the arguments for such an approach to IETF document publishing is that editing can inadvertently alter, and thereby introduce errors to, text. But the RFC Editors understand that.

Alice said changes to text can be problematic, “partly because of the technical content and partly because it is a group process. It’s agreed-upon text. The idea is how precious the text is and how a slight change can make a large difference.”

Sandy agreed, adding that “for as many changes that get pushed back upon, there are many that make it through the process: for as many people as look at the document before it gets to us, there are things that escape them; there is often missing text, missing words.” According to Alice, with working group documents, people often focus on getting the technical ideas right, but nobody has read the text from beginning to end. In addition, many in the community are not native English speakers. It all comes back to the consistency and professionalism of the output of the series.

RFC Editing Process

As the RFC Series has grown, achieving consistency has required the creation and refining of processes. “When Joyce and I took over,” said Bob Braden, “we built the website and regularized a lot of things, and the community began to ask, ‘Why do you do it that way?’” In response, the editors started publicizing the *Style Manuals* they used. Joyce and Bob generated a lot of rules that have become institutionalized.

Of course, there is continuing evolution. Bob Braden noted that the addition of errata was his idea, although “it has turned out to be a much, much bigger deal than ever imagined, as is often the case,” he said, laughing. “Now we’re talking about adding image files to solve the problem of incorporating graphics in an ASCII RFC. John Klensin and I generated a plausible solution for that, and we hope to get it installed soon.”

It is important to note that there are some edits the RFC Editor will not make. According to Sandy, the RFC Editor tries to ensure consistency of terminology and to make recommendations that improve consistency within a document, both in a technical sense and within the series. “We don’t change the active/passive voice,” she said.

“We might suggest it, but we are concerned that it would affect the author’s intent.” Being conservative is critical. Sandy said she was surprised by how “simple grammatical changes can have a serious technical effect; placement of a comma can make a big difference in how people read the document and what they implement.”

Working with authors is an important part of making the editing process successful. Innovations such as having the *RFC Editor Help Desk* at IETF meetings and making the AUTH48^[5] (final check of the RFC Editor’s edits) more of an interactive dialogue have helped build community and create awareness of how to build a better document that conveys the meaning as intended. “It is extremely useful to get discussions started earlier, which lessens problems during AUTH48,” said Alice. She added that it has also been useful to have face time with the developers of community-created tools, such as *xml2rfc*^[6] and the *Augmented Backus–Naur Form* (ABNF) checker, which have been instrumental in improving RFC production. Office hours, building relationships, and face time “all help make it about working together,” said Sandy.

Looking forward, Sandy said she would like to see the RFC editing process (and series) “grow and continue to be more consistent, with better community relations and more transparency so authors can look at our site and better understand the process, instead of thinking their document has gone into a black box.”

On the Verge of Major Change

As this article is written, the RFC world is on the brink of major structural change. Following IAB-led community discussion, there is a new model for recognizing the components of activity that make up the RFC activities. ISI is handing off the RFC Editor activity, which will be taken up by separate organizations working together. In February 2010, the IAB appointed Nevil Brownlee as the *Independent Submissions Editor* (ISE) and Glenn Kowack as the *Transitional RFC Series Editor* (RSE). In October 2009, *Association Management Solutions* (AMS) was awarded 2-year contracts to manage the RFC Production Center and the RFC Publisher.

Sandy will be joining AMS as RFC Production Center director and Alice will be joining as senior editor and information technology development project manager. To the question of whether the current RFC advisory board will carry forward in the current format or will change, Bob Braden answered, “The current board serves two functions: It provides a supply of experienced people who review independent submissions, but it also gives the RFC Editor advice on policy matters. Some members of the advisory board are very strong members of the IETF in terms of policy advice. In forming the board, I tended to identify a subset of people within the IETF who have long IETF and publishing experience. In the new world there will be an *RFC Series Advisory Group* (RSAG), which will take over the policy discussions that are currently being conducted by the editorial board. In practice it will be the same people, at least for a while, but with separate duties. That separation is useful.”

In considering the change of organizations, Sandy said the biggest thing in moving to AMS is that it is a more service-oriented environment. “In the new model,” she said, “it is important that the ISE and RSE be respected individuals who are granted some of the independence the RFC Editor had at ISI.”

Alice added that the institutional memory of the RFC Editor function will not be lost with the move to AMS. “Sandy has worked side by side with Bob Braden for 10 years, and much of the process is written down in the document series. I’m confident that the continuity of the series won’t be lost by the move to AMS.”

Bob Hinden offered another perspective. “I think one of the positive things that has come out of the new model that has gotten lost is this: A lot of people in the IETF didn’t understand where the series had come from, or why the IETF chose to use it,” he said.

“It is the formalization that there are different streams that have different rules. Before, this was confused with the IETF standards process. Going forward we’ll have the opportunity to use the RFC Series for other relevant Internet publication streams that have not been part of IETF. Now we have a framework that would allow that.”

Although it is on the verge of major changes, the RFC Series and RFC Editor functions are clearly continuing what has been a long process of constant evolution and change. This transition is just a new chapter in the history of the series.

[Ed.: This article is composed of interviews conducted by Leslie Daigle and Lucy Lynch, and notes compiled by Mat Ford. The original version was published in *The IETF Journal*, Volume 5, Issue 3, January 2010 and has been updated for use in IPJ. *The IETF Journal* can be obtained from <http://isoc.org/ietfjournal/>]

For Further Reading

- [1] RFC Editor, “40 Years of RFCs,” RFC 5540, April 2009.
- [2] L. Daigle, Ed., Internet Architecture Board, “The RFC Series and RFC Editor,” RFC 4844, July 2007.
- [3] O. Kolkman, Ed., IAB, “RFC Editor Model (Version 1),” RFC 5620, August 2009.
- [4] J. Klensin and D. Thaler, Eds., “Independent Submissions to the RFC Editor,” RFC 4846, July 2007.
- [5] <http://www.rfc-editor.org/pubprocess.html>
- [6] Marshall T. Rose and Carl Malamud, “Writing Internet Drafts and RFCs Using XML,” *The Internet Protocol Journal*, Volume 10, No. 1, March 2007.



Alice Hagens, Bob Braden, and Sandy Ginoza are recognized at IETF 76 for their work with the RFC Editor. (Photo: Internet Society)

IETF Outcomes Wiki Launched

As an organization, the *Internet Engineering Task Force* (IETF) measures its success by its publication of RFCs (see previous article). It does not explicitly ask itself whether published work is adopted and used by the greater Internet community. The IETF's dialogue about success started to change with the production of RFC 5218, "What Makes for a Successful Protocol?"^[1] which documented case studies and empirical data about some of the factors that appear to correlate with success, in terms of community uptake for IETF work.

Taking a different approach in assessing long-term IETF impact, another tool is now available: A wiki that lets community participants list the success or failure of significant standards. The *Outcomes Wiki*^[2] divides listings according to the "areas" used for managing technical work in the IETF, such as Applications or Transport. Outcomes are rated according to a 6-point scale, ranging from "complete failure" to "massive adoption, plus extensive derivative work."

The wiki began in June 2009, as an independent effort among a small set of IETF participants, to test its feasibility and evolve its design. For example, it quickly became clear that the single attribute of success vs. failure needed to be qualified by another attribute that indicates who the work is intended for, called "Target Segment." Work that is intended to support the internal operations of an *Internet Service Provider* (ISP) is not necessarily visible to the billions of Internet users and will, at best, be part of only a few thousand organizations. In terms of Internet scale, that is considered minuscule. However wide adoption of a tool among ISPs can have substantial benefit, and thereby qualify as "massive adoption."

The wiki can serve both as a means of recording the IETF's track record of successes and failures, as well as providing a means of encouraging community dialogue about the quality of different IETF efforts. In addition, it can provide a window onto completed IETF work for the broader Internet community.

[1] D. Thaler and B. Aboba, "What Makes for a Successful Protocol?" RFC 5218, July 2008.

[2] <http://trac.tools.ietf.org/misc/outcomes/>

Final Phase of Four-byte AS Number Policy Begins in APNIC Region

From 1 January 2010, the *Asia Pacific Network Information Centre* (APNIC) ceased to make a distinction between four-byte only and two-byte only *Autonomous System* (AS) numbers. Instead, all AS numbers are now considered to be four-byte AS numbers.

This change marks the third phase of the transition to four-byte AS numbers. For more information on the implementation phases of the four-byte AS number policy, please see “Policies for Autonomous System number management in the Asia Pacific region,” section 6.3, “Timetable for moving from two-byte only AS numbers to four-byte AS numbers,” available from:

<http://www.apnic.net/policy/asn-policy.html#6.3>

To learn more about how the transition to four-byte AS numbers may affect your network, see: <http://icons.apnic.net/asn>

Charting the Course for Future Internet Leaders

As the importance of the Internet grows in all aspects of modern life, so too do the challenges of those in positions of leadership and responsibility.

Responding to the need for well-qualified leadership, the *Internet Society* (ISOC) is now accepting applications from people seeking to join the new generation of Internet leaders to address the critical technology, policy, business, and education challenges that lie ahead.

Successful candidates in ISOC’s *Next Generation Leaders Program* will gain a wide range of skills in a variety of disciplines, as well as the ability and experience to work with people at all levels of society.

This program, under the patronage of the European Commission, blends course work and practical experience to help prepare young professionals (aged from 20 to 40) from around the world to become the next generation of Internet technology, policy, and business leaders.

“The Internet Society’s Next Generation Leaders Program is a unique opportunity to identify potential Internet leaders and help them accelerate their careers,” said Bill Graham, responsible for strategic global engagement at ISOC.

The key to the Internet’s success lies in the Internet Model of decentralized architecture and distributed responsibility for development, operation, and management. That model also creates important leadership opportunities, especially in those spaces where technology, policy, and business intersect.

“We have designed the Next Generation Leaders Program to prepare young professionals for leadership, bridging the boundaries between business, technical development, policy, and governance on local, regional, and international levels,” said Graham.

Full details of the Next Generation Leaders Program are available at: <http://www.isoc.org/leaders/>

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ contains standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2010 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

June 2010

Volume 13, Number 2

A Quarterly Technical Publication for
Internet and Intranet Professionals

In This Issue

From the Editor	1
Address Sharing	2
Implementing DNSSEC	16
Book Review.....	27
Fragments	30
Call for Papers.....	35

FROM THE EDITOR

Protocol changes are never easy, especially when they involve something as fundamental as the *Internet Protocol* (IP). This journal has published numerous articles about the depletion of IPv4 addresses and several articles about IPv6, including methods for a gradual transition from v4 to v6. A lot of energy has gone into the development, promotion, and deployment of IPv6, but in reality only a small fraction of the global Internet currently supports IPv6. Meanwhile, the *Internet Assigned Numbers Authority* (IANA) and the *Regional Internet Registries* (RIRs) will “soon” (12 to 24 months from now is predicted) run out of IPv4 addresses to allocate. Although this situation has some serious implications for new entrants to the *Internet Service Provider* (ISP) market, it does not spell the end of the Internet as we know it. Numerous *Network Address Translation* (NAT) solutions are already widely deployed, and the IETF is discussing other solutions. One example is *Address Sharing* as explained by Geoff Huston in our first article.

Changes to the *Domain Name System* (DNS) are also underway. The *Domain Name System Security Extensions* (DNSSEC) are being gradually deployed in the global Internet. As with any complex technology, implementation of DNSSEC is not without problems. Our second article, by Torbjörn Eklöv and Stephan Lagerholm, is a step-by-step guide for those considering implementing DNSSEC in their network.

By now you will be aware that we have implemented a renewal system for subscribers and will not be automatically extending your subscription unless you contact us via e-mail or use the online tool to renew your subscription. You can find your subscription ID and expiration date either on the back page of your copy or on the envelope that it came in. In order to access your record, click the “Subscriber Services” link on our webpage at www.cisco.com/ipj, and enter your e-mail address and the subscription ID. The system will send you a link that allows direct access to your record, and you will be able to update your address and renew your subscription. If you no longer have access to the e-mail you used when you subscribed, or have forgotten your subscription ID, just send a message to ipj@cisco.com and we will make the necessary changes for you.

—Ole J. Jacobsen, Editor and Publisher

ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

NAT++: Address Sharing in IPv4

by Geoff Huston, APNIC

In this article I examine the topic that was discussed in a session at the 74th meeting of the *Internet Engineering Task Force* (IETF) in March 2009, about *Address Sharing* (the SHARA BOF)^[0], and look at the evolution of *Network Address Translation* (NAT) architectures in the face of the forthcoming depletion of the unallocated IPv4 address pool.

Within the next couple of years we will run out of the current supply of IPv4 addresses. As of the time of writing this article, the projected date when the *Internet Assigned Numbers Authority* (IANA) pool will be depleted is August 3, 2011, and the first *Regional Internet Registry* (RIR) will deplete its address pool about March 20, 2012.

Irrespective of the precise date of depletion, the current prediction is that the consumption rate of addresses at the time when the free pool of addresses is exhausted will probably be running at some 220 million addresses per year, indicating a deployment rate of some 170–200 million new services per year using IPv4. The implication is that the Internet will exhaust its address pool while operating its growth engines at full speed.

How quickly will IPv6 come to the rescue? Even the most optimistic forecast of IPv6 uptake for the global Internet is measured in years rather than months following exhaustion, and the more pessimistic forecasts extend into multiple decades.

For one such analysis using mathematical modelling techniques, refer to Jean Camp's work^[1]. One of the conclusions from that 2008 study follows: "There is no feasible path which results in less than years of IPv4/IPv6 co-existence. Decades is not unreasonable."

The implication of this conclusion is that we will need to operate a dual-stack Internet for many years to come, and the associated implication is that we will have to make the existing IPv4 Internet span a billion or more new deployed services—and do so with no additional address space.

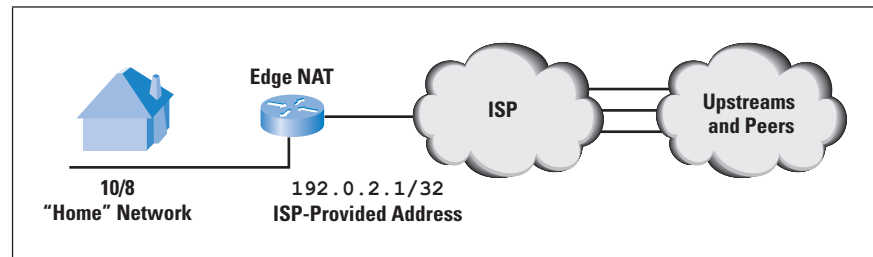
So how are we going to make the IPv4 address pool stretch across an ever-larger Internet?

Given that the tool chest we have today is the only one available, there appears to be only one answer to this question: Use *Network Address Translators*, or NATs.

For a description of how NATs work and some of the terminology used to describe NAT behavior, refer to the article "Anatomy: A Look Inside Network Address Translators," published in this journal^[2].

Today NATs are predominately edge devices that are bundled with DSL modems for residential access, or bundled with routing and security firewall equipment for small to midsize enterprise use as an edge device. The generic model of NAT deployment currently is a small-scale edge device that generally has a single external-side public IP address and an internal-side private IP network address (often network 10). The NAT performs address and port translation to map all currently active sessions from the internal addresses to ports on the public IP address. This NAT deployment assumes that each edge customer has the unique use of a public IP address (refer to Figure 1).

Figure 1: Conventional NAT Deployment



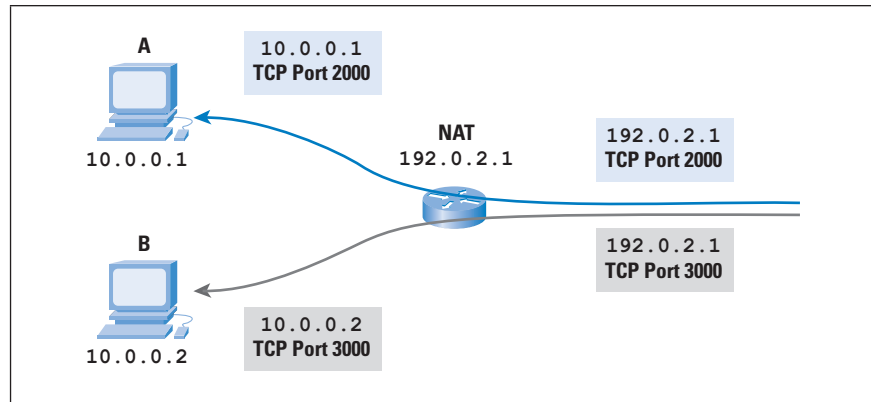
The question provoked by IPv4 address exhaustion is what happens when there are no longer sufficient IPv4 addresses to provide this 1:1 mapping between customers and public IPv4 addresses? In other words, what happens when there are simply not enough IPv4 addresses to allow all customers to have exclusive use of their own unique IPv4 address?

This question has only two possible answers. One is for no one to use IPv4 addresses at all, on the basis that the entire Internet has migrated to use IPv6. But this answer appears to be an uncomfortable number of decades away, so we need to examine the other answer: If there are not enough addresses to go around, then we will have to *share* them.

But isn't sharing IP addresses impossible in the Internet architecture? The IP address in a packet header determines the destination of the packet. If two or more endpoints share the same address, then how will the network figure out which packets go to which endpoint? It is here that NATs and the transport layer protocols, the *Transmission Control Protocol* (TCP) and the *User Datagram Protocol* (UDP), come together. The approach is to use the *port address* in the TCP and UDP header as the distinguishing element.

For example, in Figure 2, incoming TCP packets with TCP port address 2000 may need to be directed to endpoint A, while incoming TCP packets with TCP port address 3000 need to be directed to endpoint B. The incoming TCP packets with a port address of 2000 are translated to have the private IP address of endpoint A, and incoming TCP packets with a port address of 3000 are translated to have the private address of endpoint B.

Figure 2: Address Sharing with NATs



As long as you restrict yourself to applications that use TCP or UDP, you don't rely on receiving *Internet Control Message Protocol* (ICMP) packets, and you don't use applications that contain IP addresses in their payload, then you might expect this arrangement to function.

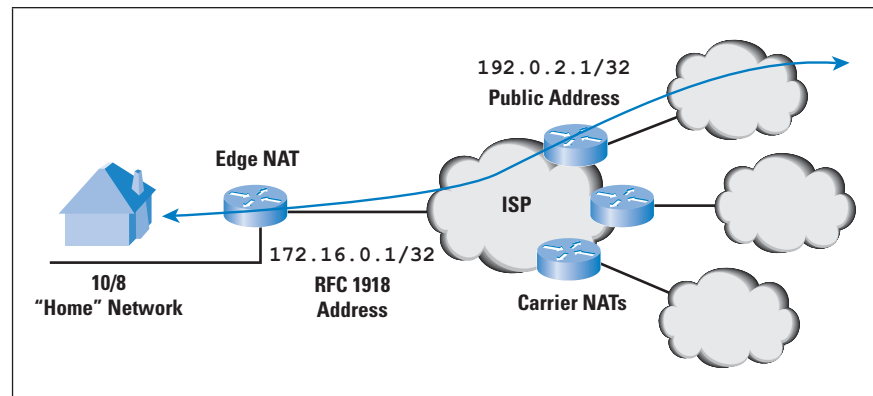
ICMP is a problem because the ICMP packet does not contain a TCP or UDP transport layer. All that a NAT sees in the ICMP packet is its own external address as the destination IP address. To successfully deliver an ICMP packet through a NAT, the NAT needs to perform a more complex function that uses the ICMP-encapsulated IP header to select the original outbound combined IP + TCP header or IP + UDP header in the ICMP payload. The source IP address and transport protocol port address in the ICMP payload are then used to perform a lookup into the NAT binding table and then perform two mappings: one on the ICMP header to map the destination IP address to the internal IP address, and the second on the payload header where the source IP address and port number are changed to the interior-side values, and the checksums altered as appropriate. Now in most cases ICMP really is not critical, and a conservative NAT implementation may elect to avoid all that packet inspection and simply discard all incoming ICMP messages, but one message that is important is the ICMP *packet-too-large-and-fragmentation-disabled* message used in IPv4 *Path MTU Discovery*^[3].

Sharing IP addresses is fine in theory, but how can we achieve it in practice? How can many customers, already using NATs, share a single public IP address?

Carrier-Grade NATs

One possible response is to add a further NAT into the path. In theory the *Internet Service Provider* (ISP) could add NATs on all upstream and peer connections, and perform an additional NAT operation as traffic enters and leaves the ISP's network. Variations of this approach are possible, placing the ISP NATs at customer aggregation points within the ISP's network, but the principle of operation of the ISP NAT is much the same.

Figure 3: Carrier NATs



The edge NATs translate between private address pools at each customer's site and an external address provided by the ISP, so nothing has changed there. The change in this model is that the ISP places a further NAT in the path within the ISP network, so that a set of customers is then sitting behind a larger NAT inside the ISP's network, as shown in Figure 3.

This scenario implies that the external address that the ISP provides to the customer is actually yet another private address, and the ISP's NAT performs yet another transform to a public address in this second NAT. In theory this NAT is just a larger version of an existing NAT with larger NAT binding space, higher packet-processing throughputs, and a comprehensive specification of NAT binding behavior. In practice it may be a little more complicated because at the network edge the packet rates are well within the processing capability of commodity processors, whereas in the core of the network there is an expectation of higher levels of robust performance from such units. Because it is intended that such a NAT handle thousands of customers and large numbers of simultaneous data flows and peak packet rates, it requires a performance level well beyond what is seen at the customer edge and, accordingly, such a NAT has been termed a *Carrier-Grade NAT* (CGN), or a *Large-Scale NAT* (LSN).

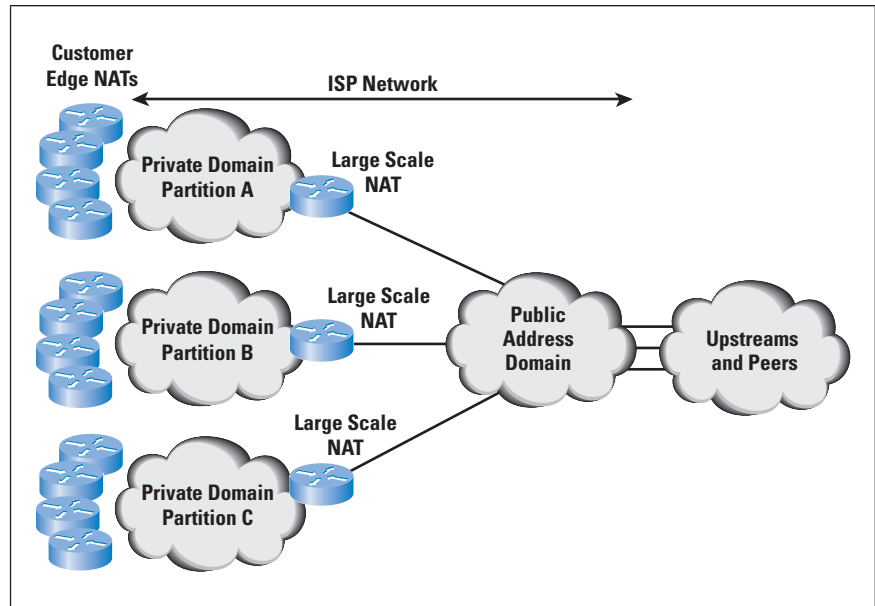
From the inside of the two NATs, not much has changed with the addition of the CGN in terms of application behavior. It still requires an outbound packet to trigger a binding that allows a return packet through to the internal destination, so nothing has changed there. Other aspects of NAT behavior, notably the NAT binding lifetime and the form of *Cone Behavior* for UDP, take on the more restrictive of the two NATs in sequence. The binding times are potentially problematic in that the two NATs are not synchronized in terms of binding behavior. If the CGN has a shorter binding time, it is possible for the CGN to misdirect packets and cause application-level problems. However, this situation is not overly different from a single-level NAT environment where aggressively short NAT binding times also run the risk of causing application-level problems when the NAT drops the binding for an active session that has been quiet for an extended period of time.

However, one major assumption is broken in this structure, namely that an IP address is associated with a single customer. In this model a single public IP address may be used simultaneously by many customers at once, albeit on different port numbers. This scenario has obvious implications in terms of some current practices in filters, firewalls, “black” and “white” lists, and some forms of application-level security and credentials where the application makes an inference about the identity and associate level of trust in the remote party based on the remote party’s IP address.

This approach is not without its potential operational problems as well. For the ISP, service resiliency becomes a critical concern in so far as moving traffic from one NAT-connected external service to another will cause all the current sessions to be dropped, unless the internal ISP network architecture uses a transit access network between the CGNs and the external transit providers. Another concern is one of resource management in the face of potentially hostile applications. For example, an end host infected with a virus may generate a large amount of probe packets to a large range of addresses. In the case of a single edge NAT, the large volumes of bindings generated by this behavior become a local resource management problem because the customer’s network is the only affected site. In the case where a CGN is deployed, the same behavior starts to consume binding space on the CGN and, potentially, can starve the CGN of external address bindings. If this problem is seen to be significant, the CGN would need to have some form of external address rationing per internal client in order to ensure that the entire external address pool is not consumed by a single errant customer application. This “rationing” would have the unwanted effect of forcing the ISP to deny access to its customers.

The other concern here is one of scalability. Although the greatest leverage of the CGN in terms of efficiency of usage of external addresses occurs when the greatest numbers of internal edge-NAT-translated clients are connected, there are some real limitations in terms of NAT performance and address availability when an ISP wants to apply this approach to networks where the customer population is in the millions or larger. In this case the ISP is required to use an IPv4 private address pool to number every client. But if all customers use network 10 as their “internal” network, then what address pool can the ISP use for its private address space? One of the few answers that come to mind is to deliberately partition the network into numerous discrete networks, each of which can be privately numbered from the smaller private address pool of **172.16.0.0/12**, allowing for some 600,000 or so customers per network partition, and then use a transit network to “glue” together the partitioned elements, as shown in Figure 4.

Figure 4: Multiple Carrier NAT
Deployment Using Network
Partitioning



The advantage of the CGN approach is that for the customer nothing changes. Customers do not need to upgrade their NAT equipment or change them in any way, and for many service providers this motivation is probably sufficient to choose this path. The disadvantages of this approach lie in the scaling properties when looking at very large deployments, and the problems of application-level translation, where the NAT attempts to be “helpful” by performing deep packet inspection and rewriting what it thinks are IP addresses found in packet payloads. Having one NAT do this rewriting is bad enough, but loading them up in sequence is a recipe for trouble!

Are there alternatives?

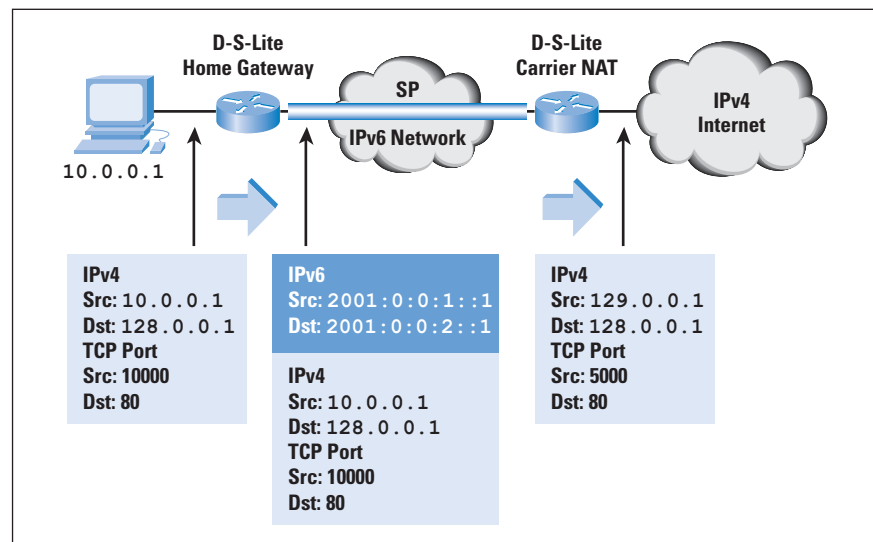
Dual-Stack Lite and Carrier-Grade NATs

One rather elegant alternative is described by Alain Durand and others in an Internet Draft “Dual-stack lite broadband deployments post IPv4 exhaustion”^[4]. The assumption behind this approach is that the ISP’s network infrastructure needs to support IPv6 running in native mode in any case, so is there a way in which the ISP can continue to support IPv4 customers without running IPv4 internally?

Here the customer NAT is effectively replaced by a tunnel ingress/egress function in the *Dual-Stack Lite Home Gateway*. Outgoing IPv4 packets are not translated, but are encapsulated in an IPv6 packet header, where the IPv6 packet header contains a source address of the carrier side of the home gateway unit and a destination address of the ISP’s gateway unit. From the ISP’s perspective, each customer is no longer uniquely addressed with an IPv4 address, but instead is addressed with a unique IPv6 address. The customer’s interface to the ISP network, the Home Gateway, is configured with this IPv6 address as the customer end of the IPv4-in-IPv6 tunnel, where the other end of the tunnel is the IPv6 address of the ISP’s Dual-Stack Lite Gateway unit.

The service provider's Dual-Stack Lite gateway unit performs the IPv6 tunnel termination and a NAT translation using an extended local binding table. The “interior” NAT address is now a 4-tuple of the IPv4 source address, protocol ID, and port, plus the IPv6 address of the home gateway unit, while the external address remains the triplet of the public IPv4 address, protocol ID, and port. In this way the NAT binding table contains a mapping between interior “addresses” that consist of IPv4 address and port plus a tunnel identifier and public IPv4 exterior addresses. This way the NAT can handle a multitude of network 10 addresses, because the addresses can be distinguished by different tunnel identifiers. The resultant output packet following the stripping of the IPv6 encapsulation and the application of the NAT function is an IPv4 packet with public source and destination addresses. Incoming IPv4 packets are similarly transformed, where the IPv4 packet header is used to perform a lookup in the Dual-Stack Lite gateway unit, and the resultant 4-tuple is used to create the NAT-translated IPv4 packet header plus the destination address of the IPv6 encapsulation header (refer to Figure 5).

Figure 5: Dual-Stack Lite



The advantage of this approach is that now only a single NAT is needed in the end-to-end path because the functions of the customer NAT are now subsumed by the carrier NAT. This scenario has some advantages in terms of those messy “value-added” NAT functions that attempt to perform deep packet inspection and rewrite IP addresses found in data payloads. There is also no need to provide each customer with a unique IPv4 address, public or private, so the scaling limitations of the dual-NAT approach are also eliminated. The disadvantages of this approach lie in the need to use a different *Customer Premises Equipment* (CPE) device, or at least one that is reprogrammed. The device now requires an external IPv6 interface and at a minimum an IPv4 or IPv6 tunnel gateway function. The device can also include a NAT if desired, but it is not required in terms of the basic Dual-Stack Lite architecture.

This approach pushes the translation into the middle of the network, where the greatest benefit can be derived from port multiplexing, but it also creates a critical hotspot for the service itself. If the carrier NAT fails in any way, the entire customer base is disrupted. It seems somewhat counter intuitive to create a resilient network with stateless switching environments and then place a critical stateful unit in the middle! So is there an approach that can push this translation back to the edges while avoiding a second NAT in the carrier's network?

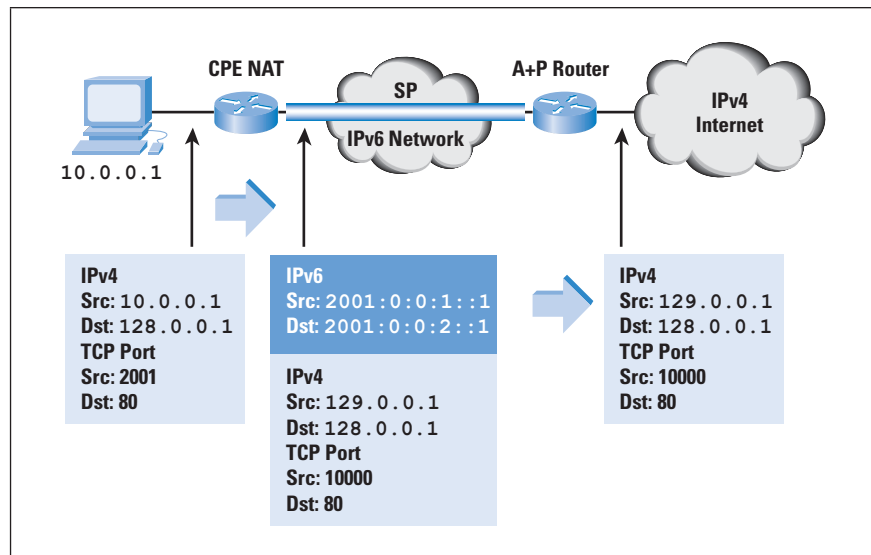
The Address Plus Port Approach

The observation here is that CPE NATs currently map connections into the 16-bit port field of the single external address. If the CPE NAT could be coerced into performing this mapping into 15 bits of the port field, then the external address could be shared between two edge CPE devices, with the leading bit of the port field denoting which CPE device. Obviously, moving the bit marker across the port field would allow more CPE devices to share the one address, but it would reduce the number of available ports for each CPE device in the process.

The theory is again quite simple. The CPE NAT is dynamically configured with an external address, as happens today, and a port range, which is the additional constraint. The CPE NAT performs the same function as before, but it is now limited in terms of the external ports it can use in its NAT bindings to those that lie within the provided port range, because some other CPE may be concurrently using the same external IP address with a different port range.

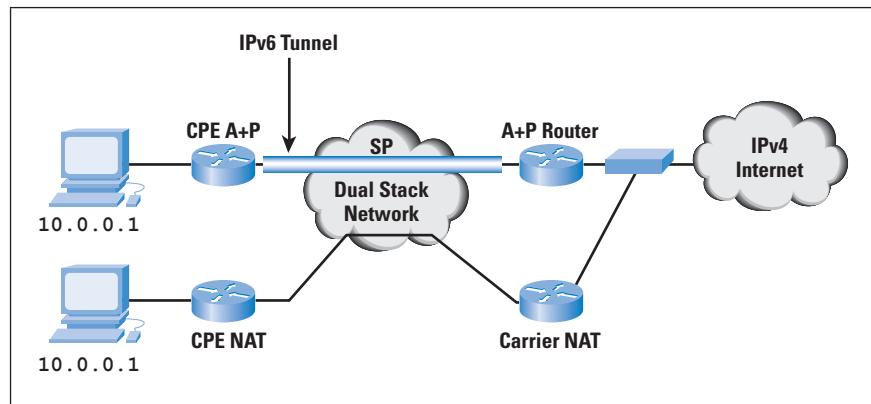
For outgoing packets this limitation implies only a minor change to the network architecture, in that the RADIUS^[9] exchange to configure the CPE now must also provide a port range to the CPE device. However, the case of incoming packets is more challenging. Here the ISP must forward the packet based not only on the destination IP address, but also on the port value in the TCP or UDP header. A convenient way to forward the packet is to take the Dual-Stack Lite approach and use an IPv4-in-IPv6 tunnel between the CPE and the external gateway (Figure 6). This gateway, or *Address Plus Port* (A + P) router, needs to be able to associate each address and port range with the IPv6 address of a CPE device, which it can learn dynamically as it decapsulates outgoing packets. Corresponding incoming packets are encapsulated in IPv6 using the IPv6 destination address that it has learned previously. In this manner the NAT function is performed at the edge, much as it is today, and the interior device is a more conventional form of tunnel server.

Figure 6: Address Plus Port Framework



This approach relies on every CPE device being able to operate using a restricted port range, to perform IPv4-in-IPv6 tunnel ingress/egress functions, and to act as an IPv6 provisioned endpoint for the ISP network, which is perhaps an unrealistic hope. Further modifications to this model (Figure 7) propose the use of an accompanying CGN operated by the ISP to handle those CPE devices that cannot support these Address Plus Port functions.

Figure 7: Combined Address Plus Port and Carrier Grade NAT



If the port range assigned to the CPE is from a contiguous range of port values, then this approach could exacerbate some known problems with infrastructure protocols. There are *Domain Name System* (DNS) problems with guessable responses. The so-called “Kaminsky Attack” on the DNS^[5, 6] is one such example where the attack can be deflected, to some extent, by using a randomly selected port number for each DNS query. Restricting the port range could mitigate the efficacy of such measures under certain conditions.

However, despite such concerns, the approach has some positive aspects. Pushing the NAT function to the edge has some considerable advantage over the approach of moving the NAT to the interior of the network.

The packet rates are lower at the edge, allowing for commodity computing to process the NAT functions across the offered packet load without undue stress. The ability for an end-user's application to request a particular NAT binding behavior by speaking directly with the local NAT using the *Internet Gateway Device Protocol*, as part of the *Universal Plug and Play* (UPnP)^[7] framework, will still function in an environment of edge NATs operating with restricted port ranges. Aside from the initial provisioning process to equip the CPE NAT with a port range, the CPE, and the edge environment is largely the same as in today's CPE NAT model.

That is not to say that this approach is without its negative aspects, and it is unclear as to whether the perceived benefits of a "local" NAT function outweigh the problems associated with this model of address sharing. The concept of port "rationing" is a very suboptimal means of address sharing, given that after a CPE device has been assigned a port range those port addresses are unusable by any other CPE. The prudent ISP would assign to each CPE device a port address pool equal to some estimate of peak demand, so that, for example, each CPE device would be assigned 1,000 ports, allowing a single external IP address to be shared across only 60 such CPE clients. Neither the Carrier-Grade NAT approach nor the Dual-Stack Lite approach attempts this form of rationed allocation, allowing the port address pool to be treated as a common resource, with far higher levels of usage efficiency through dynamic management of the port pool.

The difference here is that in the dynamically managed approach any client can use the currently unused port addresses, whereas in the rationed approach each client has access to a fixed pool of port addresses that cannot be shared with any other client—even when the client does not need them. The difference here parallels the difference in network efficiency between time-division multiplexed synchronous circuits and asynchronous packets at Layer 2 in the network model. In the Address Plus Port framework the leverage obtained in terms of making efficient use of coopting these additional 16 bits of port address into the role of additional bits of client identifier address space is reduced by the imposition of a fixed boundary between customer and ISP use in the port address plan. The central NAT model of a CGN effectively pools the port address range and facilitates far more efficient sharing of this common port address pool across a larger client base.

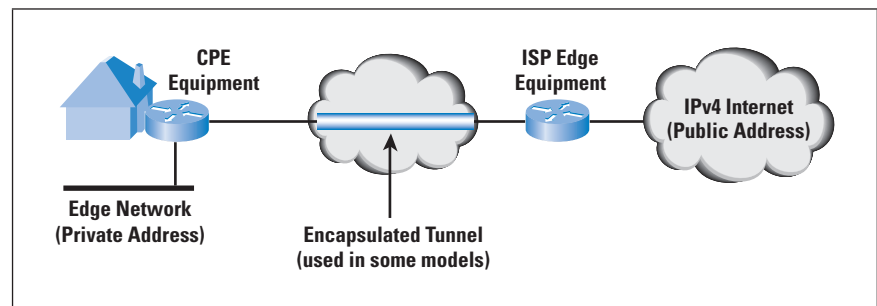
Alain Durand reported to IETF 74 on a data-collection experiment using a *Cable Modem Termination System* (CMTS) with 8,000 subscribers where the peak port consumption level was 40,000 ports, or a maximum average port consumption of 5 ports per subscriber in each direction. As Alain noted, this average value needs to be compared with the hundreds of ports consumed by a single client browsing a Web 2.0 or *Asynchronous Java and XML* (AJAX) site, but he also noted that a central model of port sharing does yield far higher levels of address-sharing efficiency than the Address Plus Port advanced allocation model.^[8]

The other consideration here is that this approach constitutes a higher overhead for the ISP, in that the ISP must support both “conventional” CPE and Address Plus Port equipment. In other words, the ISP must deploy a CGN and support customer CPE using a two-level NAT environment in addition to operating the Address Plus Port infrastructure. Unless customers would be willing to pay a significant price premium for such an Address Plus Port service, it is unlikely that this option would be attractive for the ISP as an additional cost after the CGN cost.

General Considerations with Address Sharing

The basic elements of any such approach to address sharing involve the CPE equipment at the edge, optionally some form of tunneling of traffic between the CPE and the carrier equipment, and carrier-provided equipment at the edge of the carrier’s network (refer to Figure 8).

Figure 8: Generic Architecture for Address Sharing



A variety of technical solutions here involve these basic building blocks, so it is not true to say that this challenge is technically significant. But few ISPs have decided to proceed with large-scale deployment of any form of address-sharing technology for their IPv4 network infrastructure. So what is the problem here?

I suspect that the real concern is the consideration of the relevant business model that would guide this deployment. Today’s Internet is large. It encompasses some 1.7 billion human users, a larger pool of devices, and hundreds of millions of individual points of control. If we want to change this deployed system, we will need copious quantities of money, time, and unity of purpose. So do we have money, time, and unity of purpose?

Money is missing: It could be argued that we have left the entire IPv6 transition effort to this late stage because of a lack of money. The main advantage of the Internet was that it was cheap. Packet sharing is intrinsically more efficient than circuit sharing, and the shift in functions of network service management from the network to the customer-owned and -operated endpoints implied further cost savings for the network operator. So the Internet model gained ascendancy because for consumers it represented a cost-effective choice. It was cheap.

But what does IPv6 offer consumers? For existing Internet consumers it appears that IPv6 does not offer anything that they don't already have with IPv4—it offers mail, the web, various forms of voice services, and games. So consumers are not exactly motivated to pay more for the same services they already enjoy today.

In addition, it would appear that the ISP must carry this cost without incremental revenue from its customer base. But the ISP industry has managed to shave most of its revenue margins in a highly competitive industry, and at the same time lose control of services, their delivery, and their potentially lucrative revenue margins. Thus the ISP industry has been collectively idle in this area not because it cannot see the problem in terms of the imminent exhaustion of IPv4, but because it has little choice because of financial constraints that have prevented it from making the necessary longer-term investments in IPv6. So if the ISP industry has been unwilling to invest in IPv6 so far, then what incentive is there for it to invest in IPv6 and at the same time also invest in these IPv4 address-sharing measures? Is the lure of new, low-margin customers sufficient incentive to make such investments in this carrier-grade equipment? Or is the business case still insufficiently attractive?

Time is missing: The unallocated IPv4 address pool is already visibly waning. Without any form of last-minute rush, the pool will be around for the next 2 years, or until 2012 or so. But with any form of typical last-minute rush, this pool could be depleted in the coming months rather than in the coming years. Can we do what we need to do to get any of these approaches to a state of mass-market deployment in the next few months? All these approaches appear to be at the early stages of a timeline that starts with research and then moves on to development, prototyping, and trials; then to standards activity and industry engagement to orchestrate supply lines for end user equipment, ISP equipment, and definition of operational practices; then to product and service development; and finally, to deployment. For an industry that is the size of the Internet, “technical agility” is now an obsolete historical term. Even with money and unity of purpose this process will take some years, and without money—or even the lure of money—it becomes a far more protracted process, as we have seen already with IPv6 deployment.

And do we have *unity of purpose* here? Do we agree on an approach to address sharing that will allow players to perform their tasks? That will allow consumer product vendors to develop the appropriate product? That will allow application developers to develop applications that will operate successfully in this environment? That will allow the end user platform vendors to incorporate the appropriate functions in the operating system stacks? That will allow ISPs to integrate vendors' productions into their operational environments? Right now it is pretty clear that what we have is a set of ideas, each of which has relative merits and disadvantages, and no real unity of purpose.

It is easy to be pessimistic at this stage, given that the real concerns here appear to be related more to the factors associated with a very large industry attempting to respond to a very challenging change in the environment in which it operates. The question here is not really whether Address Plus Port routing is technically inferior to Dual-Stack Lite, or whether Carrier-Grade NATs are technically better or worse than either of these approaches. The question here is whether this industry as a whole will be able to sustain its momentum and growth across this hiatus. And, from this perspective, I believe that such pessimism about the future of the Internet is unwarranted.

The communications industry has undergone significant technological changes over the years, and this change is one more in the sequence. Some of these transformations have been radical in their effect, such as the introduction of the telephone in the late nineteenth century, whereas others have been more subtle, such as in the introduction of digital technology to telephony in the latter part of twentieth century, replacing the earlier analogue circuit model of telephony carriage. Some changes have been associated with high levels of risk, and we have seen a myriad of smaller, more agile players enter the market to lead the change while the more risk-averse enterprises stand back. On the other hand, other changes require the leverage of economies of scale, and we have seen market consolidation behind a smaller number of highly capitalized players.

My personal opinion is that the Dual-Stack Lite approach is the best one, because it appears to be technically elegant. I suspect, however, that the lowest-common-denominator fall-back position that this somewhat conservative industry will adopt will rely strongly on Carrier-Grade NATs, and the industry is likely to eschew the more complex support mechanisms required by the various permutations of Address Plus Port routing.

Further Reading

- [0] The Address Sharing BOF was held at IETF 74 in March 2009. The presentations and a summary of the session can be found as part of the proceedings of that meeting:
<http://www.ietf.org/proceedings/09mar/shara.html>
- [1] http://www.ripe.net/ripe/meetings/ripe-56/presentations/Camp-IPv6_Economics_Security.pdf
- [2] Geoff Huston, "Anatomy: A Look Inside Network Address Translators," *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [3] Jeff Mogul and Steve Deering, "Path MTU Discovery," RFC 1191, November 1990.
- [4] `draft-ietf-softwire-dual-stack-lite-00.txt`
- [5] http://www.doxpara.com/DMK_BO2K8.ppt
- [6] <http://unixwiz.net/techtips/iguide-kaminsky-dns-vuln.html>

- [7] http://en.wikipedia.org/wiki/Universal_Plug_and_Play
- [8] <http://www.ietf.org/proceedings/09mar/slides/shara-8/shara-8.htm>
- [9] C. Rigney, S. Willens, A. Rubens, W. Simpson, “Remote Authentication Dial In User Service (RADIUS),” RFC 2865, June 2000.
- [10] Egevang, K., and P. Francis, “The IP Network Address Translator (NAT),” RFC 1631, May 1994.
- [11] Srisuresh, P., and D. Gan, “Load Sharing Using IP Network Address Translation (LSNAT),” RFC 2391, August 1998.
- [12] Srisuresh, P., and M. Holdrege, “IP Network Address Translator (NAT) Terminology and Considerations,” RFC 2663, August 1999.
- [13] Tsirtsis, G., and P. Srisuresh, “Network Address Translation—Protocol Translation (NAT-PT),” RFC 2776, February 2000.
- [14] Hain, T., “Architectural Implications of NAT,” RFC 2993, November 2000.
- [15] Srisuresh, P., and K. Egevang, “Traditional IP Network Address Translator (Traditional NAT),” RFC 3022, January 2001.
- [16] Holdrege, M., and P. Srisuresh, “Protocol Complications with the IP Network Address Translator,” RFC 3027, January 2001.
- [17] D. Senie, “Network Address Translator (NAT)-Friendly Application Design Guidelines,” RFC 3235, January 2002.
- [18] Srisuresh, P., J. Kuthan, J. Rosenberg, A. Molitor, and A. Rayhan, “Middlebox Communication Architecture and Framework,” RFC 3303, August 2002.
- [19] Daigle, L., and IAB, “IAB Considerations for Unilateral Self-Address Fixing (UNSAF) Across Network Address Translation,” RFC 3424, November 2002.
- [20] Rosenberg, J., Weinberger, J., Huitema, C., and R. Mahy, “STUN—Simple Traversal of User Datagram Protocol (UDP) Through Network Address Translators (NATs),” RFC 3489, March 2003.

GEOFF HUSTON holds a B.Sc. and a M.Sc. from the Australian National University. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. The author of numerous Internet-related books, he is currently the Chief Scientist at APNIC. He was a member of the Internet Architecture Board (IAB) from 1999 until 2005, and served on the Board of the Internet Society from 1992 until 2001.

E-mail: gih@apnic.net

Operational Challenges When Implementing DNSSEC

by Torbjörn Eklöv, Interlan Gefle AB, and Stephan Lagerholm, Secure64 Software Corp.

As a reader of *The Internet Protocol Journal*, you are probably familiar with the *Domain Name System* (DNS) “cache poisoning” techniques discovered a few years ago. And you have most likely heard that *Domain Name System Security Extensions* (DNSSEC)^[0, 13, 14, 15] is the long-term cure. But you might not know exactly what challenges are involved with DNSSEC and what experience the early adopters have gathered and documented. Perhaps you waited with your own rollout until you could gather more documentation about operational experiences when rolling out DNSSEC.

Stephan Lagerholm and Torbjörn Eklöv are DNS architects with significant DNSSEC experience. Torbjörn lives in Sweden and has helped several municipalities, as well as other organizations, sign their zones. Stephan Lagerholm lives in Dallas, Texas, and has been involved in implementing DNSSEC at several U.S. federal agencies. This article summarizes their experiences, including lessons learned from implementing the technology in production environments, and discusses associated operational concerns.

Background

A plethora of information about DNSSEC and cache poisoning attacks is available on the Internet^[16], so we will not repeat it, but we think it is important to state where DNSSEC is today.

During the last few years the number of deployments, as well as the size and importance of the signed domains, has increased significantly. One of the main reasons for adoption of the DNSSEC during the past year was that the U.S. *Office of Management and Budget* (OMB) issued a mandate requiring the signing of the **.gov** domain in the beginning of the year. U.S. federal agencies were mandated to sign their domains by the end of 2009. Some agencies have already implemented the technology, whereas others are still working on it.^[1]

Acceptance of DNSSEC technology is also reaching outside of the U.S. government. *Top Level Domains* (TLDs) around the globe have announced DNSSEC initiatives. To mention a few, Afilias signed **.org** and Neustar recently announced signing of **.us**. Several *County Code TLDs* (ccTLDs), including **.nl** and **.de**, announced that DNSSEC implementation is a work in progress. VeriSign has announced that it is working on signing the largest TLDs, namely **.com** and **.net**. Finally, the *Internet Corporation for Assigned Names and Numbers* (ICANN) along with VeriSign released a timeline for signing the root zone. And of course, the pioneer **.se** is on its fourth year as a signed TLD.

Several vendors have released software and products to support and make the signing of zones easier. A range of different products is now available on the market.

DNS professionals now have a broad choice of technology—from collections of open-source signing scripts to advanced systems with full automation and support for *Federal Information Processing Standard* (FIPS)-certified cryptography.

Operational Challenges

DNSSEC might significantly affect operations unless it is carefully implemented because it requires some changes to the underlying DNS protocol. Those changes are, in fact, the first significant changes that have been made to the DNS protocol since it was invented. Those changes might sometimes fool old systems into believing that the packets are illegal. DNSSEC also introduces new operational tasks such as rolling the keys and resigning the zone. Such tasks must be performed at regular intervals. Furthermore, as with any new technology, there are misconceptions about how to interpret the RFC standard.

The First Bug Reported

Late summer 2007, Torbjörn Eklöv convinced the municipality of Gävle in Sweden of the benefits of DNSSEC. He proudly signed what is believed to be the first municipality zone in the world, **gavle.se**. At first, everything worked fine. A week or so later, Gävle received reports from citizens who could not reach the municipality's websites. It turned out that a new version of *Berkeley Internet Name Domain* (BIND) was rolled out by a large service provider and that this version of BIND introduced a rather odd bug that affected DNSSEC. The result of the bug was that home users with some home routers and firewalls could not reach any signed domains.

Some people who heard about the problem at **gavle.se** wrongly believed that DNSSEC caused the problem and that DNSSEC is broken. However, this assumption is not true; DNSSEC worked as expected, but a bug in a particular version of BIND caused the problem. The problem triggered some research on how home routers handle DNSSEC. *Stiftelsen för Internetinfrastruktur*, the organization that runs the **.se** TLD, issued a report describing how commonly used home routers and firewalls handled the new protocol changes in DNS^[2]. Later, Nominet, which administers the **.uk** TLD, issued a similar report^[3]. In addition, DENIC, which administers the **.de** TLD, researched the same subject^[4]. The results are all discouraging; only 9 out of 38 tested home gateways supported DNSSEC correctly in the most recent reports.

A *Birds of a Feather* (BoF) session was held at the 76th meeting of the *Internet Engineering Task Force* (IETF) in Hiroshima to discuss the problems involving home gateways^[5]. We look forward to seeing progress in this area.

Preparing Your Firewall for DNSSEC

Most problems with DNSSEC are related to firewalls. Make sure to involve your security and networking administrators so that they can make the required changes before taking DNSSEC into production.

Two types of firewall problems are most common:

The first involves the *Transmission Control Protocol* (TCP). There is a misconception among firewall vendors and security administrators that DNS queries use the *User Datagram Protocol* (UDP) and that zone transfers use TCP. Unfortunately, this assumption is not entirely true. DNS queries first try UDP, but revert to TCP if no response is received for the initial UDP query or if the response lacks important information because it is truncated. The possibility of something in the path blocking the response to the initial query is much higher with DNSSEC because of the increased size of the responses.

For DNSSEC to work correctly, it is mandatory that you open your firewall for both TCP and UDP over port 53.

The second problem is related to the *IP Buffer Reassembly* size. The authors of the DNSSEC standard realized that a potential problem might exist with TCP queries. TCP puts a higher burden on the DNS servers. (TCP is much more expensive to process than UDP.) To avoid too much TCP traffic, the authors made the EDNS0 extension mandatory for DNSSEC. EDNS0 is one of the *Extension Mechanisms for DNS* (EDNS), a standard that, among other things, allows a client to signal that it is capable of receiving DNS replies over UDP that are larger than the previous limit of 512 bytes. Some firewalls are not aware of the fact that the EDNS0 standard allows for larger packets and they either block any DNS packet using EDNS0, or block any DNS packet larger than the 512 bytes regardless of the EDNS0 signaling.

Other firewalls allow for the large packets by default, whereas a few vendors require the firewall to be manually configured to do so. Any device in the path that does packet inspection at the application layer must be aware of the EDNS0 standard to be able to make a correct decision about whether to forward the packet or not. ICANN has summarized the status of EDNS0 support in some commonly used firewalls^[6].

Note that it is not enough to test that your firewall allows large incoming DNS replies by sending DNS queries to the Internet^[7]. You must also test that an external source can receive large DNS replies that your DNS server is sending. One way of doing so is to use an open DNSSEC-aware resolver^[8, 9].

Test and configure your firewall to allow for use of EDNS0 and for DNS packets larger than 512 bytes over UDP.

Preparing Your Slaves

Setting up DNSSEC involves substantial changes to the master name server so it can sign and serve the signed data. However, it is easy to foresee that the slaves must be upgraded, too. The slaves are much easier to upgrade and operate because they never produce signatures.

They are secondary systems that transfer data from the primary server and respond to DNS queries. But the slaves must understand how to respond to queries requesting signed data.

Slaves must be upgraded to BIND 9.3 or better to understand the *Next Secure* (NSEC)^[14] standard. NSEC is a method to provide authenticated denial of existence for DNS resource records. The newer *Next Secure 3* (NSEC3)^[10] standard introduces some additional requirements for the slaves. If you use NSEC3, you must upgrade the slaves to BIND 9.6 or later. Version 3 of *Name Server Daemon* (NSD)^[17] and any version of *Secure64 DNS Authority/Signer*^[18] can do both NSEC and NSEC3. Windows Server 2008 R2 for the x86-64 architecture supports DNSSEC as a master, slave, and validating resolver. However, we recommend limiting the use of the Windows platform to slaves and for domains using NSEC. Our opinion is that it is very hard to implement DNSSEC on Windows, and we suggest that you wait until Microsoft offers a sensible *Graphical User Interface* (GUI) and support for NSEC3. Note that the Itanium version of Windows 2008 R2 supports neither DNS nor DNSSEC.

Make sure your slaves can handle the version of DNSSEC you intend to use.

If the slaves are administered by another party, contact the administrator before you begin DNSSEC implementation. Make sure the slaves are running a version capable of DNSSEC. Stephan helped a large U.S. federal agency sign its domains. The agency used one of the major federal contractors to run its slave servers. After multiple attempts to reach somebody that understood DNS and DNSSEC, Stephan finally learned that the slaves were running BIND 9.2.3 and that the contractor had no plans to upgrade. The only alternative for the agency was to in-source the slaves and run them itself.

If your slaves are administered by another party, make sure you know if and what version of DNSSEC that party supports before you start implementing.

Communicate with Your Parent

TLDs allow you to communicate with them in two ways:

- *Registrant–Registrar–Registry Model:* In this, the most common model, the registrant (**example.org**) does not communicate directly with the registry (**.org**). Instead, a third-party registrar handles all communication related to DNS and DNSSEC. This model is, for example, used by the **.se** and **.org** TLDs.
- *Registrant–Registry Model:* This model is normally used by smaller TLDs such as **.gov**. It allows direct communication between the registrant (**agency.gov**) and the registry (**.gov**). The TLD acts as both a registrar and a registry in this model.

Most problems described in the following paragraphs apply to both models, but those involving multiple registries are obviously applicable only to the Registrant–Registrar–Registry model.

Establishing a *Chain of Trust* in DNSSEC involves uploading one or more public keys to the parent. Ultimately the parent publishes a *Delegation Signer* (DS) record, a smaller fingerprint that can be constructed from the DNSKEY record. To upload your keys, you must use a registrar that supports DNSSEC. If your registrar does not support DNSSEC, you need to move your domains to another registrar (or convince your current registrar to start supporting DNSSEC). It usually takes a few days or up to a week to move a domain from one registrar to another.

Make sure that your registrar supports DNSSEC. If it does not, move your domain to a registrar that supports DNSSEC before you begin signing your zone.

Some registrars allow registration under multiple TLDs. However, just because a registrar handles DNSSEC for one TLD does not mean that it handles DNSSEC for all TLDs it serves. For example, several registrars in Sweden support DNSSEC for `.se` but not for `.org` or `.us`.

Make sure that your registrar handles DNSSEC under the TLD in question.

Most registrars offer you the opportunity to use their name server instead of your own. The service is either offered for free or for an additional cost. The registrar typically provides a web interface where you can change your zone data. This service is a good and useful choice if your domains are uncomplicated and small. Larger and more complex domains are better operated on your own servers.

Some registrars that provide this type of service can handle DNSSEC only if you use their name servers and not your own name servers. These registrars can establish the chain of trust with the parent only if the zone is under their control. They lack a user interface for uploading a DS key that you generate on your own name servers.

If you intend to use your own name servers, make sure that your registrar supports this deployment model, and allows you to upload a DS record for further distribution to the registry.

In theory, the child zone system should create the DS record fingerprint and upload it to the parent. In practice, some registrars require you to upload the DNSKEY record to them. They then create the DS record for you. (This practice is bad because the registrar must know the hash algorithm used to construct the DS record, which it might not know.) The DNSKEY record comes in several different formats, depending on the platform you used to create the keys (BIND, Microsoft, NSD, Secure64, etc.). The formats have minor differences, and you might have to convert the DNSKEY into a format that the registrar accepts.

Not everything works smoothly, even with the correct DNSKEY format. The logic at one registrar's website was to deny uploading of DNSKEYs unless the optional *Time To Live* (TTL) field existed. (The TTL value is useless in the DNSKEY context because the parent overrides this value with its own TTL). You may have to manually change your DNSKEY before uploading it to comply with the checks that the registrar performs.

If your registrar requires you to upload the DNSKEY, make sure that your solution can generate the requested format. If not, you need to manually change the fields with a text editor.

As noted previously, some registrars are performing too many checks and irrelevant checks before accepting and creating the secure delegation. Other registrars do not check at all or have limited checks that do not work as expected. For example, some registrars assume that your key is created using a certain algorithm, and they do not double-check it prior to creating a DS record. One registrar created a bogus DS record if you uploaded a DNSKEY with upper-case characters in the domain name. The bogus DS record looked valid, and troubleshooting to find this error took hours.

Another example is keys created with *Webmin*^[11], a graphical tool that you can use for signing zones. Webmin defaults to using the less-common *Digital Signature Algorithm* (DSA) for its DNSKEYs. The registrar did not complain when uploading the Webmin key, and it created a bogus DS record by assuming that it was an RSA key.

It is hard for a registrant to do anything about errors at the registrars. The best you can do is to make sure that you upload the correct key with the correct parameters such as algorithm, key length, key-id, etc. If something goes wrong, you might have to change the keys in production. Rolling the keys to the same algorithm and key length is relatively easy—but changing your keys to another algorithm adds extra complexity. It is an interesting exercise to change to another algorithm in production, but it is something we recommend avoiding if possible.

Double-check the DNSKEY/DS so that it is created with the correct parameters prior to uploading it.

Communicate with Your Children

If you have sub-domains in your domain, you must make sure that you can accept and publish the DS records that your children upload to you. This situation is not a problem if you use zone files in text format—you can simply insert the DS record using your favorite editor. But it might be a problem if you are using an *Internet Protocol Address Management* (IPAM) system. In that case make sure that it can insert DS records into the zones that are managed by the system. Some IPAM systems do not support insertion of DS records correctly.

Make sure that your IPAM system can insert DS records into your zones.

A common strategy among organizations with high-availability requirements for their critical servers is to use a global load balancer, which is basically a DNS server that responds differently depending on the status of the service in question. For example, assume a load balancer can respond to a question for `www.example.com` with `192.0.2.1` and `192.0.2.2` if both web servers are up. If `.1` becomes unavailable, the load balancer notices a failure and responds only with `.2`. In order to use a global load balancer, you must delegate `www` as a sub-domain to its own DNS process.

When DNSSEC is implemented, you must make sure that the load balancer can handle DNSSEC (and not that many do); otherwise it is impossible to sign the responses for those resources. Unfortunately, these resources are the most critical ones for your environment and would benefit the most from DNSSEC signing.

Make sure that your load balancers support DNSSEC. If they do not, have an alternative strategy.

Rolling the Keys

You should change the DNSKEYs regularly and when you think the keys are compromised. The process of doing so is called *rolling the keys*. There are normally two different keys in DNSSEC, the *Key Signing Keys* (KSKs) and the *Zone Signing Keys* (ZSKs). Rolling the ZSK is an internal process and does not require communication with the parent. Rolling the KSK, on the other hand, requires the parent to publish a new DS record.^[12]

There is no standard yet that describes how the communication between the parent and the child should occur when a key is rolled. Early DNSSEC-capable registrants used a web interface that allowed their registrants to upload and manipulate the DNSSEC information. With a web interface, each domain must be handled separately and there is no easy way to automate the interaction.

The web interface works for a handful of domains but becomes very cumbersome when you have many domains. For those types of organizations, it is important to make sure that there is some kind of *Application Programming Interface* (API) or script access to the registrar. This interface allows the organization to upload new DS records during the rollover in a convenient way.

Make sure that your registrar supports automation through an API if you have many domains.

Scripting with an API as described previously is one way of communicating with the registrar. Another way of achieving the same type of automation is for the parent (or registrar) to monitor the child for any changes to the DNSKEY records.

Note that the chain of trust is still intact during a nonemergency rollover. The parent can securely poll the child and grab the new DNSKEY records and convert them into DS records. The polling from the parent to each signed child needs to occur regularly so that a rollover is picked up quickly. This regularity of polling makes the scheme best for domains with fewer delegations (in the order of thousands, not millions—consider how much bandwidth an hourly polling of 15 million children would require).

Automation is a good thing, but make sure you understand the implications when opting for automatic detection of key rollovers. The automation scripts are not fail-safe. It has been reported that early versions of such scripts under some circumstances wrongly assumed that a key rollover occurred and deleted the DS record, thus breaking the chain of trust.

Understand the implication when opting for automatic detection, addition, and deletion of DS records.

Management of DNSSEC

Without DNSSEC, you are not bound to any particular registrar; you can switch to a new registrar fairly easily. With DNSSEC, this situation changes. First of all, if you let the registrar sign the zone on your behalf, the registrar will be in charge of the key used to sign your zone. Extracting your key so that it can be imported to another registrar is not always straightforward (also remember that there is really no incentive for your previous registrar to help you because you just discontinued its service). An alternative is to unsign the zone before you change registrars, but that option might not always be a viable one. The lack of standards makes it hard to change registrars on a signed domain that is in production.

You must tell your new registrar that you are using DNSSEC, and you must make sure that the registrar supports it. If not, the registrar might accept the transfer but be unable to publish the DNSKEY records. The result would be a DS record published by the registry but no corresponding DNSKEY records at the child, making the zone “security lame” and causing failed validation.

The same types of problems exist if you are running your own name servers. If you change your master server, make sure that you transfer the secret keys as well. Signing with new keys will not work unless you flush out the old keys with rollovers and upload a new DS record to your parent.

Have a plan ready for how to transfer your keys to a new master server.

Timers

It is important to adjust your signature validity periods and the *Start of Authority* (SOA) timers so that they match your organizational requirements and operational practices. SOAs expire and signature validity periods all too often are too short.

Unless you are restricted by guidelines saying otherwise, you should strive to set the timers reasonably high. Set the timers so that your zones can cope with an outage as long as the longest period that the system might be unattended.

For example, if you know that your top DNS administrator usually has three weeks of vacation in July, you could consider setting the times so that the zone can survive four weeks of downtime. If you are confident in your signing solution and are monitoring your signatures carefully, you might set it a little bit lower.

Signature lifetime is a trade-off between security (low signature lifetimes) and convenience (high signature lifetimes). Setting a really high signature lifetime is convenient from an operational perspective but is less secure. Some organizations such as the IETF use an excessive signature lifetime of one year (`dig ietf.org DNSKEY +dnssec | grep RRSIG`). This lifetime is clearly not recommended, and they should know better.

Carefully set your signature lifetimes and SOA times to reflect your organization's operational requirements and practices.

A Note on Validation

This article has focused on the authoritative part of DNSSEC. That part includes signing resource records and serving DNS data. The operational challenges with signing data are much greater than the challenges of validating data. To validate data, the only thing you need to do regularly is update your trust anchor file. Make sure you do so. Torbjörn reports several outages when the `.se` DNSKEY used in the `.se` trust anchor expired in January 2010. We look forward to the work being done in this area to automate the process.

Summary

DNSSEC has been deployed and taken in production for several large and critical domains. It is not hard to implement DNSSEC, but doing so introduces some operational challenges. Those challenges exist both during the implementation phase when the zone is being signed for the first time and during the operation of the zone. Make sure you understand the possible effects of implementation and plan ahead. The following checklist summarizes the most important pitfalls with DNSSEC:

- Open your firewall for EDNS0 signaling and allow large DNS packets using UDP and TCP over port 53.
- Check the DNSSEC capabilities of all your masters and slave servers.
- Check the DNSSEC capabilities of your registrar and understand their requirements for the public key you are uploading.
- Make sure your IPAM system can handle secure delegations.

- Plan how to handle load balancers.
- Develop an automation strategy if you have a lot of zones.
- Plan how you will transfer your keys to a new master server if a disaster occurs.
- Implement a policy for DNSSEC timer settings.

Happy signing!

For Further Reading

- [0] Miek Gieben, “DNSSEC: The Protocol, Deployment, and a Bit of Development,” *The Internet Protocol Journal*, Volume 7, No. 2, June 2004.
- [1] Carolyn Duffy Marsan, “80% of Government Web Sites Miss DNS Security Deadline,” *Network World*, January 21, 2010, <http://www.networkworld.com/news/2010/012010-dns-security-deadline-missed.html>
- [2] Jaokim Ålund and Patrik Wallström, “DNSSEC—Tests of Consumer Broadband Routers,” http://www.iis.se/docs/Routertester_en.pdf
- [3] Ray Bellis and Lisa Phifer, “Test Report: DNSSEC Impact on Broadband Routers and Firewalls,” September 2008, <http://download.nominet.org.uk/dnssec-cpe/DNSSEC-CPE-Report.pdf>
- [4] Thorsten Dietrich, “DNSSEC-Unterstützung durch Heimrouter,” http://www.denic.de/fileadmin/Domains/DNSSEC/DNSSEC_20100126_Dietrich.pdf
- [5] Broadband Home Gateway BoF, <http://tools.ietf.org/agenda/76/homegate.html>
- [6] ICANN DNS Root Server System Advisory Committee (RSSAC) and Security and Stability Advisory (SSAC), “Testing Firewalls for IPv6 and EDNS0 Support,” January 2007. <http://www.icann.org/en/committees/security/sac016.htm>
- [7] Domain Name System Operations Analysis and Research Center (OARC)’s DNS Reply Size Test Server: <https://www.dns-oarc.net/oarc/services/replysizetest>
- [8] OARC’s Open DNSSEC Validating Resolver: <https://www.dns-oarc.net/oarc/services/odvr>
- [9] Comcast DNSSEC Information Center, <http://www.dnssec.comcast.net/>

- [10] Torbjörn Eklöv, “DNSSEC: Will Microsoft Have Enough Time?” *CircleID*, January 2010, http://www.circleid.com/posts/dnssec_will_microsoft_have_enough_time/
- [11] <http://www.webmin.com/>
- [12] George Michaelson, Patrik Wallström, Roy Arends, and Geoff Huston, “Rolling over DNSSEC Keys,” *The Internet Protocol Journal*, Volume 13, No. 1, March 2010.
- [13] Roy Arends, Rob Austein, Dan Massey, Matt Larson, and Scott Rose, “DNS Security Introduction and Requirements, RFC 4033, May 2005.
- [14] Roy Arends, Rob Austein, Matt Larson, Dan Massey, and Scott Rose, “Resource Records for the DNS Security Extensions,” RFC 4034, May 2005.
- [15] Roy Arends, Rob Austein, Dan Massey, Matt Larson, and Scott Rose, “Protocol Modifications for the DNS Security Extensions,” RFC 4035, May 2005.
- [16] United States Computer Emergency Readiness Team (US-CERT), “Multiple DNS Implementations Vulnerable to Cache Poisoning,” July 2008, <http://www.kb.cert.org/vuls/id/800113>
- [17] <http://nlnetlabs.nl/projects/nsd/>
- [18] <http://www.secure64.com/secure-DNS>

STEPHAN LAGERHOLM is a Senior DNS Architect with Secure64 Software Corporation—a software company offering high-performance DNS server software that makes the DNS trustworthy and secure. Secure64 DNS applications include key management and zone-signing software that make it easy to deploy DNSSEC securely and correctly as well as DNS server software that is always available. Stephan is a DNS and security expert with more than 11 years of international experience in the field. His background includes leadership positions at the largest networking and security system integrator in Scandinavia, and responsibility for designing hundreds of complex IT networks. Stephan is one of the few persons in the United States to have integrated DNSSEC into production environments. Stephan is CISSP-certified and holds a Master of Science degree in Computer Science and Mathematics from Uppsala University in Sweden. E-mail: Stephan.Lagerholm@secure64.com

TORBJÖRN EKLÖV is the founder and partner of Interlan Gefle AB, an IT consulting company in Sweden with 20 employees. He is a DNSSEC and IPv6 pioneer. All internal and external services at Interlan use both IPv6 and IPv4, and the company hosts about 200 DNSSEC-signed domains. Torbjörn has worked with Internet communication and security for 15 years, and is the founder and manager of Secure End User Connection (SEC), or Säker KundAnslutning (SKA) in Swedish, an organization that certifies products and broadband networks to protect subscribers from spoofing and hijacking. His favorite homepage is <http://test.ipv6.tk>. You can reach him at Torbjorn.Eklov@interlan.se

Book Review

The Art of Scalability

The Art of Scalability: Scalable Web Architecture, Processes, and Organizations for the Modern Enterprise, by Martin L. Abbott and Michael T. Fisher, ISBN-13: 978-0-13-703042-2, Pearson Education, 2010.

It is often claimed that the primary lesson of the Internet is one of “scaling.” So the title of this book bodes well for relevance to Internet designers. A reader would likely expect discussion of hashing algorithms, fast-path coding, protocol latencies and chattiness, distributed redundancy design, and similar guidance for handling a billion users. The reader would largely be wrong, although some of the book is dedicated to technical performance. What is easily missed in the title is the word “organizations.” It does not mean organization of modules. It means organizations within a *company*.

This book is very much a holistic one. It takes the painfully realistic position that well-designed protocols and software modules matter only if the company structure or team operation is tuned to growing and running a large-scale service. The book is comprehensive and primarily tailored for highly formal management, with substantial, bureaucratic procedures designed to ensure thorough consideration of scalability needs and implications. It is loaded with discussion of many different organizational and technical management tools that assist in making diligent decisions. For most readers and most companies, attempting to apply this level of formality is dramatic overkill. However, knowing about it is not.

The book is 533 pages, with 33 chapters and 3 appendices. The writing style is reasonably clean, but pedantic. Don't expect the type of entertainment-oriented writing that is common these days. The authors' experiences include *eBay* and *PayPal*, so scaling matters have been within their direct work responsibilities. As holds for any book attempting this kind of breadth, from technology design to organization management, discussion frequently is superficial and will be obvious to some readers, while the specific detail will in places be irrelevant to many others. Although these characteristics might be taken as negatives, they actually serve to demonstrate the utility of the book as an introduction and basic reference to the topic of scaling. A quick scan of the book helps the reader see how many different aspects of an organization's activities can aid or hinder large-scale operations. Exploring specific chapters can explain concepts and topics and suggest particular tools to help in planning or analysis.

Organization

Part I, “Staffing a Scalable Organization,” comprises six chapters. It provides a tutorial on classic problems in structuring and staffing an organization for growth. Little is taken for granted. So there is guidance about the characteristics needed in a CEO, CFO, or CTO for aiding leadership in working to scale the company and the company’s products. It even has a chapter on “Leadership 101.”

For the most part, this section is likely to be useful only for readers with no management background, because the material is extremely basic. What distinguishes it is only the constant consideration of the way its topics are relevant to scaling. The likely utility of the section is in helping employees “manage up” so they can interact with management better when seeking support for changes needed to implement or maintain scalable development or operations. On the other hand, an interesting discussion explored why some simple and entirely logical choices for organizing a company work against accountability and scaling.

Part II, “Building Processes for Scale,” at nearly 200 pages is 40 percent of the book. Whereas the first part concerned the people, this one concerns what they do. The first half of this part strongly emphasizes processes for anticipating and responding to scaling problems and for judiciously allocating limited resources. Hence there is even a chapter that considers “build versus buy.” Technical topics discussed here are conceptual rather than concrete. They concern risk, performance, capacity, and failure recovery. Each is treated as a planning and design concern, with estimates and procedures. A warning: The word “architecture” shows up in the title of several middle paragraphs in this section, but don’t be confused. It refers to groups that do architecture, not to the technical details of architecture.

Part III is “Architecting Scalable Solutions.” Now at last, techies will start to get their geek fix. But perhaps with more abstraction than they will expect? Again, this book is more about properly organizing things than about algorithms. The section introduces “technology-agnostic design,” with consideration of fault isolation and various growth factors, including repeated attention to cost, risk, scalability, and availability. There are chapters on database scaling and the use of caching for performance. The authors are fond of asynchronous and state-free interaction, with the view that it is more robust. The precise reason for this conclusion was not entirely clear to me, but presumably it is because it is easier to recover and retarget an exchange after an outage occurs during an interaction.

Two chapters of this part of the book are devoted to the “AKF Scale Cube,” and indeed the Index has a large number of citations to it. (AKF refers to the authors’ company.) For this analytic tool, the x-axis “...represents cloning of services and data with absolutely no bias.” In other words, these graphs are pure replications of equivalent, parallel components or activities, used to distribute load. The y-axis “... represents a separation of work responsibility by either the type of data, the type of work performed for a transaction, or a combination of both... We often refer to these as service or resource oriented splits.” The nature of the z-axis is described as “...biased most often by the requestor or customer... focused on data and actions that are unique to the person or system performing the request.” I took this as meaning that the axis divides work according to tailored attributes.

Part IV is the catchall for remaining topics, with some requisite discussion of clouds and grids, application monitoring, and data center planning.

Summary

The book will be useful for architects who need to understand how to scale their own work and how to support their organization for long-term growth. It will also be useful for technical, operations, and other managers who need to understand the technical and operations scaling problems, support their own architects, and work with the rest of their organization to anticipate and satisfy scaling requirements.

—Dave Crocker, *Brandenburg Internet Working*
dcrocker@bbiw.net

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. Contact us at ipj@cisco.com for more information.

Call for Candidates for Itojun Service Award

The *Itojun Service Award* is presented every year to an individual or a group who has made outstanding contributions in service to the IPv6 community. The deadline for nominations for this year's award is July 12, 2010. The award will be presented at the 79th meeting of the *Internet Engineering Task Force* (IETF) to be held in November 2010 in Beijing, China.

The Itojun Service Award, established by the friends of Itojun and administered by the *Internet Society* (ISOC), recognizes and commemorates the extraordinary dedication exercised by Itojun over the course of IPv6 development. The award includes a presentation crystal, a US\$3,000 honorarium, and a travel grant.

The award is focused on pragmatic technical contributions, especially through development or operation, with the spirit of servicing the Internet. With respect to the spirit, the selection committee seeks contributors to the Internet as a whole; open source developers are a common example of such contributors, although this is not a requirement for expected nominees. While the committee primarily considers practical contributions such as software development or network operation, higher-level efforts that help those direct contributions will also be appreciated in this regard. The contribution should be substantial, but could be immature or ongoing; this award aims to encourage the contributors to continue their efforts, rather than just recognizing well-established work. Finally, contributions of a group of individuals will be accepted as deployment work is often done by a large project, not just a single outstanding individual.

The award is named after Dr. Jun-ichiro "Itojun" Hagino, who passed away in 2007, aged just 37. Itojun worked as a Senior Researcher at *Internet Initiative Japan Inc.* (IIJ), was a member of the board of the *Widely Integrated Distributed Environment* (WIDE) project, and from 1998 to 2006 served on the groundbreaking KAME project in Japan as the "IPv6 Samurai." He was also a member of the *Internet Architecture Board* (IAB) from 2003 to 2005.

For additional information on the award, please visit:

<http://www.isoc.org/awards/itojun/>

Less than 10% of IPv4 Addresses Remain Unallocated, says NRO

The *Number Resource Organization* (NRO), the official representative of the five *Regional Internet Registries* (RIRs) that oversee the allocation of all Internet number resources, recently announced that less than 10 percent of available IPv4 addresses remain unallocated. This small pool of existing IP addresses marks a critical moment in IPv4 address exhaustion, ultimately impacting the future network operations of all businesses and organizations around the globe.

“This is a key milestone in the growth and development of the global Internet,” noted Axel Pawlik, Chairman of the NRO. “With less than 10 percent of the entire IPv4 address range still available for allocation to RIRs, it is vital that the Internet community take considered and determined action to ensure the global adoption of IPv6. The limited IPv4 addresses will not allow us enough resources to achieve the ambitions we all hold for global Internet access. The deployment of IPv6 is a key infrastructure development that will enable the network to support the billions of people and devices that will connect in the coming years,” added Pawlik.

The *Internet Protocol* (IP) is a set of technical rules that defines how devices communicate over a network. There are currently two versions of IP, IPv4 and IPv6. IPv6 includes a modern numbering system that provides a much larger address pool than IPv4. With so few IPv4 addresses remaining, the NRO is urging all Internet stakeholders to take immediate action by planning for the necessary investments required to deploy IPv6.

The NRO, alongside each individual RIR, has actively promoted IPv6 deployment for several years through grassroots outreach, speaking engagements, conferences and media outreach. To date, their combined efforts have yielded positive results in the call to action for the adoption of IPv6.

Given the less than 10 percent milestone, the NRO is continuing its call for Internet stakeholders, including governments, vendors, enterprises, telecoms operators, and end users, to fulfill their roles in IPv6 adoption, specifically encouraging the following actions:

- The business sector should provide IPv6-capable services and platforms, including web hosting and equipment, ensuring accessibility for IPv6 users.
- Software and hardware vendors should implement IPv6 support in their products to guarantee they are available at production standard when needed.
- Governments should lead the way by making their own content and services available over IPv6 and encouraging IPv6 deployment efforts in their countries. IPv6 requirements in government procurement policies are critical at this time.
- Civil society, including organizations and end users, should request that all services they receive from their ISPs and vendors are IPv6-ready, to build demand and ensure competitive availability of IPv6 services in coming years.

The NRO’s campaign to promote the next generation of Internet Protocol continues to positively impact the Internet community. IPv6 allocations increased by nearly 30% in 2009, as community members continued to recognize the benefits of IPv6.

“Many decision makers don’t realize how many devices require IP addresses—mobile phones, laptops, servers, routers, the list goes on,” said Raul Echeberria, Secretary of the NRO. “The number of available IPv4 addresses is shrinking rapidly, and if the global Internet community fails to recognize this, it will face grave consequences in the very near future. As such, the NRO is working to educate everyone, from network operators to top executives and government representatives, about the importance of IPv6 adoption,” added Echeberria.

IP addresses are allocated by the *Internet Assigned Numbers Authority* (IANA), a contract operated by the *Internet Corporation for Assigned Names and Numbers* (ICANN). IANA distributes IP addresses to RIRs, who in turn issue them to users in their respective regions. “This is the time for the Internet community to act,” said Rod Beckstrom, ICANN’s President and Chief Executive Officer.

“For the global Internet to grow and prosper without limitation, we need to encourage the rapid widespread adoption of the IPv6 protocol,” he added.

The NRO is the coordinating mechanism for the five RIRs. The RIRs—Afrinic, APNIC, ARIN, LACNIC, and the RIPE NCC—ensure the fair and equitable distribution of Internet number resources (IPv6 and IPv4 addresses and *Autonomous System* (AS) numbers) in their respective regions. The NRO exists to protect the unallocated Internet number resource pool, foster open and consensus-based policy development, and provide a single point of contact for communication with the RIRs.

Learn more about the NRO at www.nro.net/media

The five RIRs that make up the NRO are independent, not-for-profit membership organizations that support the infrastructure of the Internet through technical coordination. The IANA allocates blocks of IP addresses and ASNs, known collectively as *Internet number resources*, to the RIRs, who then distribute them to users within their own specific service regions. Organizations that receive resources directly from RIRs include *Internet Service Providers* (ISPs), telecommunications organizations, large corporations, governments, academic institutions, and industry stakeholders, including end users. The RIR model of open, transparent participation has proven successful at responding to the rapidly changing Internet environment. Each RIR holds one or two open meetings per year, as well as facilitating online discussion by the community, to allow the open exchange of ideas from the technical community, the business sector, civil society, and government regulators.

The five RIRs are:

- AfriNIC: <http://www.afrinic.net>
- APNIC: <http://www.apnic.net>
- ARIN: <http://www.arin.net>
- LACNIC: <http://www.lacnic.net>
- RIPE NCC: <http://www.ripe.net>

ISOC Funds Projects to Support Internet Access, Security, and Policy Development

The *Internet Society* (ISOC) recently announced it is funding community-based projects around the world addressing issues such as Internet leadership, education, core infrastructure, local governance, and policy development, with a strong focus on currently underserved communities.

“The diversity of projects awarded highlights the profound importance of the Internet in so many aspects of our lives, in all parts of the world,” said Jon McNerney, Chief Operating Officer of the Internet Society. “The passion and creativity of those developing the projects within their communities drives the Internet Society’s commitment to help bring the benefits of the Internet to people everywhere.”

As part of the *ISOC Community Grants Program*, each project will receive up to US\$10,000 for efforts that promote the open development, evolution, and use of the Internet for the benefit of all people throughout the world.

Projects funded in this round include:

- Training programs to build digital literacy within safe environments in India and Uganda
- Village-operated telecommunication services in East Timor
- Support for development of core Internet time infrastructure
- Policy and practical action in Kenya to improve online safety for women
- Online support for NGOs in Tunisia and more effective local governance in India
- Promotion of Internet leadership in Ecuador
- Development of important public policy resources in Georgia and Australia

ISOC Community Grants are awarded twice each year. The next round of the program will open on September 1, 2010. Additional information about the Community Grants Program and this round of award-winning projects can be found here:

<https://www.isoc.org/isoc/chapters/projects/index.php>

<https://www.isoc.org/isoc/chapters/projects/awards.php?phase=11>

RIPE Community Statement on the Internet Address Management System

At the May 2010 *Réseaux IP Européens* (RIPE) meeting in Prague, Czech Republic, the RIPE community issued the following statement:

“The RIPE community supports all efforts to assist in the deployment of IPv6, especially in developing countries.

However, we note concerns being expressed within the ITU by a few members, most recently in the ITU IPv6 Group, that the current address management system is inadequate.

The RIPE community mandates the RIPE NCC to work with the ITU IPv6 Group, individual ITU members, and the community to clearly identify these concerns and to find ways to address them within the current IP address management system.”

This statement will be sent to the *International Telecommunications Union* (ITU) to reiterate the RIPE community’s belief that the current address management system works. The RIPE NCC will continue to participate actively in the ITU IPv6 Group and report back to the RIPE community.

For more information see:

<http://www.itu.int/ITU-T/othergroups/ipv6/>

<http://ripe.net/ripe/index.html>

<http://www.nro.net/documents/nro51.html>

Upcoming Events

The *North American Network Operators’ Group* (NANOG) will meet in San Francisco, California, June 13–16, 2010.

See <http://nanog.org>

The *Internet Corporation for Assigned Names and Numbers* (ICANN) will meet in Brussels, Belgium, June 20–25, 2010.

See <http://icann.org>

The *Internet Engineering Task Force* (IETF) will meet in Maastricht, The Netherlands, July 25–30, 2010 and in Beijing, China, November 7–12, 2010. See <http://www.ietf.org/>

APNIC, the *Asia Pacific Network Information Centre*, will hold its Open Policy meeting in the City of Gold Coast, Australia, August 24–28, 2010. See <http://www.apnic.net/meetings/30/>

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ contains standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2010 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

September 2010

Volume 13, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
PMIPv6	2
Happy Eyeballs.....	16
Letter to the Editor	22
Fragments	24

FROM THE EDITOR

Technology advances—such as improvements in display technology, battery life, processor capabilities, and communications systems—have all contributed to making *mobile devices* the most important area for Internet growth. In order to fully support these devices, the IETF developed *Mobile IP* many years ago, and it has continued to work on the general area of IP mobility. We have covered some of this work in previous issues of IPJ, and this time we look at *Proxy Mobile IPv6* (PMIPv6), which is being standardized by the IETF. The article is by Ignacio Soto, Carlos J. Bernardos, María Calderón, and Telemaco Melia.

Deployment of IPv6 is progressing, albeit slowly. In several upcoming articles we will examine some transition technologies or implementation details that can make this deployment easier, and above all, transparent, to the end user. In our first article, Dan Wing and Andrew Yourtchenko explain the concept of “Happy Eyeballs” as applied to dual-stack IPv4/IPv6 systems.

Domain Name System Security Extensions (DNSSEC) have recently been applied to the Internet system of root servers. For details, see our “Fragments” section, where you will also find a statement from the *Number Resource Organization* (NRO) regarding the results of a recent IPv6 readiness study.

Once again, please remember to check your subscription expiration date and take the necessary steps if you wish to continue receiving this journal. It's not too late to renew and get back on the distribution list, even if your subscription expired some time ago. You can find your subscription ID and expiration date either on the back page of your copy or on the envelope that it came in. In order to access your record, click the “Subscriber Services” link on our webpage at www.cisco.com/ipj and enter your e-mail address and the subscription ID. The system will send you a link that allows direct access to your record, and you can update your address and renew your subscription. If you no longer have access to the e-mail you used when you subscribed or have forgotten your subscription ID, just send a message to ipj@cisco.com and we will make the necessary changes for you.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

PMIPv6: A Network-Based Localized Mobility Management Solution

by Ignacio Soto, Universidad Politécnica de Madrid; Carlos J. Bernardos, and María Calderón, Universidad Carlos III de Madrid; and Telemaco Melia, Alcatel Lucent Bell Labs

Traditional IP mobility procedures^[4] are based on functions residing in both the mobile terminal and the network. Recently, we have been assisting in a shift in IP mobility protocol design, mostly focusing on solutions that relocate mobility procedures from the mobile device to network components. This new approach, known as *Network-Based Localized Mobility Management* (NetLMM), allows conventional IP devices (for example, devices running standard protocol stacks) to roam freely across wireless stations belonging to the same local domain. This property is appealing from the operator's viewpoint because it allows service providers to enable mobility support without imposing requirements on the terminal side (for example, software and related configuration). For this purpose the *Internet Engineering Task Force* (IETF) has standardized *Proxy Mobile IPv6* (PMIPv6)^[1].

This article details the Proxy Mobile IPv6 protocol, providing a general overview and an exhaustive description of a few selected functions.

Why Network-Based Localized Mobility?

The ability to move while being connected to a communication network is very attractive for users, as demonstrated by the success of cellular networks. However, while designing the IP stack, mobility was not retained as a requirement and, as a consequence, IP does not natively support mobility. The reason is a very basic design choice adopted in IP, both in IPv4^[2] and in IPv6^[3], namely that addresses have two roles: they are used as locators and identifiers at the same time.^[16]

IP addresses are *locators* that specify, by means of the routing system, how to reach the node (more properly, the *network interface*) that is using a specific destination address. The routing system keeps information about how to reach different sets of addresses that have a common network prefix, thus improving scalability of the system itself. However, IP addresses are also *identifiers* used by upper-layer protocols (for example, the *Transmission Control Protocol* [TCP]) to identify the endpoints of a communication channel. Additionally, names of nodes are translated by the *Domain Name System* (DNS) to IP addresses (which, in that way, play the role of node identifiers).

The linking of these two roles (*locators* and *identifiers*) is appealing because name resolution of the peer with whom we want to communicate and location finding translate to the same problem (that is, no translation mechanism is needed). However, the negative side effect is that supporting mobility becomes difficult.

Mobility implies separating the identifier role from the location one. From the identification standpoint, the IP address of a node should never change, but from the location point of view the IP address should change each time the node moves, showing its current location within the routing hierarchy (that is, the IP subnet to which the node is currently attached).

The IETF has studied the problem of terminal mobility in IP networks for a long time. It has developed IP-layer solutions for both IPv4 (Mobile IPv4^{[4], [5]}) and IPv6 (Mobile IPv6^[6]), enabling the movement of terminals and providing transparent service continuity. These solutions, being IP-based, are independent of the Layer 2 technologies. They provide Mobile Nodes with a permanent address (the *Home Address* [HoA]) to be used as identifier, and a temporal address (the *Care-of Address* [CoA]) to be used as locator. The CoA changes in each IP subnet visited by the Mobile Node. An entity in the network, the *Home Agent*, binds both addresses with the help of signaling generated by the Mobile Node. The Home Agent serving a Mobile Node must be placed in the subnet where the Home Address of that Mobile Node is topologically correct (the home network).

Although Mobile IP enables a host to move (that is, change the point of attachment in an IP network) while keeping session continuity, this ability is not sufficient for true mobility. Enabling efficient hand-offs is an additional and critical requirement. Because the IP handoff latency is affected by the time required to exchange signaling between the Mobile Node and the Home Agent, a new family of solutions proposes to use a local Home Agent (that is, a Home Agent closer to the Mobile Node) to provide mobility in a local domain; that is, to provide localized mobility support. Changing the point of attachment within the local domain requires only signaling to the local Home Agent, allowing faster signaling messages exchange because it is limited within the local domain. This approach is attractive because users typically move in localized environments (for example, they commute between their living homes and their work places) that can be covered with localized domains. Examples of these types of solutions are “Regional Registrations for IPv4”^[7] or “Hierarchical Mobile IPv6 for IPv6”^[8]. Note that the term “localized” refers to a particular area from the point of view of the IP network topology, but depending on the access technology, geographically the area can be large, as happens when applying a localized mobility approach to cellular networks.

A common feature of Mobile IP and the localized mobility proposals mentioned previously is that all of them are *host-based*. Mobile Nodes must signal themselves to the network when their location changes and must update routing states in the Home Agent, in the local Home Agent, or in both. This situation also raises the problem of complex security configurations to authenticate those signaling exchanges and modifications of routing states.

Therefore, the IETF decided to work on a solution for NetLMM^[10, 11], compounding the advantages of a network-based approach with the benefits of localized mobility management strategies. In NetLMM the network provides mobility support, although the Mobile Node does not participate in IP mobility procedures. That is, network operators can provide mobility support without requiring additional software and complex security configuration in the Mobile Nodes. Thus the deployment of network-based mobility solutions is greatly facilitated. Moreover, the Mobile Node can implement any global mobility solution, because the localized one is transparent and independent from it.

There are several target scenarios for Network-Based Localized Mobility Management^[9]:

- Large campus networks with *Wireless Local-Area Network* (WLAN) access: Users move with IP standard devices (that is, no additional hardware or software is required) within a campus that provides WLAN access and mobility support.
- Advanced beyond-third-generation (3G) networks: Cellular operators have been important promoters in the development of the NetLMM solution in the IETF. *Universal Mobile Telecommunications System* (UMTS) and *General Packet Radio Service* (GPRS) networks use a proprietary network-based localized mobility mechanism to provide mobility support for user data traffic (typically IP). This mechanism is based on the GPRS Tunneling Protocol^[11], a special-purpose solution developed for *Third-Generation Partnership Project* (3GPP) networks that uses TCP/IP application layer tunnels. A standardized NetLMM protocol for the Internet has important advantages:
 - Reduced costs in network management and in equipment supporting the technology (because of economy of scale)
 - Easier extension of mobility support to other technologies
 - Easier integration with other networks
- Other more-complex scenarios involving network mobility, as in automotive scenarios^[12], could benefit from a NetLMM approach to support mobility.

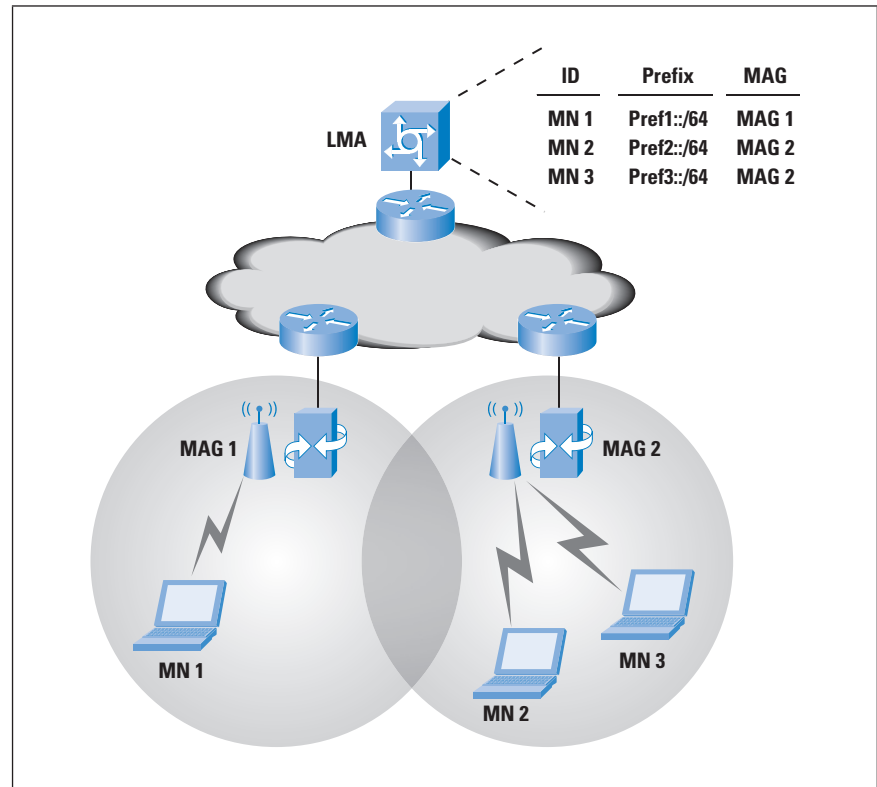
With these advantages in mind, the IETF has standardized a protocol to provide Network-Based Localized Mobility support in IP networks, the *Proxy Mobile IPv6* (PMIPv6) protocol.

Operation of Proxy Mobile IPv6

The main idea of PMIPv6 is that the mobile node is not involved in any IP layer mobility-related signaling. The Mobile Node is a conventional IP device (that is, it runs the standard protocol stack). The purpose of PMIPv6 is to provide mobility to IP devices without their involvement. This provision is achieved by relocating relevant functions for mobility management from the Mobile Node to the network.

PMIPv6 provides mobility support within a localized area, the *Localized Mobility Domain* (LMD) or PMIPv6 domain. While moving within the LMD, the Mobile Node keeps its IP address, and the network is in charge of tracking its location. PMIPv6 is based on *Mobile IPv6* (MIPv6), reusing the Home Agent concept but defining nodes in the network that must signal the changes in the location of a Mobile Node on its behalf.

Figure 1: Network Entities in Proxy Mobile IPv6

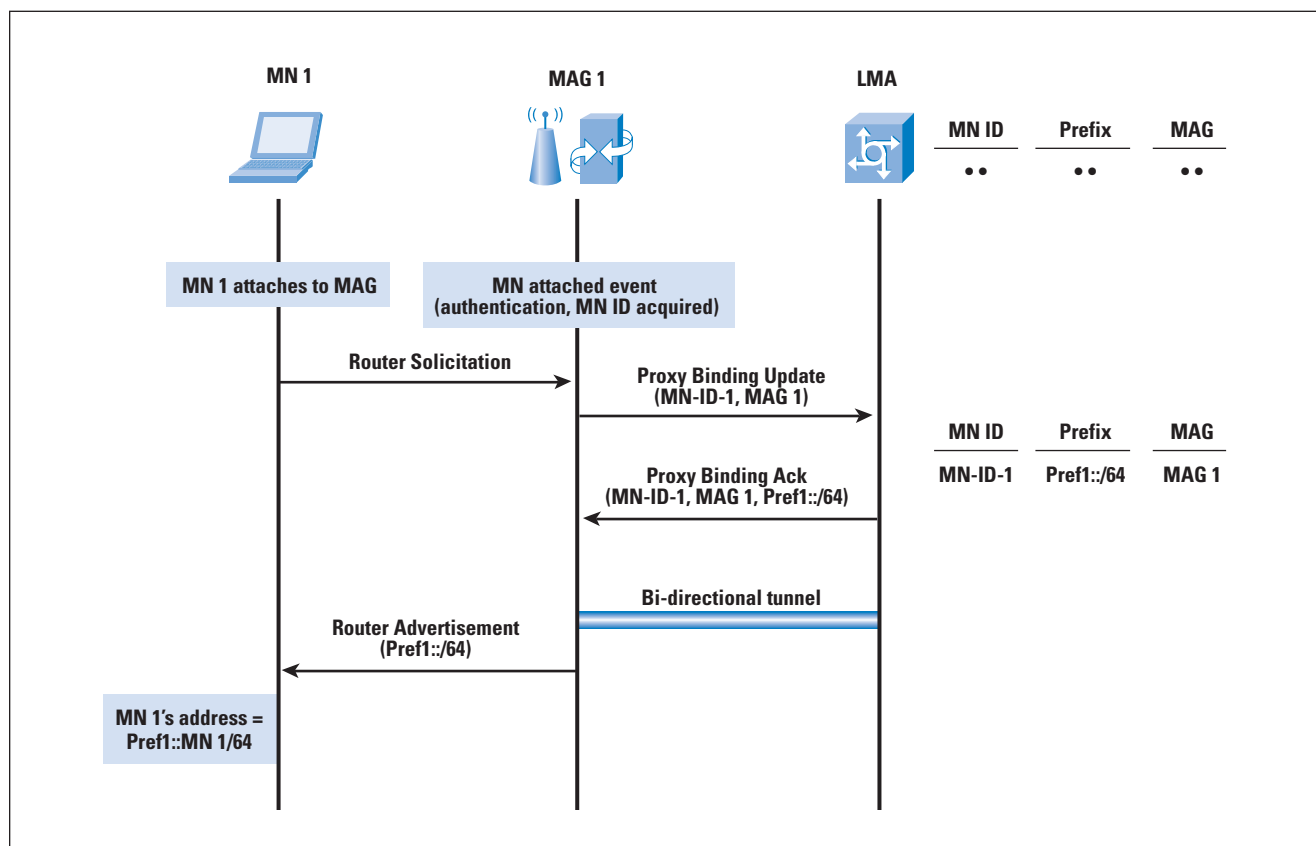


The functional entities in the PMIPv6 network architecture (refer to Figure 1) include the following:

- *Mobile Access Gateway* (MAG): This entity performs the mobility-related signaling on behalf of the Mobile Nodes attached to its access links. The MAG is usually the access router for the Mobile Node, that is, the first-hop router in the Localized Mobility Management infrastructure. It is responsible for tracking the movements of the Mobile Node in the LMD. An LMD has multiple MAGs.
- *Local Mobility Anchor* (LMA): This entity within the core network maintains a collection of routes for each Mobile Node connected to the LMD. The routes point to MAGs managing the links where the Mobile Nodes are currently located. Packets sent to or from the Mobile Node are routed through tunnels between the LMA and the corresponding MAG. The LMA is a topological anchor point for the addresses assigned to Mobile Nodes in the LMD, meaning that packets with those addresses as destination are routed to the LMA.

The basic operation of PMIPv6 follows. When a Mobile Node enters a PMIPv6 domain, it attaches to an access link provided by a MAG. The MAG proceeds to identify the Mobile Node, and checks if it is authorized to use the network-based mobility management service. If it is, the MAG performs mobility signaling on behalf of the Mobile Node (see in Figure 2 the signaling when the Mobile Node enters the PMIPv6 domain). The MAG sends to the LMA a *Proxy Binding Update* (PBU) associating its own address with the identity of the Mobile Node (for example, its *Media Access Control* [MAC] address or an identifier related to its authentication in the network). Upon receiving this request, the LMA allocates a prefix to the Mobile Node. Then the LMA sends to the MAG a *Proxy Binding Acknowledgment* (PBA) including the prefix allocated to the Mobile Node. It also creates a *Binding Cache* entry and establishes a bidirectional tunnel to the MAG. The MAG sends *Router Advertisement* messages to the Mobile Node, including the prefix allocated to the Mobile Node, so the Mobile Node can configure an address (stateless autoconfiguration). The Mobile Node can alternatively use stateful address autoconfiguration mechanisms. For simplicity, we assume in the rest of the article that the stateless address autoconfiguration mechanism is used, except when indicated otherwise.

Figure 2: Signaling When a Mobile Node Connects to the PMIPv6 Domain



Whenever the Mobile Node moves, the new MAG updates the location of the Mobile Node in the LMA and advertises the same prefix to the Mobile Node (through Router Advertisement messages), thereby making the IP mobility transparent to the Mobile Node. In this way the Mobile Node keeps the address configured when it first enters the LMD, even after changing its point of attachment within the network, and the LMD appears as a single link from the perspective of the Mobile Node. It should be noted that all the MAGs configure the same link local address for a specific Mobile Node. That is, the Mobile Node will never see a change in its default route configuration.

The bidirectional tunnel between the LMA and the MAG and associated routing states in both LMA and MAG manage the Mobile Node data plane. Downlink packets sent to the Mobile Node from outside of the LMD arrive to the LMA, which forwards them through the tunnel to the serving MAG. The MAG, after decapsulation, sends the packets to the Mobile Node directly through the access link. Uplink packets that originated in the Mobile Node are sent to the LMA from the MAG through the tunnel, and then are forwarded to the destination by the LMA. Traffic originated inside the LMD and directed to a Mobile Node also inside the LMD follows a similar procedure, going through two tunnels from the originating MAG, to the LMA, and then to the destination MAG. It should be noted that PMIPv6 allows a MAG to short-circuit the tunneling in case two mobile nodes directly communicate through any of its interfaces.

Protocol Details

We next describe the PMIPv6 primary functions. Because PMIPv6 is based on the Mobile IPv6 protocol format, we will highlight the differences and extensions to MIPv6. Readers interested in knowing all protocol details should refer to the RFC^[1].

Entering a PMIPv6 Domain

The Mobile Node enters the PMIPv6 domain by attaching to an access link. PMIPv6 defines a new functional entity, the MAG, typically residing in the access router. The MAG detects the attachment of the Mobile Node to the access link. The only access link types supported in PMIPv6 are point-to-point links; other types of links can be used as long as they are configured to emulate point-to-point links.

The MAG, upon detecting a Mobile Node attachment, verifies if the Mobile Node is eligible to the network-based mobility management service. Specific procedures to achieve this verification are out of the scope of the PMIPv6 standard. A Mobile Node that uses the mobility support service is identified by the network entities using a *Mobile Node Identifier* (MN-ID). The MN-ID must be stable and unique for the Mobile Node throughout the PMIPv6 domain, but the exact nature of this identifier is not specified. Possible examples are the Mobile Node MAC address or an identifier obtained as part of the Mobile Node authentication procedure.

After the MAG identifies the Mobile Node, authorizes its use of the NetLMM service, and acquires its Mobile Node Identifier, the MAG sends a PBU to the LMA; that is, it sends a registration request on behalf of the Mobile Node to the LMA. The PBU message is based on the MIPv6 *Binding Update* (BU) message with some extensions, but whereas the BU is sent by the Mobile Node, the PBU is sent by the MAG on behalf of the Mobile Node. A flag in the message is used to indicate that it is a PBU and not a BU. The PBU has as source address (and also in the alternate CoA option, if present) the global address configured in the egress interface of the MAG. This address is called *Proxy-CoA* in PMIPv6 terminology and is used by the LMA as locator of the Mobile Node. In the PBU, unlike in the BU, a Home Address destination option is not present; instead a *Mobile Node Identifier Option*^[13] has to be included with the Mobile Node Identifier, which is used to identify the Mobile Node throughout the PMIPv6 domain.

The PBU also contains additional information, such as the access link technology, a handoff indicator, the requested lifetime for the registration, and other optional data. The *handoff indicator* is a new mobility option defined in PMIPv6 that allows the MAG to signal the LMA whether the PBU originated upon network attachment or upon handover of a Mobile Node (if known by some unspecified mechanisms), and that information could be useful to support advanced functions such as multihoming. Examples of values of the handoff indicator include: a Mobile Node entering the PMIPv6 domain, a reregistration to update the registration lifetime, a handoff between MAGs, or a handoff between interfaces of the Mobile Node.

Upon sending the PBU, the MAG creates a Binding Update List entry^[6] for the Mobile Node. Note that this data structure in Mobile IPv6 is maintained by the Mobile Node to keep track of its bindings, but consequently to the PMIPv6 philosophy, the MAG maintains a *Binding Update List* (BUL) storing the bindings of the Mobile Nodes attached to it. The information in the Binding Update List allows the MAG to link the information about the Mobile Node, the interface in the MAG to which the Mobile Node is connected, and the LMA serving it, among others.

When the LMA receives the PBU sent by the MAG, it first checks that the message is correct according to the PMIPv6 specification, rejecting the registration otherwise. If the LMA accepts the PBU, it has to verify if its *Binding Cache* contains an entry for the Mobile Node identified in the PBU. When a Mobile Node first enters the PMIPv6 domain, the LMA cannot find an entry in its Binding Cache and has to create a new one. The Binding Cache entry is an extended version of the data structure defined for the Binding Cache entries in Mobile IPv6^[6].

The entry in the Binding Cache has a flag to indicate that it is a proxy registration, and it links all the information related to the Mobile Node, including its identification and the MAG serving it; that is, the location of the Mobile Node. If there is no entry for the Mobile Node in the Binding Cache (that is, the Mobile Node is entering the PMIPv6 domain), the LMA allocates one or more network prefixes to the Mobile Node. These prefixes are called *Home Network Prefixes*, and it must be noted that at least one network prefix is assigned per Mobile Node.

If the LMA cannot allocate a network prefix to a Mobile Node, it has to reject the registration. The address(es) that the Mobile Node uses while inside the PMIPv6 domain are configured from those Home Network prefixes. The decision of allocating one or more network prefixes depends on a global policy in the PMIPv6 domain or a per-Mobile Node policy. When the registration request is accepted, the LMA creates a *Binding Cache Entry* (BCE) with the accepted values for the registration, including the Mobile Node Identifier, the Proxy CoA (the address of the MAG serving the Mobile Node), and the Home Network prefix(es) allocated to the Mobile Node.

Upon BCE creation, the LMA creates an IPv6-in-IPv6 bidirectional tunnel, if one does not already exist, to the MAG sending the PBU. The LMA sets up forwarding routes through the tunnel for any traffic received that is addressed to the Home Network prefixes of the Mobile Node. Finally, the LMA creates a *Proxy Binding Acknowledgment* (PBA) and sends it to the corresponding MAG. The PBA message is based on the MIPv6 *Binding Acknowledgment* (BA) message with a few more extensions, including a flag that indicates that the message is a Proxy Binding Acknowledgement. The PBA informs the MAG about the registration request result, if it has been rejected (and why, using a status code) or accepted. The PBA contains the Mobile Node Identifier and the Home Network prefixes allocated to the Mobile Node. Unlike the Binding Acknowledgment, the PBA does not include a type 2 routing header (that in the Binding Acknowledgment includes the Home Address of the Mobile Node). Also the PBA is received and processed by the MAG, and not by the Mobile Node.

If the PBA confirms that the registration request has been accepted for the Mobile Node, the MAG creates an IPv6-in-IPv6 bidirectional tunnel, if one does not already exist, to the LMA. The MAG sets up forwarding routes, through the tunnel, for uplink or downlink packets received or sent from or to the Mobile Node. The MAG also updates the Binding Update List entry to reflect the accepted binding registration values.

Upon network attachment and during the PBU or PBA procedure, the Mobile Node can send a *Router Solicitation* in the access link as part of the standard neighbor discovery procedures. The MAG should not reply to this Router Solicitation until the registration in the LMA has been successfully completed. When the MAG receives the PBA indicating a successful registration, the MAG sends a Router Advertisement to the Mobile Node announcing the Home Network prefix(es). The Mobile Node can then apply the stateless address autoconfiguration mechanism or the stateful one (using the *Dynamic Host Configuration Protocol* [DHCP]) according to the indication in the Router Advertisement. For supporting DHCP, a DHCP relay agent has to be present in every MAG in the domain, and the relay agent must include in the link-address field of the *Relay Forward* message an IPv6 address from the Home Network prefix, to indicate to the DHCP server the range of addresses it can assign.

The PMIPv6 specification, as mentioned previously, supports only point-to-point access links with the Mobile Nodes. An interesting use case is to have a broadcast access link and to emulate point-to-point links with the Mobile Nodes to be able to apply the PMIPv6 specification. This case raises the problem of sending Router Advertisements that should be received only by the corresponding Mobile Node, and not by other Mobile Nodes present in the broadcast link. There are several ways to send these advertisements. The Router Advertisements could be sent to the IPv6 link-local address of the Mobile Node that the MAG can learn from the source address of router solicitations sent by the Mobile Node, or by some other unspecified means. Another possibility is to send Router Advertisements to the all-nodes multicast address at the IP layer but to the Link Layer 2 address of the Mobile Node.

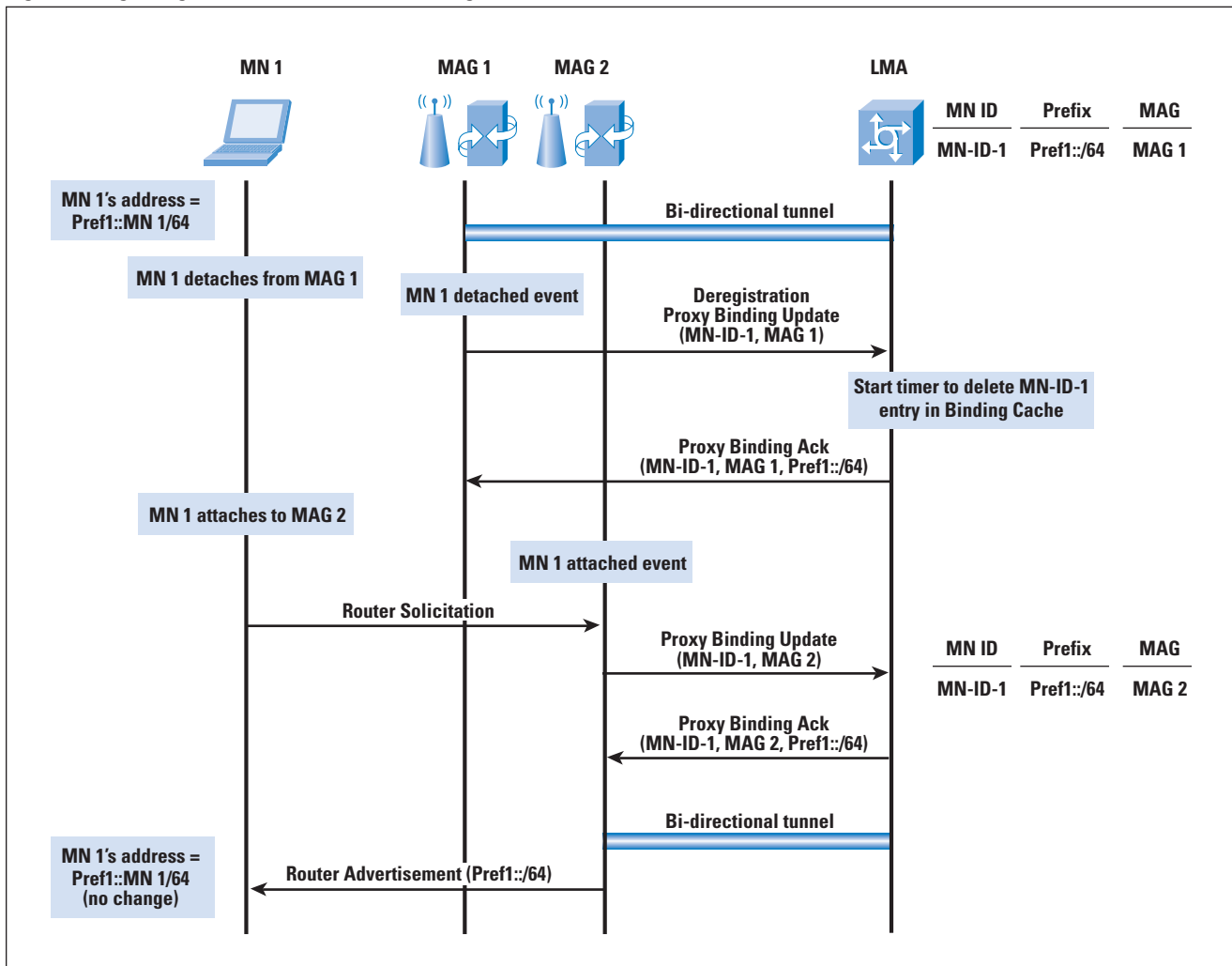
Changing MAG in a PMIPv6 Domain

The complete signaling for supporting the change of attachment by a Mobile Node in a PMIPv6 domain is described in Figure 3.

When a Mobile Node leaves a link, the event is detected by the corresponding MAG. The mechanism for Mobile Node movement detection is not specified in PMIPv6, but some possible options are link-layer events or an *IPv6 Neighbor Unreachability Detection* event. The MAG that detects that the Mobile Node has left the link must send a PBU with a Mobile Node de-registration request to the LMA. Upon receiving a PBA replying to the PBU or after a timer, the MAG deletes all the states associated with a specific Mobile Node.

When the LMA receives a PBU with a de-registration request for a Mobile Node with a valid entry in the Binding Cache, it sends the corresponding PBA and starts a timer. During the period defined by the timer the LMA drops any packets received for the Mobile Node. The use of this timer allows the LMA to receive a PBU from a new MAG updating the location of the Mobile Node. If the PBU is not received during that time, the LMA deletes the state associated with the Mobile Node.

Figure 3: Signaling When a Mobile Node Changes Point of Attachment



In a handoff situation the Mobile Node, after leaving a link, attaches to a new access link associated with a new MAG. The new MAG detects the Mobile Node and sends a PBU to the LMA on behalf of the Mobile Node. The LMA receives and processes the PBU, and detects that there is already a Binding Cache entry for that Mobile Node (the same Mobile Node Identifier). The LMA updates the Binding Cache entry with the new information, in particular with the Proxy CoA (egress IPv6 address) of the new MAG, updating also the tunnel and routing information for handling the traffic from or to the Mobile Node. The LMA sends a PBA to the new MAG in which it includes the Home Network prefix(es) already assigned to the Mobile Node. This scenario allows the new MAG to send a Router Advertisement with the same network prefix information as the Mobile Node received from the previous MAG. As stated before, the Mobile Node does not detect a link change and it keeps the same address(es). To make the change of link completely transparent to the Mobile Node, it must also continue receiving the Router Advertisements from the same link-local and link layer address; otherwise the Mobile Node would detect a change of default router. We describe how this problem is addressed in the next section.

Home Network Emulation and Address Uniqueness

MAGs must ensure that Mobile Nodes do not detect link changes when moving in a PMIPv6 domain; that is, MAGs must provide a home network emulation to the Mobile Nodes. To achieve this emulation, all the MAGs in the PMIPv6 domain must send, to a particular Mobile Node, Router Advertisements with the same network prefix information, as described previously. Additionally, the source IPv6 link-local address and the source link layer address in Router Advertisements sent to a Mobile Node must never change, independently of the MAG sending them. Therefore, the PMIPv6 specification requires all the MAGs to use, in any access link to which a particular Mobile Node attaches, the same link-local and link layer address.

PMIPv6 proposes two ways to meet this requirement:

- Configure a fixed link-local and link layer address to be used in all the access links in a PMIPv6 domain.
- Generate at the LMA the link-local address to be used by MAGs with a particular Mobile Node, and send it to the serving MAG through PMIPv6 signaling messages.

Both of these configuration methods are also helpful to guarantee address uniqueness in the access links of the PMIPv6 domain. The global addresses are always unique because all links are point-to-point and only one Mobile Node uses unicast global addresses over that link. Link-local addresses are used by the MAG and the Mobile Node on the link and a collision is possible. However, because the PMIPv6 specification requires that the link-local address used by the different MAGs with a particular Mobile Node is always the same while the Mobile Node moves across the PMIPv6 domain, the collision problem can happen only when the Mobile Node enters the PMIPv6 domain.

When a Mobile Node enters the domain, we must rely on *Duplicate Address Detection* (DAD) to detect a collision. If we use a globally unique link-local address for all the MAGs in the PMIPv6, then it is easy for the MAGs to respond to DAD requests from Mobile Nodes, because MAGs always know the address they must defend. If the link-local address to be used by the MAG with a Mobile Node is generated in the LMA, then it is desirable that the MAG learns that link-local address (that is, completes the PMIPv6 registration procedure) to defend it before the Mobile Node carries out the DAD procedure. You can ensure the MAG can learn this address by ensuring that the Layer 2 attachment is not completed until finishing the PMIPv6 signaling registration, or by configuring the PMIPv6 registration procedure in such a way that it is likely to be completed before the default waiting time of a DAD procedure.

Security Considerations

As with Mobile IPv6 signaling, PMIPv6 signaling is very sensitive to security threats, because it changes routing states of nodes in the network on behalf of the Mobile Nodes. PMIPv6 specification recommends using *IP Security* (IPsec) to protect the signaling exchanges between the MAGs and the LMA. A security association is needed between MAGs and the LMA, but how it is created is not defined. Two cases are possible:

- The network elements (LMA and MAGs) belong to the same operator.
- The elements belong to different operators with an agreement for roaming support.

In both scenarios, creating the security association is an affordable problem.

Traffic Handling in a PMIPv6 Domain

Traffic sent to any address belonging to a Home Network prefix is received by the LMA, the anchor point for those addresses. The LMA forwards the traffic through the tunnel to the MAG serving the Mobile Node, and the MAG decapsulates the packets and forwards them to the Mobile Node through the access link. Packets sent by the Mobile Node are forwarded by the MAG through the tunnel to the LMA. The LMA decapsulates the packets and forwards them to the destination. If a MAG has data traffic that originated in one of its access links and is destined to another of its access links, it can forward the traffic locally to avoid the forwarding through the LMA. This forwarding is done according to a policy configured in the MAG.

Performance Considerations

PMIPv6 presents two performance advantages compared with MIPv6. First, the LMA is a local network entity, so in principle the delay of sending signaling to the LMA will be lower than sending signaling to a remote Home Agent. And second, because the tunnel required to handle the traffic is terminated in the MAG instead of in the Mobile Node (as happens in MIPv6), we avoid the overhead of having a tunnel (two IP headers) over the radio interface. This overhead avoidance is relevant because bandwidth resources are scarcer over the air interface than in the backhaul network.

IPv4 Support Considerations

PMIPv6 acknowledges the existence of a dual-stack mobile host. To this end there are ongoing efforts to standardize IPv4 support for PMIPv6 operations. The extensions defined in [14] specify how to assign an IPv4 Home Address to a mobile host accessing the PMIPv6 domain. That is, the MAG—upon Mobile Node detection attachment and verification that the Mobile Node is eligible for PMIPv6 service—inserts in the PBU an “IPv4 Home Address Request Option.”

The LMA, upon reception of the PBU message, assigns an IPv6 *Home Network Prefix* (HNP) or an IPv4 Home Address by attaching an “IPv4 Home Address Reply Option” to the PBA. How the information is delivered to the Mobile Node depends on the interface between the Mobile Node and the MAG, possible examples being DHCP or *Internet Key Exchange Version 2* (IKEv2). The Mobile Node—independent of the method deployed—configures the HNP and the IPv4 Home address assigned by the LMA, thus supporting both IPv4- and IPv6-based applications.

Conclusions

PMIPv6 is a promising specification that allows network operators to provide localized mobility support without relying on mobility functions or configuration present in the mobile nodes. This reality greatly eases the deployment of the solution.

The IETF is currently working in the *Network-Based Mobility Extensions* (netext) Working Group on extending the PMIPv6 specification to add functions such as enhanced multihoming and intertechnology handoff support, and localized routing for traffic between MAGs to avoid going through the LMA. Additionally, the *Multicast Mobility* (multimob) Working Group is working on the support of multicast in PMIPv6.

References

- [1] S. Gundavelli (Ed.), K. Leung, V. Devarapalli, K. Chowdhury, and B. Patil, “Proxy Mobile IPv6,” RFC 5213, August 2008.
- [2] Jon Postel, “Internet Protocol,” RFC 791, September 1981.
- [3] Stephen E. Deering and Robert M. Hinden, “Internet Protocol, Version 6 (IPv6) Specification,” RFC 2460, December 1998.
- [4] William Stallings, “Mobile IP,” *The Internet Protocol Journal*, Volume 4, Number 2, June 2001.
- [5] Charles E. Perkins, “IP Mobility Support for IPv4,” RFC 3344, August 2002.
- [6] David B. Johnson, Charles E. Perkins, and Jari Arkko, “Mobility Support in IPv6,” RFC 3775, June 2004.
- [7] E. Fogelstroem, A. Jonsson, and C. Perkins, “Mobile IPv4 Regional Registration,” RFC 4857, June 2007.
- [8] H. Soliman, C. Castelluccia, K. ElMalki, and L. Bellier, “Hierarchical Mobile IPv6 (HMIPv6) Mobility Management,” RFC 5380, October 2008.
- [9] J. Kempf (Ed.), “Problem Statement for Network-Based Localized Mobility Management (NETLMM),” RFC 4830, April 2007.
- [10] J. Kempf (Ed.), “Goals for Network-Based Localized Mobility Management (NETLMM),” RFC 4831, April 2007.

- [11] 3GPP TS 29.060, “GPRS Tunnelling Protocol (GTP) across the Gn and Gp interface,” 2009. Available at: <http://www.3gpp.org/ftp/Specs/html-info/29060.htm>
- [12] Ignacio Soto, Carlos J. Bernardos, Maria Calderon, Albert Banchs, and Arturo Azcorra, “NEMO-Enabled Localized Mobility Support for Internet Access in Automotive Scenarios,” *IEEE Communications Magazine*, Vol. 47, No. 5, May 2009.
- [13] A. Patel, K. Leung, M. Khalil, H. Akhtar, and K. Chowdhury, “Mobile Node Identifier Option for Mobile IPv6 (MIPv6),” RFC 4283, November 2005.
- [14] R. Wakikawa and S. Gundavelli, “IPv4 Support for Proxy Mobile IPv6,” RFC 5844, May 2010.
- [15] Carlos J. Bernardos, Ignacio Soto, and María Calderón, “IPv6 Network Mobility,” *The Internet Protocol Journal*, Volume 10, Number 2, June 2007.
- [16] Dave Meyer, “The Locator Identifier Separation Protocol (LISP),” *The Internet Protocol Journal*, Volume 11, Number 1, March 2008.

IGNACIO SOTO received a telecommunication engineering degree in 1993, and a Ph.D. in telecommunications in 2000, both from the University of Vigo, Spain. He was a research and teaching assistant in telematics engineering at the University of Valladolid from 1993 to 1999. In 1999 he joined University Carlos III of Madrid, where he was an associate professor from 2001 until 2010. In 2010, he joined Universidad Politécnica de Madrid as associate professor. His research activities focus on mobility support in packet networks and heterogeneous wireless access networks. E-mail: isoto@dit.upm.es

CARLOS J. BERNARDOS received a telecommunication engineering degree in 2003, and a Ph.D. in telematics in 2006, both from the University Carlos III of Madrid, where he worked as a research and teaching assistant from 2003 to 2008, and since then as an associate professor. His Ph.D. thesis focused on route optimization for mobile networks in IPv6 heterogeneous environments. He has published more than 30 scientific papers in prestigious international journals and conferences, and he also contributes to the IETF. He served as TPC chair of WEEDEV 2009 and as guest editor of *IEEE Network*. E-mail: cjbc@it.uc3m.es

MARÍA CALDERÓN is an associate professor at the Telematics Engineering Department of University Carlos III of Madrid. She received a computer science engineering degree in 1991 and a Ph.D. degree in computer science in 1996, both from the Technical University of Madrid. She has published more than 40 papers in the fields of advanced communications, reliable multicast protocols, programmable networks, and IPv6 mobility. E-mail: maria@it.uc3m.es

TELEMACO MELIA received his Informatics Engineering degree in 2002 from the Polytechnic of Turin, Italy, and his Ph.D. in Mobile Communications from the University of Goettingen in April 2007. From June 2002 to December 2007 he worked at NEC Europe Ltd. in Heidelberg, Germany, in the Mobile Internet Group. He worked on IPv6-based Mobile Communication focusing on IP mobility support across heterogeneous networks and resource optimization control. In September 2008 he joined Alcatel Lucent Bell Labs. He is currently working on interworking architectures spanning 3GPP, WiMAX forum, and IETF standardization bodies. His main research interests include wireless networking and next-generation networks. He is author of more than 20 publications and he actively contributes to the IETF. E-mail: telemaco.melia@alcatel-lucent.com

Improving User Experiences with IPv6 and SCTP

by Dan Wing and Andrew Yourtchenko, Cisco Systems

To be successful, new technologies must improve the user experience. In the process of finding the best way to deploy a new technology, several approaches are typically conceived, written down, tried, and possibly discarded. This article addresses two such approaches for *Internet Protocol Version 6* (IPv6) and the *Stream Control Transmission Protocol* (SCTP)^[10].

Modern web browsers, web servers, and operating systems support IPv4 and IPv6, and several major content providers already support IPv6, including Google, NetFlix, and Facebook. However, their properties are not generally available over IPv6 because of a conflict between IPv6 technology and their business realities.

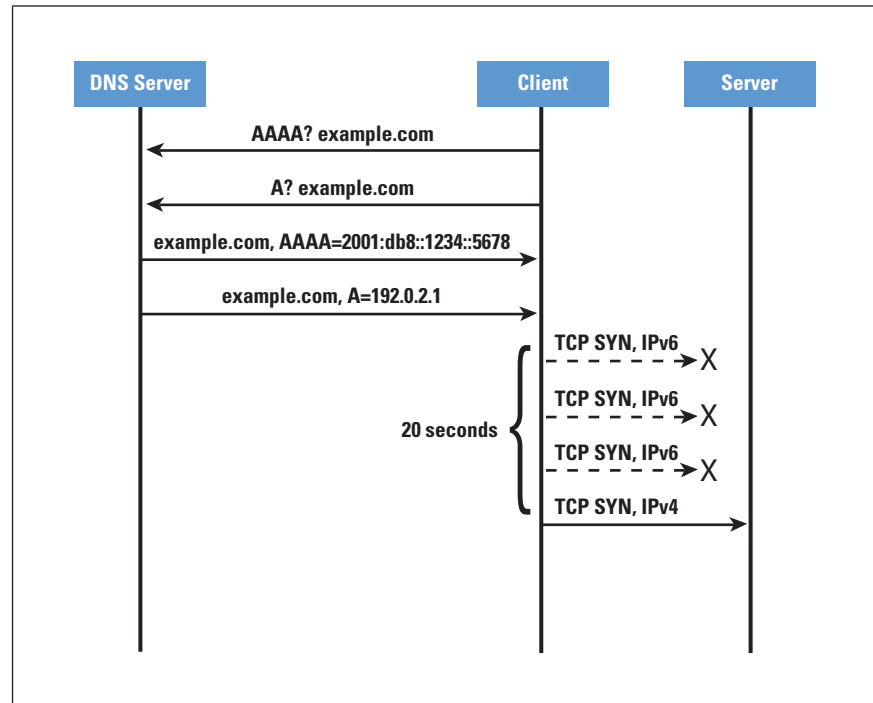
The technology in web browsers and operating systems involves doing *Domain Name System* (DNS) queries for AAAA and A resource records and then attempting to connect to the resulting IPv6 and IPv4 addresses *sequentially*. If the IPv6 path is broken (or slow), this connection can take a long time before it falls back to trying IPv4. This process is especially painful on typical websites that retrieve objects from different hosts—each failure incurs a delay. The combination of operating system and web browser results in delays from 20 seconds to several minutes if the IPv6 path is broken^[2]. The typical message flow of a TCP client is shown in Figure 1. Clearly, this delay is unacceptable to users. Users avoid this delay by disabling IPv6^[3] or avoiding IPv6-enabled websites.

The problem of broken IPv6 networks is relatively widespread^[6]. Providing content is a business—either directly (for example, streaming movies) or indirectly (for example, selling advertising). If users suffer delays viewing IPv6-enabled content (because of the technology reasons described previously), they will have an incentive to visit other websites. This scenario means lost revenue and is unacceptable to the business. Considering that all of the customers on today's Internet can reach IPv4 content, it is a business risk to enable IPv6 because some customers will suffer delays attempting to view IPv6 websites. Major content providers have been monitoring the situation and have published results^[7] showing that the IPv6 failure rate is too high to enable IPv6 AAAA for their content.

IPv6 problems have several causes. It is new technology, and monitoring of IPv6 connectivity is not yet on par with that of IPv4 because of single-point tunnels, unmanaged tunnels^[11], accidentally misconfigured firewalls, and router and link failures can more easily cause outages on IPv6. Many applications remain IPv4-only, or network administrators are relying on dual-stack equipment to transparently fail over to IPv4 during IPv6 outages.

However, such failover is never transparent to users—it takes many seconds or minutes! To avoid these problems, the content provider has only one choice: don't provide AAAA records if users might experience broken or slow IPv6.

Figure 1: Behavior of a Typical Web Browser



To work around that problem, Google implements a white list of DNS servers that it will provide AAAA records for^[8]. However, in its current incarnation, DNS white listing does not scale well because the *Internet Service Provider* (ISP) has to prove good IPv6 connectivity to Google, and then Google white lists the ISP's DNS servers to receive the AAAA records. The scaling problem is that there are thousands of ISPs around the world, and white listing and de-white listing them becomes a tiresome manual task for both ISPs and Google. Furthermore, if every content provider did DNS white listing, ISPs would have to work with several content providers in order to give value to the IPv6 network they have deployed to their subscribers! Content providers have started working together to consolidate requirements for DNS white listing and operate some sort of DNS white-listing service to slightly automate this process^[5].

Yet, DNS white listing still does not guarantee a working IPv6 network or a fast IPv6 network, because there is not a direct relationship between good IPv6 connectivity and the DNS server of a user's ISP. Even with the best of intentions and network design, there will still be instances where an IPv6 path or IPv4 path is working when the other path is broken. The result will be excessive delays for IPv4-only clients or dual-stack clients, depending on what sort of breakage occurs.

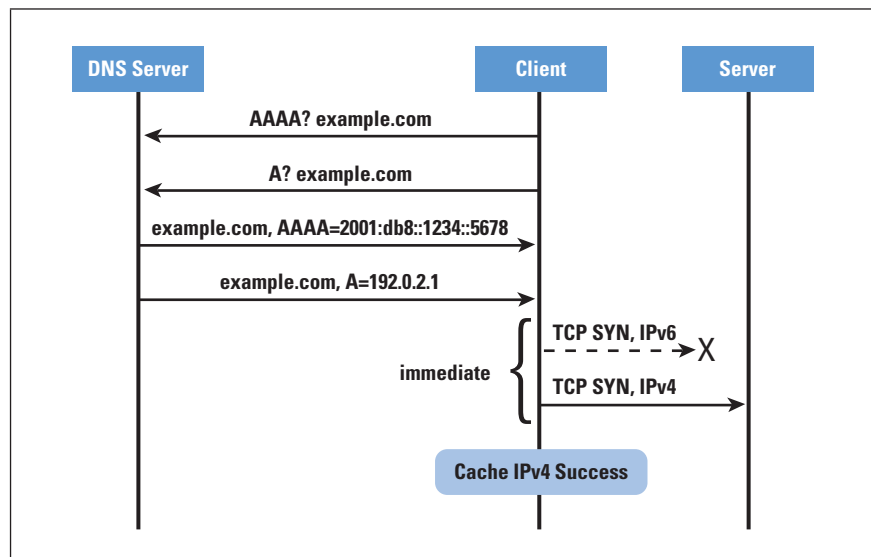
This situation contributes to the user perception that the Internet, or the particular website being accessed, is “down.” The user will visit a different site instead, possibly never returning to the site that was “down.”

Happy Eyeballs

A different approach solves these problems. In this approach, rather than an application slowly trying to make a connection on IPv6 and then on IPv4, the application makes its connection attempts more aggressively over both IPv6 and IPv4. Initially, the connection attempts are made *simultaneously* (rather than serialized), in order to provide a fast user experience.

The simultaneous connection attempts consume a little extra network bandwidth and twice the connection attempts on the server. To reduce that chatter, a cache is also maintained to store the success or failure of connecting using IPv6 or IPv4. We nickname this approach “Happy Eyeballs”^[1], because the “eyeballs” (users) are happier—their computer provides them immediate content, even if the network is suffering slow performance on IPv6 or IPv4 (Figure 2).

Figure 2: Dual-Stack Web Browser Implementing Happy Eyeballs



Obviously, sending a TCP SYN on both IPv6 and IPv4 doubles the number of connection attempts sent by the client. As discussed in [1], this chatter can be reduced by the application remembering if IPv6 (or IPv4) was successful in the previous connection attempt, and using that information for subsequent connection attempts. The sophistication of this cache is dependent on the memory (or disk) available, but even simple caching can be quite effective. When connecting to a new network (*third generation* [3G], different Wi-Fi network, or physical Ethernet), the connectivity of that new network can be determined and the cache of success or failure entirely or partially flushed, as necessary.

Thus, the doubling of connection attempts occurs only when connecting to a new network. Thereafter, initial connection attempts are delayed so that IPv6 (or IPv4) is tried first. But in all cases, significant user-noticeable delays are avoided when the IPv6 (or IPv4) is broken. The goal of Happy Eyeballs is to keep IPv6 enabled; that is, to make users unaware of IPv6 outages, so the user still visits IPv6-enabled websites without suffering any delay.

In this way, the user experiences a smooth migration from IPv4 to IPv6, and when necessary the fallback to IPv4 is almost immediate. This solution represents a significant improvement over today's web browsers. A drawback of this idea, however, is that it needs to be implemented in the application itself. Although it is a burden to upgrade those web browsers, there are only five major browsers^[9], and the browsers receive the immediate benefit of the aggressive probing. Browsers are also commonly upgraded already for faster *JavaScript* engines and other new features.

Another idea to determine if IPv6 is working is to *ping* or send another simple request to an IPv6 resource on the Internet, and disable IPv6 on the host if that IPv6 request fails. This approach interferes with IPv6 traffic within the enterprise (which may be working fine, whereas IPv6 to the Internet is broken), and disabling IPv6 would break IPv6 features deployed in OSs (for example, *DirectAccess* in Windows or *Back to My Mac* in Mac OS X). An advantage of this approach is that if IPv6 is disabled, no application suffers the IPv6 outage and associated delay to fall back to IPv4.

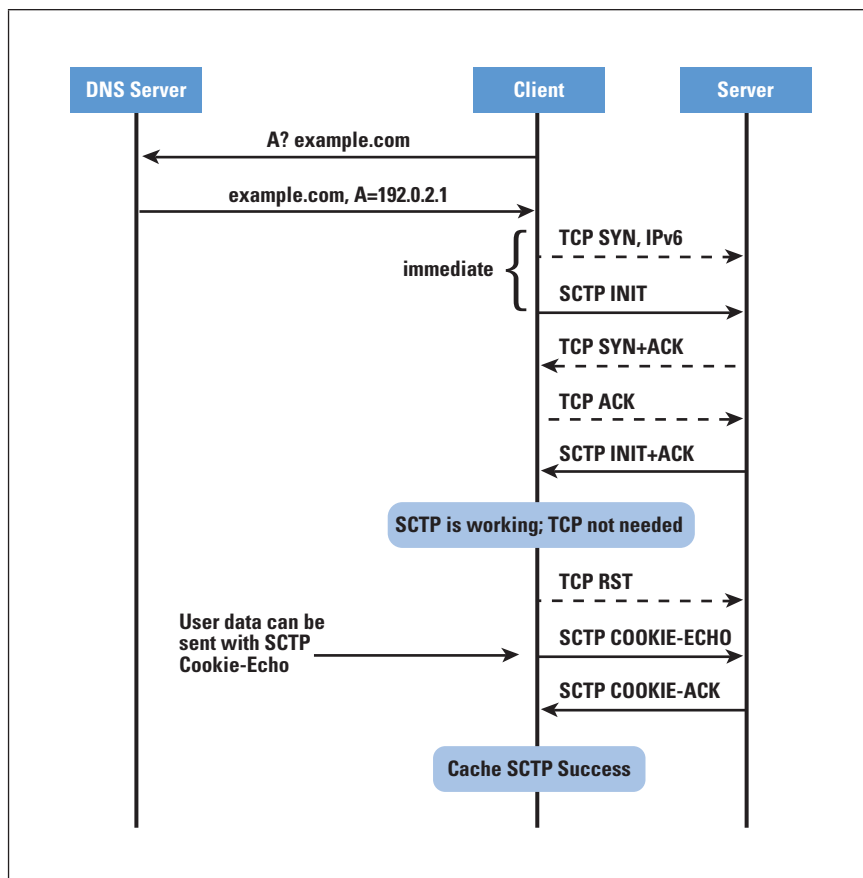
New Transport: SCTP

Besides the problem of network layer protocol selection, a similar task can be performed at the transport layer. Maybe surprisingly, one more transport protocol exists besides TCP, namely *Stream Control Transmission Protocol* (SCTP). SCTP provides significant advantages over TCP, and it was designed with some of the lessons learned by TCP implementations and deployment^[4] in mind.

Unlike IPv6 and IPv4, which have different DNS resource records (AAAA and A), we don't have a resource record to indicate that an application could, or should, use a different transport protocol. But even if we could indicate support for SCTP in DNS, the path might block it, reducing the usefulness of a DNS resource record. The path could be blocked by a NAT or firewall that expects only TCP or *User Datagram Protocol* (UDP).

Happy Eyeballs also describes a technique where a client can simultaneously try connecting using both TCP and SCTP. By necessity, this attempt is done entirely in the application, and the application would prefer the transport that responded faster and cache that information to reduce network chatter for subsequent connections to that server. This scenario is shown in Figure 3.

Figure 3: Client Implementing Happy Eyeballs for TCP/SCTP Selection



By combining the IPv6/IPv4 technique with the SCTP/TCP technique, a web browser running on a computer connected to a new dual-stack network sends four packets—an IPv4 TCP SYN, an IPv6 TCP SYN, an IPv4 SCTP INIT, and an IPv6 SCTP INIT. Based on the responses, it decides which transport protocol and which address family (IPv6 or IPv4) it prefers, and abandons the other connections. As described previously, connection information is cached for subsequent use to avoid consuming network bandwidth and server resources for subsequent network connections.

Conclusion

New technology aimed at improving user experience will be successful only if it meets expectations—an improved user experience. Because many companies are deriving all of their revenue from the Internet, any reduction in service means a loss of revenue. Thus, deploying new technology must not negatively affect the user experience. This article described one of the mechanisms that implementers can use to avoid negative effects on the user experience.

References

- [1] Dan Wing, Andrew Yourtchenko, and Preethi Natarajan, “Happy Eyeballs: Trending Towards Success (IPv6 and SCTP),” Internet-Draft, Work-in-Progress, July 2009:
<http://tools.ietf.org/html/draft-wing-http-new-tech>
- [2] “Broken IPv6 clients,” Lorenzo Colitti, June 2010:
<https://sites.google.com/site/ipv6implementors/2010/agenda>
- [3] “Google Trends”: <http://www.google.com/trends?q=enable+ipv6%2C+disable+ipv6>
- [4] P. Natarajan, “Leveraging Innovative Transport Layer Services for Improved Application Performance,” February 2009:
<http://www.cis.udel.edu/~amer/PEL/poc/pdf/NatarajanPhDdissertation.pdf>
- [5] Carolyn Duffy Marsan, “Google, Microsoft, Netflix in talks to create shared list of IPv6 users,” *Network World*, March 2010:
<http://www.networkworld.com/news/2010/032610-dns-ipv6-whitelist.html>
- [6] Tore Anderson, “IPv6 brokenness experiment, November results,” November 2009: <http://lists.cluonet.de/piper-mail/ipv6-ops/2009-December/002707.html>
- [7] Igor Gashinsky, “IPv6 & recursive resolvers: How do we make the transition less painful?” March 2010: <http://www.ietf.org/proceedings/77/slides/dnsop-7.pdf>
- [8] “Access Google services over IPv6”:
<http://www.google.com/intl/en/ipv6>
- [9] “Usage share of web browsers”: http://en.wikipedia.org/wiki/Usage_share_of_web_browsers
- [10] R. Stewart, Ed., “Stream Control Transmission Protocol,” RFC 4960, September 2007.
- [11] Gunter Van de Velde, Ole Troan, and Tim Chown, “Non-Managed IPv6 Tunnels considered Harmful,” July 2009:
<http://tools.ietf.org/html/draft-vandavelde-v6ops-harmful-tunnels>

DAN WING has a B.S. in Computer Science from Central Washington University and has co-chaired the IETF’s BEHAVE working group since 2006. He is a Distinguished Engineer at Cisco Systems, where he works on IPv6 transition technologies and has 30 patents issued or pending. E-mail: dwing@cisco.com

ANDREW YOURTCHENKO is a graduate of St. Petersburg Technical University in Russia, and has been in the networking industry since 1995. He is a Technical Leader at Cisco Systems in the network security area, and at IETF Andrew participates in the areas of security, TCP protocol, and IPv6 transition. E-mail: ayourtch@cisco.com

Letter to the Editor

In response to “NAT++: Address Sharing in IPv4,” in *The Internet Protocol Journal*, Volume 13, No. 2, June 2010:

Excellent article Geoff, so good I read it twice. While reading your article I was reminded of a recent experience that falls in the category of “unintended consequences.” Since one of your situation descriptions was similar to the one I’m in, I thought I would relay my circumstance and experience and see if I can make my point.

A couple of months ago I signed up for an IPTV trial with my provider, and it was installed with a minimum of effort. The service is based on Cisco *Dial-on-Demand Routing* (DDR) and, of course, DSL service.

It worked fine for a couple of days; video feeds were good and all my computers and server worked just as before on a wireless network within my home. Then one day it appeared that I had lost *Domain Name System* (DNS) service, because I couldn’t get name resolution to work but could route using the raw IPv4 addresses. So, I placed a trouble ticket and, of course, the provider’s first request was to cold boot the DDR device and everything in the house, which I did. Sure enough, upon bringing all components back up (except one), everything was fine.

A day or two later I had to print something and powered on my HP 6510 wireless printer, printed what I needed to print, and then discovered I had lost DNS service again. I placed a trouble call and my provider came out and replaced the DDR device I went through the cold boot process (except one device) and everything was OK until I brought the printer online and the trouble returned. By now I had this nagging memory that wouldn’t surface; something about that printer... With the printer powered off I rebooted the DDR, fired up SharkWire, and everything looked and worked OK.

Then I powered up the HP printer and saw another nagging memory; it immediately performed an *Address Resolution Protocol* (ARP) broadcast of the v4 address **169.254.65.206**—the famous black-hole address from RFC 3927^[1]. Immediately after the ARP broadcast, the printer put out the normal *Dynamic Host Configuration Protocol* (DHCP) request and was assigned one from the *Network Address Translation* (NAT) pool.

That’s when I stepped back from looking at the “trees” and gazed upon the “forest” and realized, with some embarrassment, that the public side (access side) was using a single IPv4 address with *Port Address Translation* (PAT) so the DDR box was blocking all the outbound PAT addresses attached to the single IPv4 address. I wrote down the details and e-mailed them to my provider, and had revised code pushed to the DDR the next day. Problem fixed.

All of this discussion leads me to ponder about other situations of “hard codes” in the network, either RFC-based or circumstance-based, that will falter with a switch to IPv6. Not in the core but in the customer networks. These unintended consequences could be many. Does HP run a dual stack for IPv4 and IPv6? I doubt it.

How can we get customers and vendors thinking about possible long-ago workarounds that they may have hard coded using IPv4? Any other RFCs out there like 3927? (It used to be easy when there were only a few hundred RFCs.) That could be the most expensive portion of the transition, verifying code ...

Keep up the good work; your articles make me think a lot and I really enjoy them. And, yes, I do use them for reference quite often.

Regards,

—Paul Dover
pdover@centeriem.com

- [1] S. Cheshire, B. Aboba, and E. Guttman, “Dynamic Configuration of IPv4 Link-Local Addresses,” RFC 3927, May 2005.

The author responds:

Thank you Paul for this anecdote and the important lesson behind it. Over some 30 years of intense development we’ve managed to accumulate a sizeable volume of technical specifications. Indeed, in October 2010 the RFC Editor published RFC 6068, and I’m not sure that any individual could claim a deep familiarity with every one of them, let alone claim to have a good understanding of their potential interaction. So when we look at various transitional technologies to sustain this industry through the next few years of attempting to support a comprehensive dual stack network in the face of the forthcoming hiatus of supply of IPv4 addresses, it should not come as a surprise when some devices or configurations fail in strange and unexpected ways, simply because they adhere to a technical standard that perhaps we’ve lost sight of in the flurry of generating new transitional technologies.

—Geoff Huston
gih@apnic.net

Dr. Jianping Wu Receives Postel Award

The *Internet Society* (ISOC) recently awarded its prestigious *Jonathan B. Postel Service Award* for 2010 to leading Chinese technologist Dr. Jianping Wu for the pioneering role he has played in advancing Internet technology, deployment, and education in China and Asia Pacific over the last twenty years.

Dr. Wu's best-known contribution is the development of the *China Education and Research Network* (CERNET) which he designed and developed to be the first Internet backbone network in China. Created to establish a nation-wide advanced network infrastructure to support education and research among universities, CERNET has since become the world's largest national academic network. Since 1998, Dr. Wu has also devoted his time to the design and development of a large-scale native IPv6 backbone in China that now serves to connect over 200 universities and millions of users.

The Postel Award was established by the ISOC to honour individuals or organisations that, like Jon Postel, have made outstanding contributions in service to the data communications community. Commenting on its presentation to Dr. Wu, Lynn St. Amour, President and CEO of ISOC said: "Jianping Wu has dedicated his career in China to developing a broadly accessible Internet that brings people together. Twenty years ago, Dr. Wu recognized the importance and future impact of the Internet and the pivotal role it would play in terms of its impact on social reform, technology advancement and economic growth for China. He has worked tirelessly to bring his vision to life. As a result, the networks that resulted from his determination and hard work have played an important role in driving Internet development in China and have had a significant impact on the Internet worldwide."

ISOC presented the award, including a US\$20,000 honorarium and a crystal engraved globe, during the 78th meeting of the *Internet Engineering Task Force* (IETF) in Maastricht, The Netherlands 25–30 July 2010.

DNSSEC Deployed in the Root Zone

On July 16, 2010 the U.S. Department of Commerce's *National Telecommunications and Information Administration* (NTIA) and the *National Institute of Standards and Technology* (NIST) announced the completion of an initiative with the *Internet Corporation for Assigned Names and Numbers* (ICANN) and VeriSign to enhance the security and stability of the Internet.

The announcement marks full deployment of a security technology—*Domain Name System Security Extensions* (DNSSEC)^[1]—at the Internet’s authoritative root zone, which will help protect Internet users against cache poisoning and other related cyber attacks.

“The Internet plays an increasingly vital role in daily life, from helping businesses expand to improving education and health care,” said Assistant Secretary for Communications and Information and NTIA Administrator Lawrence E. Strickling. “The growth of the Internet is due in part to the trust of its users—trust, for example, that when they type a website address, they will be directed to their intended website. Today’s action will help preserve that trust. It is an important milestone in the ongoing effort to increase Internet security and build a safer online environment for users.”

“Improving the trustworthiness, robustness and scaling of the Internet’s core infrastructure is an activity that lines up strongly with NIST’s mission, and we have been contributing to design, standardization and deployment of DNSSEC technology for several years,” said NIST Director Patrick Gallagher. “The deployment of DNSSEC at the root zone is the linchpin to facilitating its deployment throughout the world and enabling the current domain-name system to evolve into a significant new trust infrastructure for the Internet.”

The *Domain Name System* (DNS) is a critical component of the Internet infrastructure. The DNS associates user-friendly domain names (for example, www.commerce.gov) with the numeric network addresses (for example, 170.110.225.168) required to deliver information on the Internet, making the Internet easier for the public to navigate. The authenticity of the DNS data is essential to Internet use. For example, it is vital that users reach their intended destinations on the Internet and are not unknowingly redirected to bogus and malicious websites.

The DNS was not originally designed with strong security mechanisms, and technological advances have made it easier to exploit vulnerabilities in the DNS protocol that put the integrity of DNS data at risk. Many of these vulnerabilities are mitigated by the deployment of DNSSEC, which is a suite of *Internet Engineering Task Force* (IETF) specifications for securing information provided by the DNS.

A main goal of this action—DNSSEC deployment at the root zone—is to facilitate greater DNSSEC deployment throughout the rest of the global DNS hierarchy. While deployment of DNSSEC will protect Internet users from certain DNS-related cyber attacks, users must continue to exercise vigilance in protecting their information online.

ISOC Embraces DNSSEC

The *Internet Society* (ISOC) recently announced that it has deployed DNSSEC, a set of extensions to the DNS that provides a level of assurance, for its **isoc.org** domain. The announcement builds on an announcement by the *Public Interest Registry* (PIR) that they have implemented DNSSEC for the entire **.org** top-level domain.

“We are pleased to be among the first organisations in the **.org** top level domain to deploy DNSSEC, as DNSSEC provides an important building block for increasing user confidence in the Internet,” said Lynn St.Amour, President and CEO of the Internet Society. “Implementing DNSSEC for the **.org** top-level domain is an important step in ensuring the global Internet serves as a trusted channel for communication and collaboration and we applaud the PIR’s efforts in this area.”

“DNSSEC acts like tamper-proof packaging to make sure that when you type in the website name of your bank you actually get the server IP address your bank wants you to use,” said Leslie Daigle, Chief Internet Technology Officer of ISOC. “In this way, DNSSEC allows us to have more confidence in the online activities that are increasingly becoming a part of our lives at work, home, and school.”

DNSSEC technology used today is the result of careful protocol engineering and standardization within the IETF; implementation by various DNS vendors; and operational trials by DNS operators. In addition to **.org**, DNSSEC is currently implemented by several country-specific top-level domains: Brazil (**.br**), Bulgaria (**.bg**), The Czech Republic (**.cz**), Puerto Rico (**.pr**), and Sweden (**.se**).

ISOC is a non-profit organisation founded in 1992 to provide leadership in Internet related standards, education, and policy. ISOC is the organisational home of the IETF. With offices in Washington, D.C., and Geneva, Switzerland, it is dedicated to ensuring the open development, evolution, and use of the Internet for the benefit of people throughout the world. For more information see: <http://isoc.org>

DNSSEC Fund Announced

In order to speed up the process of introduction a more secure global DNS infrastructure, the Netherlands-based charity *NLnet Foundation* has announced the creation of a global fund where open source projects can apply for grants to work on *Domain Name System Security Extensions* (DNSSEC) in their Internet applications.

DNSSEC is one of the key technologies for a safer Internet, as it allows the Internet user to know for sure that he or she is being sent to the right computer or service on the Internet. “If you type the name of your bank into a browser, you want to be sure that you are actually directed to a computer of that bank,” said Michiel Leenaars, Director of Strategy at NLnet foundation. “Domain names are vital to the way we use the Internet, and without DNSSEC users are open to serious abuse.”

DNSSEC provides a cryptographic seal of authenticity that gives real proof of the validity of the domain name you use when you visit a website, chat or send an e-mail. With DNSSEC you get a *chain of trust* from the root of the Internet to the service you want to connect to—opening the way for many new exciting opportunities for humans and computers to exchange information safely. DNSSEC is being gradually introduced worldwide.

The new fund will provide grants for reengineering important software to reliably work with DNSSEC. “The signing of the root through DNSSEC is a historical moment, but in a way it is only the beginning,” said Leslie Daigle, Chief Internet Technology Office at the Internet Society. “Actual users will not fully benefit from protection in the more challenging situations as long as DNSSEC does not reach them.” A great deal of work has already been done at the infrastructure level—most DNS servers such as *BIND*, *NSD* and *Unbound* now support the new technology. However, it will take a lot of work at the user level as well: operating systems, web browsers, e-mail servers, VoIP clients, and many other pieces of software need to be able to reliably work with DNSSEC.

“Every Internet user deserves to be protected by DNSSEC, yet currently almost no end user software is ready to take full advantage of the availability of DNSSEC,” said Leenaars. “The IT community has a big responsibility in making sure that DNSSEC gets deployed across the board swiftly. We aim to accelerate the process significantly by putting some money on the table, and we invite other stakeholders to join us.”

Since there are many applications and platforms that will require work, the NLnet Foundation is very open to cooperation with others as well as to targeted donations from interested stakeholders such as governments, registries and corporations.

The NLnet Foundation is a registered Netherlands charity with a long history of supporting Internet standardization. The foundation gained its capital from selling the first Dutch Internet Service Provider.

Potential applicants and collaborators can find more information at:
<http://nlnet.nl/dnssec>

See also:

- [1] Miek Gieben, “DNSSEC: The Protocol, Deployment, and a Bit of Development,” *The Internet Protocol Journal*, Volume 7, No. 2, June 2004.
- [2] Torbjörn Eklöv, and Stephan Lagerholm, “Operational Challenges when Implementing DNSSEC,” *The Internet Protocol Journal*, Volume 13, No. 2, June 2010.
- [3] <http://www.dnssec.net/>

Call for Papers: Internet Privacy Workshop

The *Internet Architecture Board* (IAB), *World Wide Web Consortium* (W3C), *Internet Society* (ISOC) and the *Massachusetts Institute of Technology* (MIT) will hold a joint *Internet Privacy Workshop* on December 8 and 9, 2010 at MIT, Cambridge, Massachusetts on the question:

“How Can Technology Help to Improve Privacy on the Internet?”

Information about who we are, what we own, what we have experienced, how we behave, where we are located, and how we can be reached are among the most personal pieces of information about us. This information is increasingly being made more easily available electronically via the Internet, often without the consent of the subject. The question for the workshop therefore is: How can we ensure that architectures and technologies for the Internet, including the World Wide Web, are developed in ways that respects users’ intentions about their privacy?

This workshop aims to explore the experience and approaches taken by developers of Internet including Web technology, when designing privacy into these protocols and architectures. Engineers know that many design considerations need to be taken into account when developing solutions. Balancing between the conflicting goals of openness, privacy, economics, and security is often difficult, as illustrated by Clark, et al. in “Tussle in Cyberspace: Defining Tomorrow’s Internet,” see:

<http://groups.csail.mit.edu/ana/Publications/PubPDFs/Tussle2002.pdf>

As a member of the technical community, we invite you to share your experiences by participating in this important workshop. Workshop participants will focus on the core privacy challenges, the approaches taken to deal with them, and the status of the work in the field. The objective is to draw a relationship with other application areas and other privacy work in an effort to discuss how specific approaches can be generalized.

Interested parties must submit a brief contribution describing their work or approach as it relates to the workshop theme. We welcome visionary ideas for how to tackle Internet privacy problems, as well as write-ups of existing concepts, deployed technologies, and lessons-learned from successful or failed attempts at deploying privacy technologies. Contributions are not required to be original in content.

Submitters of accepted position papers will be invited to the workshop. The workshop will be structured as a series of working sessions, punctuated by invited speakers, who will present relevant background information or controversial ideas that will motivate participants to reach a deeper understanding of the subject.

The organizing committee may ask submitters of particularly topical papers to present their ideas and experiences to the workshop. We will publish submitted position papers and slides together with a summary report of the workshop. There are no plans for any remote participation in this workshop.

To be invited to the workshop, please submit position papers to privacy@iab.org by November 5, 2010. More detailed information about the workshop, including further details about the position paper requirements, is available at:

<http://www.iab.org/about/workshops/privacy/>

We look forward to your input,

Bernard Aboba (IAB)

Daniel Appelquist (W3C)

Jon Peterson (IAB)

Karen Sollins (MIT)

Trent Adams (ISOC)

Karen O'Donoghue (ISOC)

Thomas Roessler (W3C)

Hannes Tschofenig (IAB)

Organizations Urged to Stop Delaying IPv6 Deployment

The *Number Resource Organization* (NRO), the official representative of the five *Regional Internet Registries* (RIRs) that oversee the allocation of all Internet number resources, recently unveiled the findings of a global, independent survey into organizations' IPv6 readiness. Funded by the European Commission and conducted by GNKS Consult and TNO, the study reveals that the majority of organizations are taking steps toward IPv6 deployment, as the IPv4 address pool continues to deplete rapidly.

IP addresses are critical for the operation of the Internet. Every Internet-enabled device needs an IP address to connect to the rest of the network. The biggest threat facing the Internet today is that less than 6% of the current form of IP addresses, IPv4, remains and the pool is likely to be completely depleted next year. This means that organizations need to adopt IPv6, the next-generation addressing protocol. There is a far larger pool of IPv6 addresses, allowing for more devices to connect to the Internet and helping to safeguard the sustainable growth of the Internet.

The survey, which polled over 1,500 organizations from 140 countries, highlights that organizations are increasingly aware of the need to deploy IPv6: approximately 84% already have IPv6 addresses or have considered requesting them from the RIRs. Only 16% of respondents have no plans to deploy IPv6 addresses.

The study also demonstrates that there are some misconceptions around the cost of adopting IPv6. Over half of all respondents noted that the cost of deployment was a major barrier for IPv6 adoption. While organizations might delay investing in IPv6, this may ultimately result in greater costs, with last-minute deployment and poor planning likely to increase the investment required.

Of the 84% of respondents that have requested IPv6 addresses or have considered doing so, three-quarters reported the need to stay ahead of competition as the main reason for IPv6 adoption. Half of these respondents also noted that a lack of available IPv4 space was a major driver for deployment. When asked about issues they had encountered when deploying IPv6:

- 60% cited the lack of vendor support as a major barrier for deployment. However, most of the latest hardware and software support IPv6. The RIRs are strongly urging organizations to check with their suppliers to ensure that the technologies they use are IPv6 compatible.
- 45% reported a struggle to find knowledgeable technical staff to support deployment. However, all five RIRs arrange technical training to facilitate an efficient IPv6 deployment, details of which can be accessed via the NRO website.

Fifty-eight percent of all organizations polled were ISPs. It is likely that respondents to this survey are further ahead in IPv6 deployment than ISPs overall, but all organizations should ensure that their ISP offers or plans to offer services over IPv6. Out of the polled ISPs:

- Approximately 60% already offer, or plan to offer within the next year, IPv6 to consumers.
- 70% already offer, or plan to offer within the next year, IPv6 to businesses.
- Only about 10% of polled ISPs have no plans to offer IPv6 to consumers or businesses.

Axel Pawlik, Chairman of the NRO, commented: “It’s great to see that as we move toward complete IPv4 exhaustion, more organizations worldwide are waking up to the need to adopt IPv6 and are sourcing IPv6 addresses from the RIRs.”

“Yet there is still a distinct lack of Internet traffic over the next addressing protocol, with not enough ISPs offering IPv6 services and 30% of ISPs saying the proportion of this traffic is less than 0.5%. It’s critical that ISPs now take the next step in the global adoption effort by offering IPv6 services to their customers to help boost traffic over IPv6.”

Per Blixt, Head of Unit in the Information Society and Medias at the European Commission, said:

“It’s encouraging to see that so many organizations have made IPv6 adoption their priority. Still, as the Internet becomes increasingly important for global socio-economic development, it’s critical that those who are still sitting on the fence act now on IPv6. Only by ensuring that all organizations adopt IPv6 can we ensure the sustainable growth of the digital economy worldwide.”

This survey is a follow-up to a study conducted in 2009 amongst organizations in Europe, Middle East and parts of Central Asia, as well as Asia Pacific; however this year's survey polled organizations worldwide. The full research report is available at:

<http://www.nro.net/documents/GlobalIPv6SurveySummaryv2.pdf>

The NRO exists to protect the pool of unallocated Internet numbers (IP addresses and AS numbers) and serves as a coordinating mechanism for the five RIRs to act collectively on matters relating to the interests of RIRs. For further information, visit <http://www.nro.net>

The RIRs are independent, not-for-profit membership organizations that support the infrastructure of the Internet through technical coordination. There are five RIRs in the world today. Currently, the *Internet Assigned Numbers Association* (IANA) allocates blocks of IP addresses and ASNs, known collectively as *Internet Number Resources*, to the RIRs, who then distribute them to their members within their own specific service regions. RIR members include *Internet Service Providers* (ISPs), telecommunications organizations, large corporations, governments, academic institutions, and industry stakeholders, including end users

The RIR model of open, transparent participation has proven successful at responding to the rapidly changing Internet environment. Each RIR holds one to two open meetings per year, as well as facilitating online discussion by the community, to allow the open exchange of ideas from the technical community, the business sector, civil society, and government regulators. Each RIR performs a range of critical functions including: The reliable and stable allocation of Internet number resources (IPv4, IPv6 and *Autonomous System Number* resources); The responsible storage and maintenance of this registration data; The provision of an open, publicly accessible database where this data can be accessed. RIRs also provide a range of technical and coordination services for the Internet community. The five RIRs are:

AfriNIC: <http://www.afrinic.net>

APNIC: <http://www.apnic.net>

ARIN: <http://www.arin.net>

LACNIC: <http://www.lacnic.net>

RIPE NCC: <http://www.ripe.net>

This publication is distributed on an "as-is" basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2010 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

December 2010

Volume 13, Number 4

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Emergency Services.....	2
Integrating Core BGP/MPLS Networks	18
Letter to the Editor	32
Book Review.....	33
Fragments	37

FROM THE EDITOR

I have recently started using both a smartphone and a tablet device for Internet access. Like millions of other Internet users, I have discovered the wonders of mobile applications that provide everything from the traditional Internet services (e-mail and web browsing) to specialized software that can pinpoint my location on a map, provide live currency-exchange calculations, give weather forecasts, and my favorite: play radio stations from all over the world. I am old enough to remember the orange glow from pre-transistor vacuum-tube radios, so having a customizable “world radio” in the form of an “app” on a smartphone seems almost like science fiction.

But radio is not the only traditional service that is now available over the Internet. Another prominent example is telephony or *Voice over IP* (VoIP). Not only is VoIP replacing traditional land lines in many places, the original circuit-switched telephone network is itself increasingly using VoIP technology in place of an infrastructure of land lines and dedicated switching equipment. An important aspect of traditional phone service is the notion of special numbers for *emergency services*. Such systems rely on a database of phone numbers and addresses that allow emergency personnel to dispatch responders to the correct location. This location identification becomes a lot more complicated if the caller is using an Internet-based calling service rather than a hard-wired telephone. The IETF has been tackling this problem in the *Emergency Context Resolution with Internet Technology* (ECRIT) working group. Our first article, by Hannes Tschofenig and Henning Schulzrinne, is an overview of the architecture this working group is developing.

According to the ITU-T, a *Next Generation Network* (NGN) is “...a packet-based network which can provide services including Telecommunication Services and is able to make use of multiple broadband, Quality of Service-enabled transport technologies in which service-related functions are independent from underlying transport-related technologies.” Paul Veitch, Paul Hitchen, and Martin Mitchell describe the integration of a standalone core BGP/MPLS VPN network into an NGN architecture.

Please check your subscription expiration date and renew online if you wish to continue receiving this journal. Click the “Subscriber Services” link at www.cisco.com/ipj to get to the login page. If you need any assistance just send e-mail to ipj@cisco.com and we will make the necessary changes for you.

—Ole J. Jacobsen, Editor and Publisher

ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

Emergency Services for Internet Multimedia

by Hannes Tschofenig, Nokia Siemens Networks and Henning Schulzrinne, Columbia University

Summoning the police, the fire department, or an ambulance in emergencies is one of the most important functions the telephone enables. As telephone functions move from circuit-switched to Internet telephony, telephone users rightfully expect that this core feature will continue to be available and work as well as it has in the past. Users also expect to be able to reach emergency assistance using new communication devices and applications, such as instant messaging or *Short Message Service* (SMS), and new media, such as video. In all cases, the basic objective is the same: The person seeking help needs to be connected with the most appropriate *Public Safety Answering Point* (PSAP), where call takers dispatch assistance to the caller's location. PSAPs are responsible for a particular geographic region, which can be as small as a single university campus or as large as a country.

The transition to Internet-based emergency services introduces two major structural challenges. First, whereas traditional emergency calling imposed no requirements on end systems and was regulated at the national level, Internet-based emergency calling needs global standards, particularly for end systems. In the old *Public Switched Telephone Network* (PSTN), each caller used a single entity, the landline or mobile carrier, to obtain services. For Internet multimedia services, network-level transport and applications can be separated, with the *Internet Service Provider* (ISP) providing IP connectivity service, and a *Voice Service Provider* (VSP) adding call routing and PSTN termination services. We ignore the potential separation between the Internet access provider, that is, a carrier that provides physical and data link layer network connectivity to its customers, and the ISP that provides network layer services. We use the term VSP for simplicity, instead of the more generic term *Application Server Provider* (ASP).

The documents that the IETF *Emergency Context Resolution with Internet Technology* (ECRIT) working group is developing support multimedia-based emergency services, and not just voice. As is explained in more detail later in this article, emergency calls need to be identified for special call routing and handling services, and they need to carry the location of the caller for routing and dispatch. Only the calling device can reliably recognize emergency calls, while only the ISP typically has access to the current geographical location of the calling device based on its point of attachment to the network. The reliable handling of emergency calls is further complicated by the wide variety of access technologies in use, such as *Virtual Private Networks* (VPNs), other forms of tunneling, firewalls, and *Network Address Translators* (NATs).

This article describes the architecture of emergency services as defined by the IETF and some of the intermediate steps as end systems and the call-handling infrastructure transition from the current circuit-switched and emergency-calling-unaware *Voice-over-IP* (VoIP) systems to a true any-media, any-device emergency calling system.

IETF Emergency Services Architecture

The emergency services architecture developed by the IETF ECRIT working group is described in [1] and can be summarized as follows: *Emergency calls are generally handled like regular multimedia calls, except for call routing.* The ECRIT architecture assumes that PSAPs are connected to an IP network and support the *Session Initiation Protocol* (SIP)^[2] for call setup and messaging. However, the calling user agent may use any call signaling or instant messaging protocol, which the VSP then translates into SIP.

Nonemergency calls are routed by a VSP, either to another subscriber of the VSP, typically through some SIP session border controller or proxy, or to a PSTN gateway. For emergency calls, the VSP keeps its call routing role, routing calls to the emergency service system to reach a PSAP instead. However, we also want to allow callers that do not subscribe to a VSP to reach a PSAP, using nothing but a standard SIP^[2] user agent (see [3] and [4] for a discussion about this topic); the same mechanisms described here apply. Because the Internet is global, it is possible that a caller's VSP resides in a regulatory jurisdiction other than where the caller and the PSAP are located. In such circumstances it may be desirable to exclude the VSP and provide a direct signaling path between the caller and the emergency network. This setup has the advantage of ensuring that all parties included in the call delivery process reside in the same regulatory jurisdiction.

As noted in the introduction, the architecture neither forces nor assumes any type of trust or business relationship between the ISP and the VSP carrying the emergency call. In particular, this design assumption affects how location is derived and transported.

Providing emergency services requires three crucial steps, which we describe in the following sections: recognizing an emergency call, determining the caller's location, and routing the call and location information to the appropriate emergency service system operating a PSAP.

Recognizing an Emergency Call

In the early days of PSTN-based emergency calling, callers would dial a local number for the fire or police department. It was recognized in the 1960s that trying to find this number in an emergency caused unacceptable delays; thus, most countries have been introducing single nationwide emergency numbers, such as 911 in North America, 999 in The United Kingdom, and 112 in all European Union countries.

This standardization became even more important as mobile devices started to supplant landline phones. In some countries, different types of emergency services, such as police or mountain rescue, are identified by separate numbers. Unfortunately, more than 60 different emergency numbers are used worldwide, many of which also have nonemergency uses in other countries, so simply storing the list of numbers in all devices is not feasible. In addition, hotels and university campuses often use dial prefixes, so an emergency caller in some European universities may actually have to dial 0112 to reach the fire department.

Because of this diversity, the ECRIT architecture decided to separate the concept of an emergency dial string, which remains the familiar and regionally defined emergency number, and a protocol identifier that is used for identifying emergency calls within the signaling system. The calling end system has to recognize the emergency (service) dial string and translate it into an emergency service identifier, which is an extensible set of *Uniform Resource Names* (URNs) defined in RFC 5031^[5]. A common example for such a URN, defined to reach the generic emergency service, is `urn:service.sos`. The emergency service URN is included in the signaling request as the destination and is used to identify the call as an emergency call. If the end system fails to recognize the emergency dial string, the VSP may also perform this service.

Because mobile devices may be sold and used worldwide, we want to avoid manually configuring emergency dial strings. In general, a device should recognize the emergency dial string familiar to the user and the dial strings customarily used in the currently visited country. The *Location-to-Service Translation Protocol* (LoST)^[6], described in more detail later, also delivers this information.

Some devices, such as smartphones, can define dedicated user interface elements that dial emergency services. However, such mechanisms must be carefully designed so that they are not accidentally triggered, for example, when the device is in a pocket.

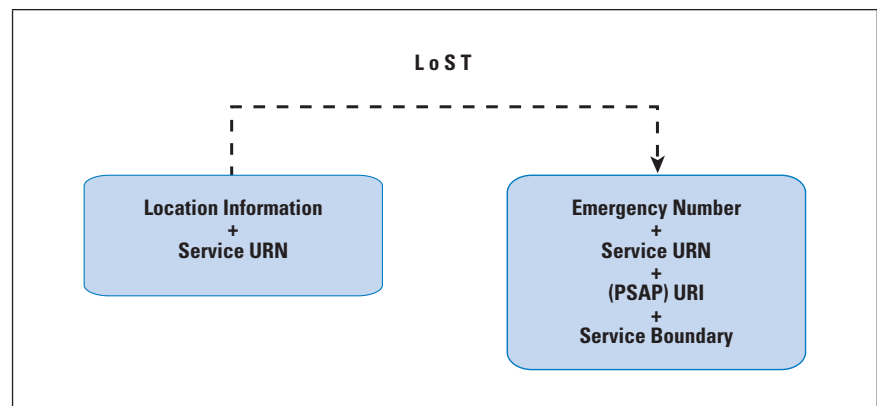
Emergency Call Routing

When an emergency call is recognized, the call needs to be routed to the appropriate PSAP. Each PSAP is responsible for only a limited geographic region, its service region, and some set of emergency services. For example, even in countries with a single general emergency number such as the United States, poison-control services maintain their own set of call centers. Because VSPs and end devices cannot keep a complete up-to-date mapping of all the service regions, a mapping protocol, LoST^[6], maps a location and service URN to a specific PSAP *Uniform Resource Identifier* (URI) and a service region.

LoST, illustrated in Figure 1, is a *Hypertext Transfer Protocol* (HTTP)-based query/response protocol where a client sends a request containing the location information and service URN to a server and receives a response containing the service URL, typically a SIP URL, the service region where the same information would be returned, and an indication of how long the information is valid. Both request and response are formatted as *Extensible Markup Language* (XML). For efficiency, responses are cached, because otherwise every small movement would trigger a new LoST request. As long as the client remains in the same service region, it does not need to consult the server again until the response returned reaches its expiration date. The response may also indicate that only a more generic emergency service is offered for this region. For example, a request for **urn:service:sos.marine** in Austria may be replaced by **urn:service:sos**. Finally, the response also indicates the emergency number and dial string for the respective service.

The number of PSAPs serving a country varies significantly. Sweden, for example, has 18 PSAPs, and the United States has approximately 6,200. Therefore, there is roughly one PSAP per 500,000 inhabitants in Sweden and one per 50,000 in the United States. As all-IP infrastructure is rolled out, smaller PSAPs may be consolidated into regional PSAPs. Routing may also take place in multiple stages, with the call being directed to an *Emergency Services Routing Proxy* (ESRP), which in turn routes the call to a PSAP, accounting for factors such as the number of available call takers or the language capabilities of the call takers.

Figure 1: High-Level Functions of Location-to-Service Translation (LoST) Protocol



Location Information

Emergency services need location information for three reasons: routing the call to the right PSAP, dispatching first responders (for example, policemen), and determining the right emergency service dial strings. It is clear that the location must be automatic for the first and third applications, but experience has shown that automated, highly accurate location information is vital to dispatching as well, rather than relying on callers to report their locations to the call taker.

Such information increases accuracy and avoids dispatch delays when callers are unable to provide location information because of language barriers, lack of familiarity with their surroundings, stress, or physical or mental impairment.

Location information for emergency purposes comes in two representations: geo(detic), that is, longitude and latitude, and civic, that is, street addresses similar to postal addresses. Particularly for indoor location, vertical information (floors) is very useful. Civic locations are most useful for fixed Internet access, including wireless hotspots, and are often preferable for specifying indoor locations, whereas geodetic location is frequently used for cell phones. However, with the advent of femto and pico cells, civic location is both possible and probably preferable because accurate geodetic information can be very hard to acquire indoors.

In almost all cases, location values are represented as *Presence Information Data Format Location Object* (PIDF-LO), an XML-based document to encapsulate civic and geodetic location information. The format of PIDF-LO is described in [7], with the civic location format updated in [8] and the geodetic location format profiled in [9]. The latter document uses the *Geography Markup Language* (GML) developed by the *Open Geospatial Consortium* (OGC) for describing commonly used location shapes.

Location can be conveyed either by value (“LbyV”) or by reference (“LbyR”). For the former, the XML location object is added as a message body in the SIP message. Location by value is particularly appropriate if the end system has access to the location information; for example, if it contains a *Global Positioning System* (GPS) receiver or uses one of the location configuration mechanisms described later in this section. In environments where the end host location changes frequently, the LbyR mechanism might be more appropriate. In this case, the LbyR is an HTTP/*Secure HTTP* (HTTPS) or SIP/*Secure SIP* (SIPS) URI, which the recipient needs to resolve to obtain the current location. Terminology and requirements for the LbyR mechanism are available in [10].

An LbyV and an LbyR can be obtained through location configuration protocols, such as the *HTTP Enabled Location Delivery* (HELD) protocol^[11] or *Dynamic Host Configuration Protocol* (DHCP)^[12, 13]. When obtained, location information is required for LoST queries, and that information is added to SIP messages^[14].

The requirements for location accuracy differ between routing and dispatch. For call routing, city or even county-level accuracy is often sufficient, depending on how large the PSAP service areas are, whereas first responders benefit greatly when they can pinpoint the caller to a particular building or, better yet, apartment or office for indoor locations, and an outdoor area of at most a few hundred meters. This detailed location information avoids having to search multiple buildings, for example, for medical emergencies.

As mentioned previously, the ISP is the source of the most accurate and dependable location information, except for cases where the calling device has built-in location capabilities, such as GPS, when it may have more accurate location information. For landline Internet connections such as DSL, cable, or fiber-to-the-home, the ISP knows the provisioned location for the network termination, for example. The IETF GEOPRIV working group has developed protocol mechanisms, called *Location Configuration Protocols*, so that the end host can request and receive location information from the ISP. The Best Current Practice document for emergency calling^[15] enumerates three options that clients should universally support: DHCP civic^[16] and geo^[12] (with a revision of RFC 3825 in progress^[17]), and HELD^[11]. HELD uses XML query and response objects carried in HTTP exchanges. DHCP does not use the PIDF-LO format, but rather more compact binary representations of locations that require the endpoint to construct the PIDF-LO.

Particularly for cases where end systems are not location-capable, a VSP may need to obtain location information on behalf of the end host^[18].

Obtaining at least approximate location information at the time of the call is time-critical, because the LoST query can be initiated only after the calling device or VSP has obtained location information. Also, to accelerate response, it is desirable to transmit this location information with the initial call signaling message. In some cases, however, location information at call setup time is imprecise. For example, a mobile device typically needs 15 to 20 seconds to get an accurate GPS location “fix,” and the initial location report is based on the cell tower and sector. For such calls, the PSAP should be able to request more accurate location information either from the mobile device directly or the *Location Information Server* (LIS) operated by the ISP. The SIP event notification extension, defined in RFC 3265^[19], is one such mechanism that allows a PSAP to obtain the location from an LIS. To ensure that the PSAP is informed only of pertinent location changes and that the number of notifications is kept to a minimum, event filters^[20] can be used.

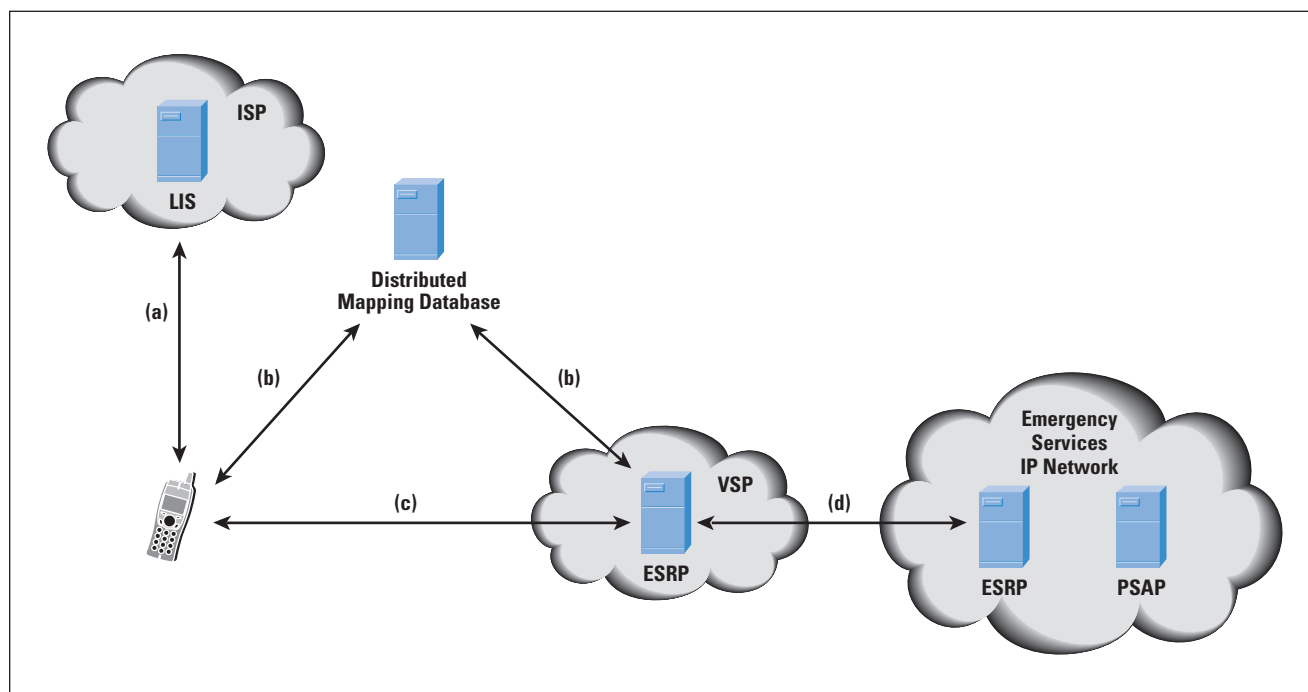
The two-stage location refinement mechanism described previously works best when location is provided by reference (LbyR) in the SIP INVITE call setup request. The PSAP subscribes to the LbyR provided in the SIP exchange and the LbyR refers to the LIS in the ISP’s network. In addition to a SIP URI, the LbyR message can also contain an HTTP/HTTPS URI. When such a URI is provided, an HTTP-based protocol can be used to retrieve the current location^[21].

Obligations

This section discusses the requirements the different entities need to satisfy, based on Figure 2. A more detailed description can be found in [15].

Note that this narration focuses on the final stage of deployment and does not discuss the transition architecture, in which some implementation responsibilities can be rearranged, with an effect on the overall functions offered by the emergency services architecture. A few variations were introduced to handle the transition from the current system to a fully developed ECRIT architecture.

Figure 2: Main Components Involved in an Emergency Call



With the work on the IETF emergency architecture, we have tried to balance the responsibilities among the participants, as described in the following sections.

End Hosts

An end host, through its VoIP application, has three main responsibilities: it has to attempt to obtain its own location, determine the URI of the appropriate PSAP for that location, and recognize when the user places an emergency call by examining the dial string. The end host operating system may assist in determining the device location.

The protocol interaction for location configuration is indicated as interface (a) in Figure 2; numerous location configuration protocols have been developed to provide this capability.

A VoIP application needs to support the LoST protocol^[6] in order to determine the emergency service dial strings and the PSAP URI. Additionally, the device needs to understand the service identifiers, defined in [5].

As currently defined, it is assumed that SIP can reach PSAPs, but PSAPs may support other signaling protocols, either directly or through a protocol translation gateway. The LoST retrieval results indicate whether other signaling protocols are supported. To provide support for multimedia, use of different types of codecs may be required; details are available in [15].

ISP

The ISP has to make location information available to the endpoint through one or more of the location configuration protocols.

In order to route an emergency call correctly to a PSAP, an ISP may initially disclose the approximate location for routing to the endpoint and give more precise location information later, when the PSAP operator dispatches emergency personnel. The functions required by the IETF emergency services architecture are restricted to the disclosure of a relatively small amount of location information, as discussed in [22] and in [23].

The ISP may also operate a (caching) LoST server to improve the robustness and reliability of the architecture. This server lowers the round-trip time for contacting a LoST server, and the caches are most likely to hold the mappings of the area where the emergency caller is currently located.

When ISPs allow Internet traffic to traverse their network, the signaling and media protocols used for emergency calls function without problems. Today, there are no legal requirements to offer prioritization of emergency calls over IP-based networks. Although the standardization community has developed a range of *Quality of Service* (QoS) signaling protocols, they have not experienced widespread deployment.

VSP

SIP does not mandate that call setup requests traverse SIP proxies; that is, SIP messages can be sent directly to the user agent. Thus, even for emergency services it is possible to use SIP without the involvement of a VSP. However, in terms of deployment, it is highly likely that a VSP will be used. If a caller uses a VSP, this VSP often forces all calls, emergency or not, to traverse an outbound proxy or *Session Border Controller* (SBC) operated by the VSP. If some end devices are unable to perform a LoST lookup, VSP can provide the necessary functions as a backup solution.

If the VSP uses a signaling or media protocol that the PSAP does not support, it needs to translate the signaling or media flows.

VSPs can assist the PSAP by providing identity assurance for emergency calls; for example, using [30], thus helping to prosecute prank callers. However, the link between the subscriber information and the real-world person making the call is weak.

In many cases, VSPs have, at best, only the credit card data for their customers, and some of these customers may use gift cards or other anonymous means of payment.

PSAP

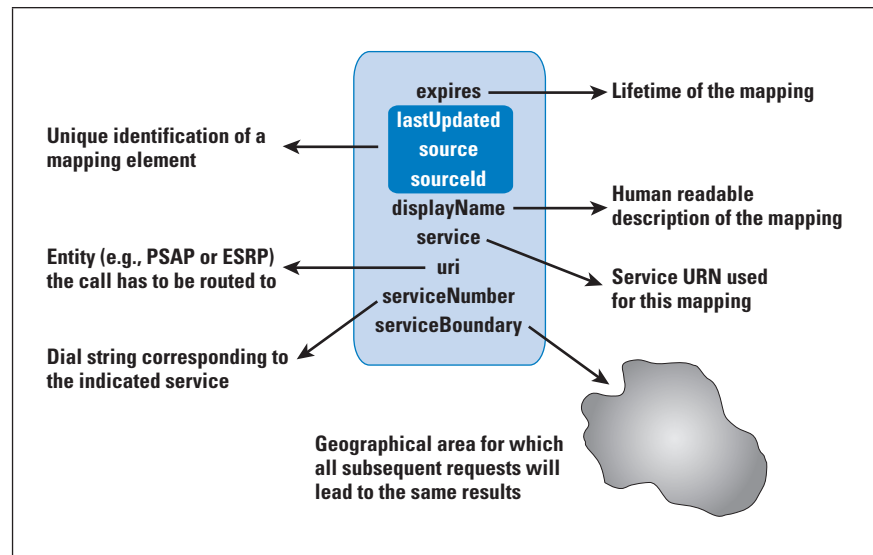
The emergency services Best Current Practice document [15] discusses only the standardization of the interfaces from the VSP and ISP toward PSAPs and some parts of the PSAP-to-PSAP call transfer mechanisms that are necessary for emergency calls to be processed by the PSAP. Many aspects related to the internal communication within a PSAP, between PSAPs as well as between a PSAP and first responders, are beyond the scope of the IETF specification.

When emergency calling has been fully converted to Internet protocols, PSAPs must accept calls from any VSP, as shown in interface (d) of Figure 2. Because calls may come from all sources, PSAPs must develop mechanisms to reduce the number of malicious calls, particularly calls containing intentionally false location information. Assuring the reliability of location information remains challenging, particularly as more and more devices are equipped with *Global Navigation Satellite Systems* (GNSS) receivers, including GPS and Galileo, allowing them to determine their own location^[24]. However, it may be possible in some cases to check the veracity of the location information an endpoint provides by comparing it against infrastructure-provided location information; for example, a LIS-determined location.

Mapping Architecture

So far we have described LoST as a client-server protocol. Similar to the *Domain Name System* (DNS), a single LoST server does not store the mapping elements for all PSAPs worldwide, for both technical and administrative reasons. Thus, there is a need to let LoST servers interact with other LoST servers, each covering a specific geographical region. Working together, LoST servers form a distributed mapping database, with each server carrying mapping elements, as shown in Figure 3. LoST servers may be operated by different entities, including the ISP, the VSP, or another independent entity, such as a governmental agency. Typically, individual LoST servers offer the necessary mapping elements for their geographic regions to others. However, LoST servers may also cache mapping elements of other LoST servers either through data synchronization mechanisms (for example, FTP or exports from a *Geographical Information System* [GIS] or through a specialized protocol^[25]) or by regular usage of LoST. This caching improves performance and increases the robustness of the system.

Figure 3: Mapping Element



A detailed description of the mapping architecture with examples is available in [29].

Steps Toward an IETF Emergency Services Architecture

The architecture described so far requires changes both in already-deployed VoIP end systems and in the existing PSAPs. The speed of transition and the path taken vary between different countries, depending on funding and business incentives. Therefore, it is generally difficult to argue whether upgrading endpoints or replacing the emergency service infrastructure will be easier. In any case, the transition approaches being investigated consider both directions. We can distinguish roughly four stages of transition (Note: The following descriptions omit many of the details because of space constraints):

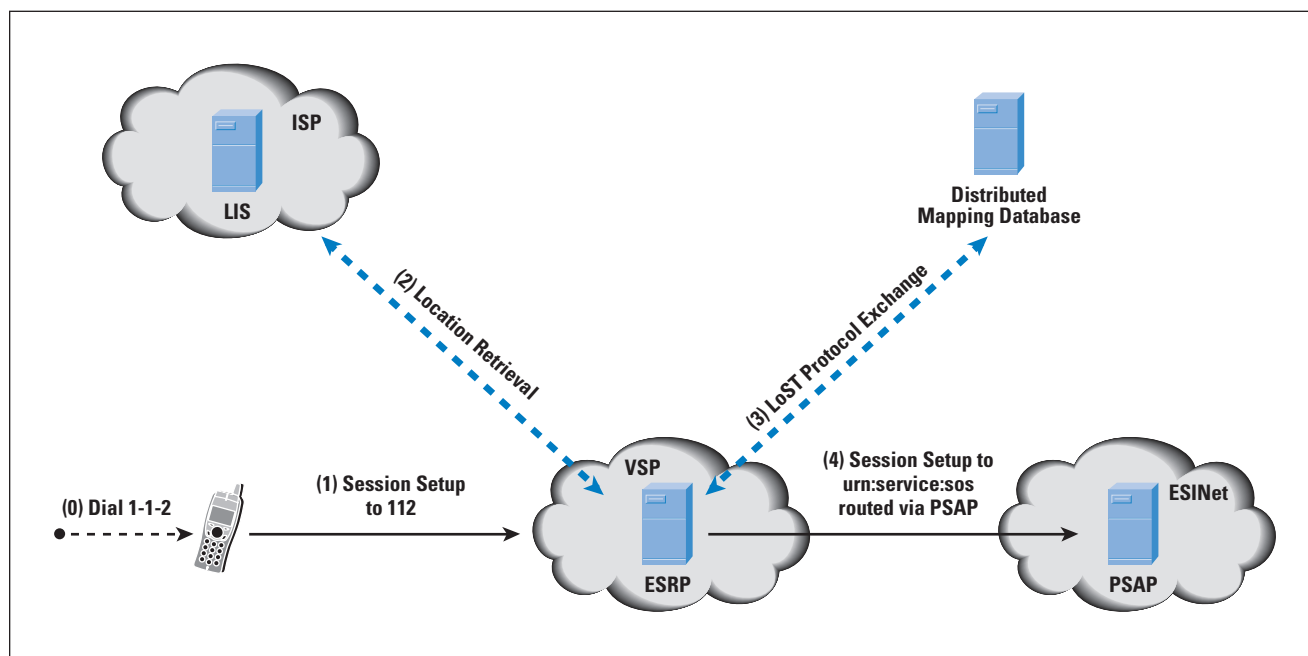
1. Initially, VoIP end systems cannot place emergency calls at all; for example, many software clients, such as *GoogleTalk*, cannot place emergency calls.
2. In a second stage, VoIP callers manually configure their location, and emergency calls are routed to the appropriate PSAP as circuit-switched calls through PSTN gateways using technologies similar to mobile calls. This level of service is now offered in some countries for PSTN-replacement VoIP services; that is, VoIP services that are offered as replacement for the home phone. In the United States, this service is known as the “NENA I2” service.
3. In a third stage, PSAPs maintain two separate infrastructures, one for calls arriving through an IP network and the traditional infrastructure.
4. In the final stage, all calls, including those from traditional cell phones and analog landline phones, reach the PSAP through IP networks, with the traditional calls converted to the ECRIT requirements by the carriers or the emergency service infrastructure.

If devices are used in environments without location services, the VSP's SIP proxy may need to insert location information based on estimates or subscriber data. These cases are described briefly in the following sections.

Traditional Endpoints

Figure 4 shows an emergency services architecture with traditional endpoints. When the emergency caller dials the Europeanwide emergency number 112 (step 0), the device treats it as any other call without recognizing it as an emergency call; that is, the dial string provided by the endpoint that may conform to RFC 4967^[26] or RFC 3966^[27] is signaled to the VSP (step 1). Recognition of the dial string is then left to the VSP for processing or sorting; the same is true for location retrieval (step 2) and routing to the nearest (or appropriate) PSAP (step 3). Dial-string recognition, location determination, and call routing are simpler to carry out using a fixed device and the voice and application service provided through the ISP than they are when the VSP and the ISP are two separate entities.

Figure 4: Emergency Services Architecture with Traditional Endpoints



There are two main challenges to overcome when dealing with traditional devices: First, the VSP must discover the LIS that knows the location of the IP-based end host. The VSP is likely to know only the IP address of that device, visible in the call signaling that arrives at the VSP. When a LIS is discovered and contacted and some amount of location information is available, then the second challenge arises, namely, how to route the emergency call to the appropriate PSAP. To accomplish the latter task it is necessary to have some information about the PSAP boundaries available.

Reference [15] does not describe a complete and detailed solution but uses building blocks specified in ECRIT. Still, this deployment scenario shows many constraints:

- Only the emergency numbers configured at the VSP are understood. This situation may lead to cases where a dialed emergency number is not recognized.
- Using the IP address to find the ISP is challenging and may, in case of mobility protocols and VPNs, lead to wrong results.
- Security concerns might arise when a potentially large number of VSPs or ASPs are able to retrieve location information from an ISP. It is likely that only authorized VSP and ASPs will be granted access. Hence, it is unlikely that such a solution would work smoothly across national boundaries.
- When the user agent does not recognize the emergency call, functions such as call waiting, call transfer, three-way call, flash hold, and outbound call blocking cannot be disabled.
- The user-agent software may block callbacks from the PSAP.
- Privacy settings may not get considered and identity may get disclosed to unauthorized parties. These identity privacy features exist in some jurisdictions even in emergency situations.
- Certain VoIP call features may not be supported, such as REFER (for conference call and transfer to secondary PSAP) and *Globally Routable UA URI* (GRUU).
- User agents will not convey location information to the VSP (even if available).

Partially Upgraded End Hosts

A giant step forward in simplifying the handling of IP-based emergency calls is to provide the end host with some information about the ISP so that LIS discovery is possible. The end host may, for example, learn the ISP's domain name by using LIS discovery^[28], or might even obtain a *Location by Reference* (LbyR) through the DHCP-URI option^[13] or through HELD^[11]. The VSP can then either resolve the LbyR in order to route the call or use the domain to discover a LIS using DNS.

Additional software upgrades at the end device may allow for recognition of emergency calls based on some preconfigured emergency numbers (for example, 112 and 911) and allow for the implementation of other emergency service-related features, such as disabling silence suppression during emergency calls.

Outlook

In most countries, national and sometimes regional telecommunications regulators, such as the *Federal Communications Commission* (FCC) and individual states, or the European Union, strongly influence how emergency services are provided, who pays for them, and the obligations that the various parties have. Regulation is, however, still at an early stage: in most countries current requirements demand only manual update of location information by the VoIP user. The ability to obtain location information automatically is, however, crucial for reliable emergency service operation, and it is required for nomadic and mobile devices. (Nomadic devices remain in one place during a communication session, but are moved frequently from place to place. Laptops with Wi-Fi interfaces are currently the most common nomadic devices.)

Regulators have traditionally focused on the national or, at most, the European level, and the international nature of the Internet poses new challenges. For example, mobile devices are now routinely used beyond their country of purchase and, unlike traditional cellular phones, need to support emergency calling functions. It appears likely that different countries will deploy IP-based emergency services over different time horizons, so travelers may be surprised to find that they cannot call for emergency assistance outside their home country.

The separation between Internet access and application providers on the Internet is one of the most important differences to existing circuit-switched telephony networks. A side effect of this separation is the increased speed of innovation at the application layer, and the number of new communication mechanisms is steadily increasing. Many emergency service organizations have recognized this trend and advocated for the use of new communication mechanisms, including video, real-time text, and instant messaging, to offer improved emergency calling support for citizens. Again, this situation requires regulators to rethink the distribution of responsibilities, funding, and liability.

Many communication systems used today lack accountability; that is, it is difficult or impossible to trace malicious activities back to the persons who caused them. This problem is not new, because pay phones and prepaid cell phones have long offered mischief makers the opportunity to place hoax calls, but the weak user registration procedures, the lack of deployed end-to-end identity mechanisms, and the ease of providing fake location information increases the attack surface at PSAPs. Attackers also have become more sophisticated over time, and Botnets that generate a large volume of automated emergency calls to exhaust PSAP resources, including call takers and first responders, are not science fiction.

References

- [1] Rosen, B., Schulzrinne, H., Polk, J., and A. Newton, “Framework for Emergency Calling Using Internet Multimedia,” Internet Draft, work in progress, **draft-ietf-ecrit-framework-11**, July 2010.
- [2] Rosenberg, J., Schulzrinne, H., Camarillo, G., Johnston, A., Peterson, J., Sparks, R., Handley, M., and E. Schooler, “SIP: Session Initiation Protocol,” RFC 3261, June 2002.
- [3] Winterbottom, J., Thomson, M., Tschofenig, H., and H. Schulzrinne, “ECRIT Direct Emergency Calling,” Internet Draft, work in progress, **draft-winterbottom-ecrit-direct-02.txt**, March 2010.
- [4] Schulzrinne, H., McCann, S., Bajko, G., Tschofenig, H., and D. Kroesenberg, “Extensions to the Emergency Services Architecture for Dealing with Unauthenticated and Unauthorized Devices,” Internet Draft, work in progress, **draft-ietf-ecrit-unauthenticated-access-00.txt**, September 2010.
- [5] Schulzrinne, H., “A Uniform Resource Name (URN) for Emergency and Other Well-Known Services,” RFC 5031, January 2008.
- [6] Hardie, T., Newton, A., Schulzrinne, H., and H. Tschofenig, “LoST: A Location-to-Service Translation Protocol,” RFC 5222, August 2008.
- [7] Peterson, J., “A Presence-based GEOPRIV Location Object Format,” RFC 4119, December 2005.
- [8] Thomson, M. and J. Winterbottom, “Revised Civic Location Format for Presence Information Data Format Location Object (PIDF-LO),” RFC 5139, February 2008.
- [9] Winterbottom, J., Thomson, M., and H. Tschofenig, “GEOPRIV Presence Information Data Format Location Object (PIDF-LO) Usage Clarification, Considerations, and Recommendations,” RFC 5491, March 2009.
- [10] R. Marshall, “Requirements for a Location-by-Reference Mechanism,” RFC 5808, May 2010.
- [11] M. Barnes, “HTTP Enabled Location Delivery (HELD),” RFC 5985, September 2010.
- [12] Polk, J., Schnizlein, J., and M. Linsner, “Dynamic Host Configuration Protocol Option for Coordinate-based Location Configuration Information,” RFC 3825, July 2004.

- [13] Polk, J., “Dynamic Host Configuration Protocol (DHCP) IPv4 and IPv6 Option for a Location Uniform Resource Identifier (URI),” Internet Draft, work in progress, **draft-ietf-geopriv-dhcp-lbyr-uri-option-08**, July 2010.
- [14] Polk, J., Rosen, B., and J. Peterson, “Location Conveyance for the Session Initiation Protocol,” Internet Draft, work in progress, **draft-ietf-sipcore-location-conveyance-03**, July 2010.
- [15] Rosen, B. and J. Polk, “Best Current Practice for Communications Services in Support of Emergency Calling,” Internet Draft, work in progress, **draft-ietf-ecrit-phonebcp-15**, July 2010.
- [16] Schulzrinne, H., “Dynamic Host Configuration Protocol (DHCPv4 and DHCPv6) Option for Civic Addresses Configuration Information,” RFC 4776, November 2006.
- [17] Polk, J., Schnizlein, J., Linsner, M., and B. Aboba, “Dynamic Host Configuration Protocol Option for Coordinate-based Location Configuration Information,” Internet Draft, work in progress, **draft-ietf-geopriv-rfc3825bis-11**, July 2010.
- [18] Winterbottom, J., Thomson, M., Tschofenig, H., and R. Barnes, “Use of Device Identity in HTTP-Enabled Location Delivery (HELD),” Internet Draft, work in progress, **draft-ietf-geopriv-held-identity-extensions-04**, June 2010.
- [19] Roach, A., “Session Initiation Protocol (SIP)-Specific Event Notification,” RFC 3265, June 2002.
- [20] Mahy, R., Rosen, B., and H. Tschofenig, “Filtering Location Notifications in the Session Initiation Protocol,” Internet Draft, work in progress, **draft-ietf-geopriv-loc-filters-11**, March 2010.
- [21] Winterbottom, J., Tschofenig, H., Schulzrinne, H., Thomson, M., and M. Dawson, “A Location Dereferencing Protocol Using HELD,” Internet Draft, work in progress, **draft-ietf-geopriv-deref-protocol-01**, September 2010.
- [22] Schulzrinne, H., Liess, L., Tschofenig, H., Stark, B., and A. Kuett, “Location Hiding: Problem Statement and Requirements,” Internet Draft, work in progress, **draft-ietf-ecrit-location-hiding-req-04**, Feb 2010.
- [23] Barnes, R., and M. Lepinski, “Using Imprecise Location for Emergency Context Resolution,” Internet Draft, work in progress, **draft-ietf-ecrit-rough-loc-03**, August 2010.

- [24] Tschofenig, H., Schulzrinne, H., and B. Aboba, "Trustworthy Location Information," Internet Draft, work in progress, **draft-tschofenig-ecrit-trustworthy-location-00**, September 2010.
- [25] Schulzrinne, H., and H. Tschofenig, "Synchronizing Location-to-Service Translation (LoST) Servers," Internet Draft, work in progress, **draft-ietf-ecrit-lost-sync-10**, March 2010.
- [26] B. Rosen, "Dial String Parameter for the Session Initiation Protocol Uniform Resource Identifier," RFC 4967, July 2007.
- [27] H. Schulzrinne "The tel URI for Telephone Numbers," RFC 3966, December 2004.
- [28] Thomson, M. and J. Winterbottom, "Discovering the Local Location Information Server (LIS)," RFC 5986, September 2010.
- [29] H. Schulzrinne, "Location-to-URL Mapping Architecture and Framework," RFC 5582, September 2009.
- [30] C. Jennings, J. Peterson, and M. Watson, "Private Extensions to the Session Initiation Protocol (SIP) for Asserted Identity within Trusted Networks," RFC 3325, November 2002.

HENNING SCHULZRINNE, Levi Professor of Computer Science at Columbia University, received his Ph.D. from the University of Massachusetts in Amherst, Massachusetts. He was an MTS at AT&T Bell Laboratories and an associate department head at GMD-Fokus (Berlin) before joining the Computer Science and Electrical Engineering departments at Columbia University. He served as chair of Computer Science from 2004 to 2009. Protocols that he co-developed, such as RTP, RTSP, and SIP, are now Internet standards, used by almost all Internet telephony and multimedia applications. His research interests include Internet multimedia systems, ubiquitous computing, and mobile systems. He is a Fellow of the IEEE. E-mail: **hgs@cs.columbia.edu**

HANNES TSCHOFENIG received a Diploma degree from the University of Klagenfurt, Austria. He joined Siemens Corporate Technology, Munich, in 2001 and joined Nokia Siemens Networks in April 2007 to move to Finland in December 2007, where he focuses on standards development. Most of his time is dedicated to the participation in the Internet Engineering Task Force (IETF) where he, among other responsibilities, co-chaired the ECRIT working group from 2005 to early 2010. Additionally, he co-chairs the Next Generation 112 Technical Committee of the European Emergency Number Association (EENA) and contributes to the technical specifications developed within the National Emergency Number Association (NENA), and he co-organized the SDO emergency services workshop series. In March 2010 he joined the Internet Architecture Board (IAB). E-mail: **hannes.tschofenig@nsn.com**

Integration of Core BGP/MPLS VPN Networks

by Paul Veitch, Paul Hitchen, and Martin Mitchell, BT Innovate & Design

This article explores the architectural and operational challenges involved in integrating an existing standalone core *Border Gateway Protocol* (BGP)/*Multiprotocol Label Switching* (MPLS) VPN network onto a target *Next-Generation Network* (NGN). The rationale for consolidating and transforming multiple networks is explained, mainly in terms of potential cost savings and operational simplification achieved by the network operator. The article specifically focuses on the MPLS *Carrier-supporting-Carrier* (CsC) architectural framework, which allows the serving nodes of one MPLS VPN network to be interconnected through the serving nodes of another MPLS VPN network. The required architectural building blocks to implement CsC, the manner in which routing protocols must interact, as well as end-to-end packet flow and label encapsulation are all explained. The main design and operational challenges, including maintaining performance levels for customers, network resiliency, fault-handling, and capacity management, are also addressed in this article.

Network operators are under increasing pressure to deliver exceptional levels of customer experience and service while decreasing the capital and operational cost base of their networks. Many operators have traditionally built multiple network platforms, each of which has been uniquely designed to meet the requirements of specific services targeted at specific customer markets, such as voice, broadband IP, *Virtual Private Networks* (VPNs), etc.

In a bid to remain competitive and achieve cost reductions and operational simplifications, many operators have built all IP-based NGNs. The principal transformational benefits of an NGN with a single protocol such as IP at its heart include versatility in catering for multiple traffic requirements (for example, by employing IP *Quality-of-Service* [QoS] techniques), the ability to introduce novel and reusable services and features in a flexible manner, and the potential to maximise vendor interworking due to standards-based technology.

When a network operator builds an NGN, the challenge remains as to how to migrate *existing* networks and customers onto the new platform. The full commercial benefits of an NGN can be properly realised only after legacy networks are either consolidated or phased out completely. Many important factors must be considered, including the cost benefits, the potential effect on end customers, and the operational approach to carrying out migrations. These concerns must be weighed against the commercial and business risks associated with the alternative approach of sustaining and running multiple standalone platforms indefinitely.

This article focuses on a specific scenario: how to integrate an existing BGP/MPLS VPN network that provides VPN services to a corporate customer base with a “target” NGN. Following a brief overview of MPLS VPN services and networks, the rationale for consolidating multiple MPLS VPN networks is explained, mainly in terms of potential cost savings and operational simplification achieved by the network operator. The article then details the MPLS CsC architectural framework that allows the serving nodes or *Points of Presence* (POPs) of one MPLS VPN network to be interconnected to the serving nodes of another MPLS VPN network. The way in which routing protocols must interact and the subsequent effect on end-to-end packet forwarding across a CsC-enabled core network are explained. The principal design and operational challenges introduced by integrating core MPLS networks are then outlined, including maintaining performance levels, network resiliency, fault management, and capacity management.

The Business Case for MPLS VPN Network Consolidation

VPNs are an attractive solution to serve the enterprise networking requirements of a wide range of businesses from *Small-to-Medium Enterprises* (SMEs) to multinational “blue-chip” corporate organisations. Essentially, VPNs provide a transparent network infrastructure that allows multiple customer sites to communicate over a shared backbone network, as though they are using their own private network, regardless of geographical location. Typical applications that run across an organisation’s VPN include corporate Intranet, mail services, and *Voice-over-IP* (VoIP) telephony.

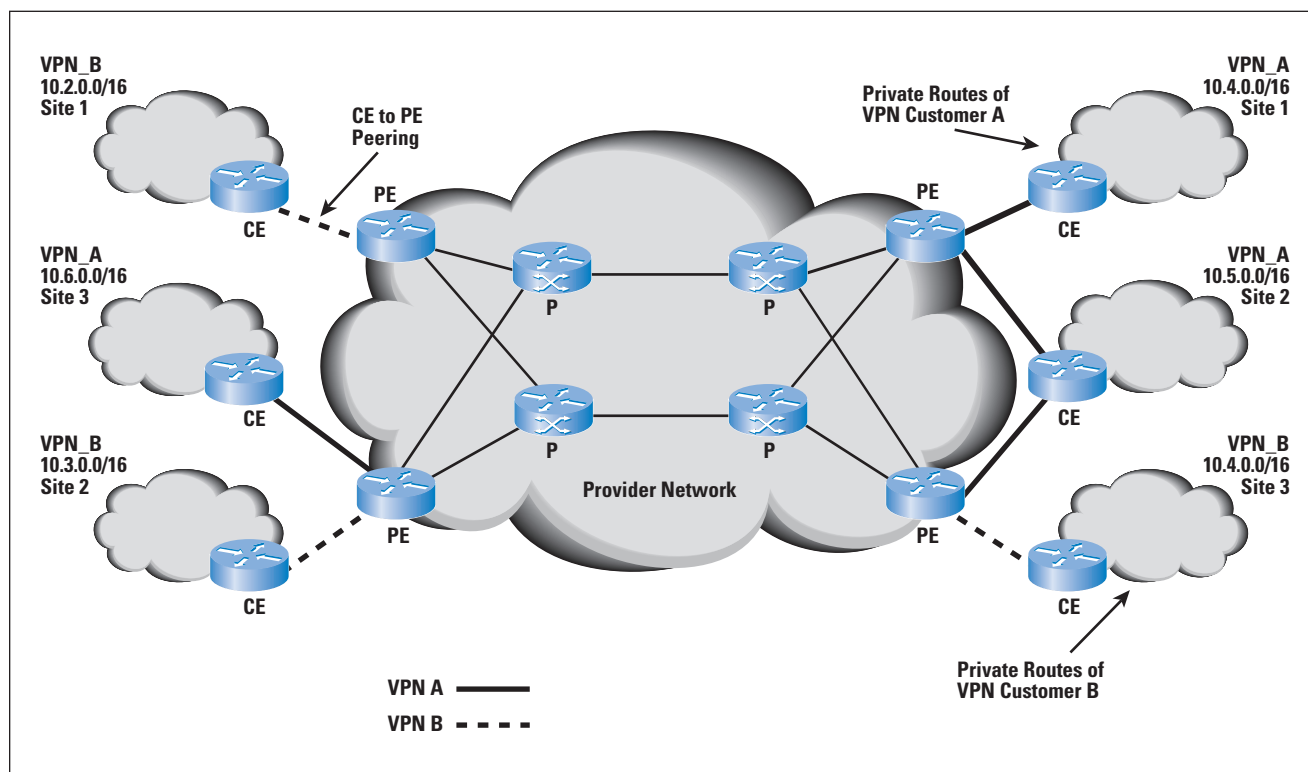
Although distinct categories of VPN networking technology exist^[1], this article focuses exclusively on “Layer 3” BGP/MPLS VPNs, as defined in RFC 4364^[2] and other related Internet Drafts. Such networks have been deployed for more than 10 years and have seen significant growth during that period.

The critical core network elements of a provider-provisioned BGP/MPLS VPN network are *Provider Edge* (PE) and *Provider Core* (P) routers, as shown in Figure 1.

PE routers terminate customer access circuits, whereas P routers perform packet forwarding and typically do not have directly connected customer access circuits. PE routers perform label encapsulation and de-encapsulation, P routers run label switching, and both operate control-plane protocols that build MPLS *Label Switched Paths* (LSPs) from each PE to each other PE. Many protocols can be used to establish these LSPs; a commonly deployed approach uses the *Label Distribution Protocol* (LDP) in conjunction with an *Interior Gateway Protocol* (IGP), such as *Open Shortest Path First* (OSPF).

When a PE forwards a VPN-addressed packet across the core, it adds an inner MPLS label to identify the VPN of which the packet is a member and then an outer MPLS label to identify the egress PE router. Any intermediate P routers switch the packet to the egress PE using the outer label only. The egress PE uses the inner label to determine which VPN or port to forward the packet to.

Figure 1: Overview of BGP/MPLS VPN Network



The *Customer Edge* (CE) router is not considered part of the provider's core network. It acts as a peer of the PE router, but not a peer of other CE routers. Each PE router supports multiple routing and forwarding tables, called *Virtual Route Forwarding* (VRF) tables. VRF routes are logically separate, and they may contain IP prefixes received from the CE router that overlap with addresses in other VRFs. (For example, in Figure 1, VPN_A, site 1 has the same private routes as VPN_B, site 3.) VPNs are formed by defining individual customer accesses to be members of a specific VRF table, with several sites formed on one PE by defining all sites to use the same VRF table or allocating each site a VRF table and controlling connectivity through selective import and export of the IP routes of each VRF table.

The PE routers use an extended variant of BGP for signaling between themselves and propagating information about the actual routes of each VPN, as well as the inner MPLS label. The extended BGP, referred to as *Multiprotocol BGP*, carries each VPN route together with two new fields, the *Route Distinguisher* (RD) and the *Route Target* (RT), a form of extended BGP Community.

The RD is added to each VPN route to ensure that routes from different customers are unique; BGP treats VPN routes as equal only if both the RD and the IP prefix mask are equal. BGP uses RTs to indicate a group of routes, thus defining VPN membership information for exchange between PEs.

Maintenance Costs of BGP/MPLS VPN Networks

As detailed in the previous section, the main core components of a VPN network based on BGP/MPLS technology are the PE and P routers. Although not shown in detail in Figure 1, another critical element of a core VPN network is the *Wide-Area Network* (WAN) topology that interconnects the P (core) routers residing in specific service nodes, also called POPs. The WAN topology is essentially the way in which transmission links—typically *Synchronous Optical Network* (SONET)/*Packet over SONET/SDH* (PoS), Gigabit Ethernet, or 10 Gigabit Ethernet—are used to interconnect the POPs together.

It follows that maintenance costs associated with a self-contained MPLS VPN network will be incurred for PE and P routers, as well as the interconnecting WAN transmission links. These maintenance costs will split into capital and operational elements.

Capital expenditures are required on an ongoing basis for all IP router infrastructure (PE and P routers), for example, to upgrade hardware to meet increasing capacity demands, replace faulty line cards and processors, or replace end-of-life hardware with newer equipment. Capital expenditures are also needed on WAN links, for example, to replace faulty line cards and optics, as well as to deploy increased capacity transmission links to cater for traffic growth across the core network. Further capital costs accrue from accommodation-related aspects such as power, racking, and air conditioning.

Additional maintenance costs reside in the operational space. For example, if an MPLS VPN network has 40 POP locations, each with a pair of P (core) routers, the 80 core routers will consume a certain amount of operational team resources for critical maintenance, scheduled maintenance activities, and ongoing monitoring and reporting of router status (processors and line cards).

Benefits of Core Integration

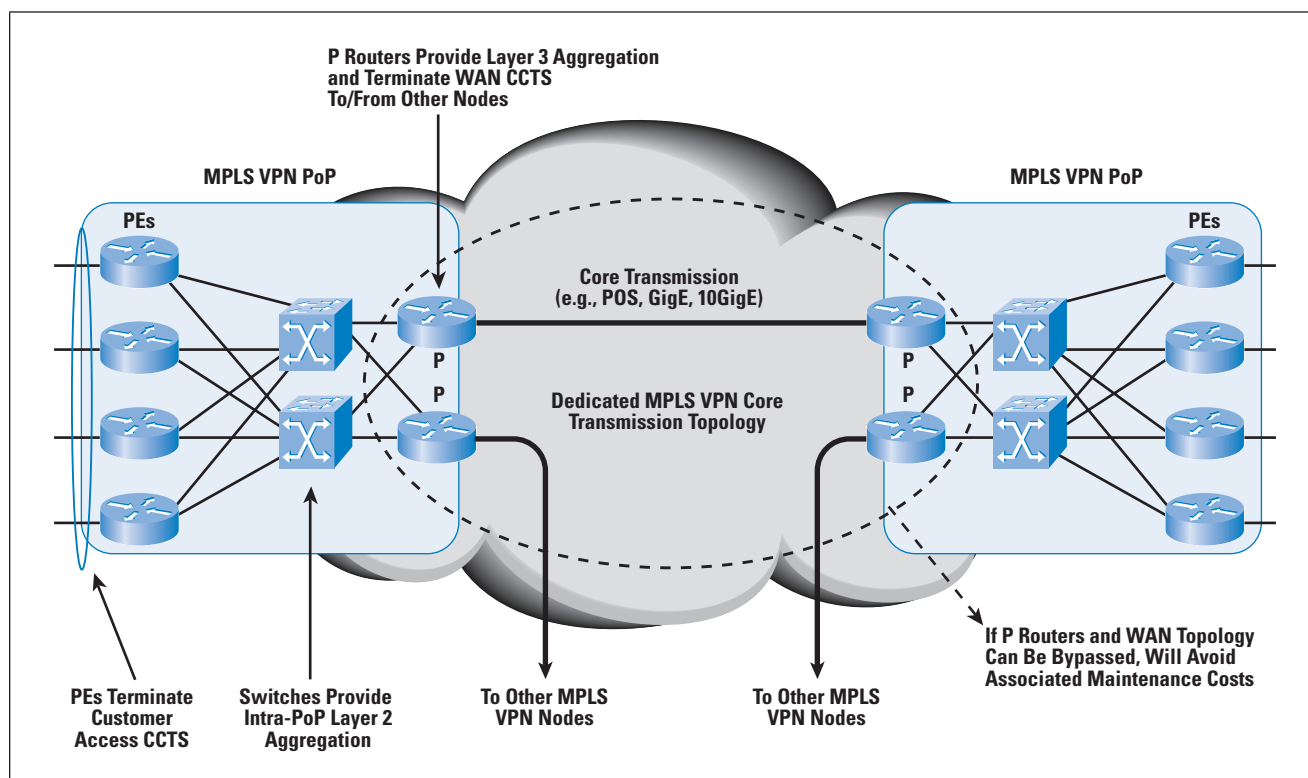
If a network operator has deployed an IP-based NGN alongside an existing MPLS VPN network, the question should be asked: can the existing MPLS VPN network be integrated onto the NGN so as to avoid some or all of the previously stated maintenance costs? One approach would be to target the P (core) routers and WAN transmission links for eventual removal (Figure 2) and replacement by suitable connectivity of the MPLS VPN nodes to the NGN network. The VPN PE routers that often terminate large volumes of customer access circuits and host the rich service-related functions for corporate VPN services can essentially be left *in situ*, minimising the effect on end customers and confining the integration of networks to the inner part of the core infrastructure. The way in which this goal can actually be achieved in practice is detailed in the next section.

The main benefits that can be accrued for the network operator are as follows:

- Substantial cost avoidance for maintaining and upgrading P (core) routers and dedicated WAN links for the existing MPLS VPN network can be achieved (Figure 2). As much as a 35-percent reduction of fixed inner core capital costs is possible.
- If the technical solution for core integration can be made as reusable as possible, then in addition to allowing integration of “same provider” core networks, the network operator could provide the capability on a wholesale basis for other service providers. This capability could be a potentially significant source of new revenue.
- From an operational perspective, integration of core networks should lend itself to a singular and much more streamlined approach to capacity planning, fault management, and network monitoring.

The combination of all these benefits can produce a compelling business case for network operators to consolidate core MPLS-based network platforms.

Figure 2: MPLS VPN Network Showing Inner Core Components Targeted for Replacement



Carrier-supporting-Carrier Framework

Carrier-supporting-Carrier (CsC) is a term used to describe a situation where one network, designated the *customer carrier*, is permitted to use a segment of another network, designated the *backbone carrier*^[3]. Although the term “Carrier of Carriers” is also used to describe the same architectural framework, this article uses Carrier-supporting-Carrier for consistency. In principle, the two “carrier” networks could belong to the same organisation, or could belong to two different organisations. Whatever the case, there is no reason why the backbone carrier cannot support multiple customer carrier networks. Furthermore, the customer carrier network itself can be either a BGP/MPLS VPN network providing Layer 3 VPN services or an *Internet Service Provider* (ISP) network^[3].

A network operator with an existing BGP/MPLS VPN network infrastructure that has also built an IP-based NGN based on BGP/MPLS technology as per RFC 4364^[2] could choose to exploit the CsC architectural framework to merge the two core networks. In such a scenario, the existing BGP/MPLS VPN network that serves the needs of VPN business customers would be viewed as the “customer carrier,” whereas the NGN network would be positioned as the “backbone carrier.”

Physical Connectivity and CsC VRF Creation

In order to integrate an existing BGP/MPLS VPN network such as that shown in Figure 2, with an NGN core belonging to the same or different organisation, the NGN network must be enabled to act as a backbone carrier. Assuming the NGN network is configured to support BGP/MPLS VPNs as per RFC 4364^[2], it comprises PE and P router core infrastructure. The PE routers of the NGN acting as the backbone carrier are denoted “CsC-PEs.” The PE routers of the existing BGP/MPLS VPN network, that is, the customer carrier network that is being itself integrated with the NGN core, are denoted “CsC-CEs.”

As shown in Figure 3, the NGN backbone carrier network provides MPLS VPN service to the customer carrier network using its own VRF table enabled on the CsC-PE. One important distinction between normal MPLS VPN service and CsC is the fact that traffic passed between the CsC-CE and CsC-PE is labeled rather than native IP^[3, 4].

The CsC architecture is designed such that the backbone carrier network—the network provider’s NGN network—needs to know only about internal routes within the customer carrier network. This setup allows formation of full “any-to-any” logical connectivity between the customer carrier routers, which in this scenario are the PE routers of the existing BGP/MPLS VPN network providing VPN services to end customers.

Furthermore, the backbone carrier routers themselves do not need to retain route prefix information for the end-customer VPNs connected to the customer carrier network because the end-customer traffic is transported over a second level of VRF tables that bear relevance only to the customer carrier itself, that is, the endpoint CsC-CEs. This *nesting* of MPLS VPN networks emphasises the inherent scalability of the CsC architecture. The CsC backbone carrier is effectively behaving like “proxy” P routers for the customer carrier network.

Figure 3: MPLS VPN “Customer Carrier” Network Connected Across NGN “Backbone Carrier”

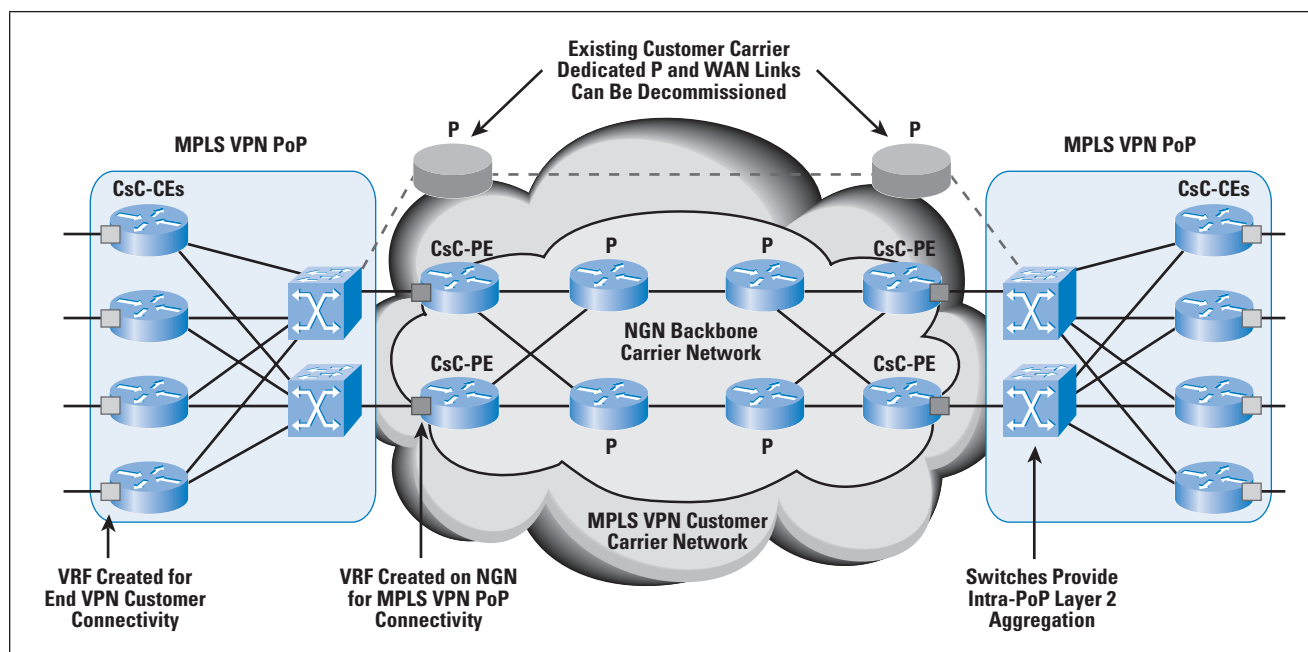


Figure 3 also shows the physical connectivity between the customer carrier network and backbone carrier NGN. Because many large-scale BGP/MPLS network deployments comprise large numbers of PE devices in the same service node or POP, there is often a Layer 2 Ethernet switch acting as an “intra-POP” aggregator. It is convenient to allow physical connectivity between the BGP/MPLS VPN service node and the CsC-PE in the NGN network using this aggregation switch. One or more *Virtual LANs* (VLANs) can be configured across this physical trunk to provide logical Layer 2 connectivity into the CsC-PE on the NGN, and be associated with the CsC VRF on that device. The Layer 2 switch also provides direct intra-POP connectivity between CsC-CEs present on the same VLANs.

Control-Plane Routing Protocols

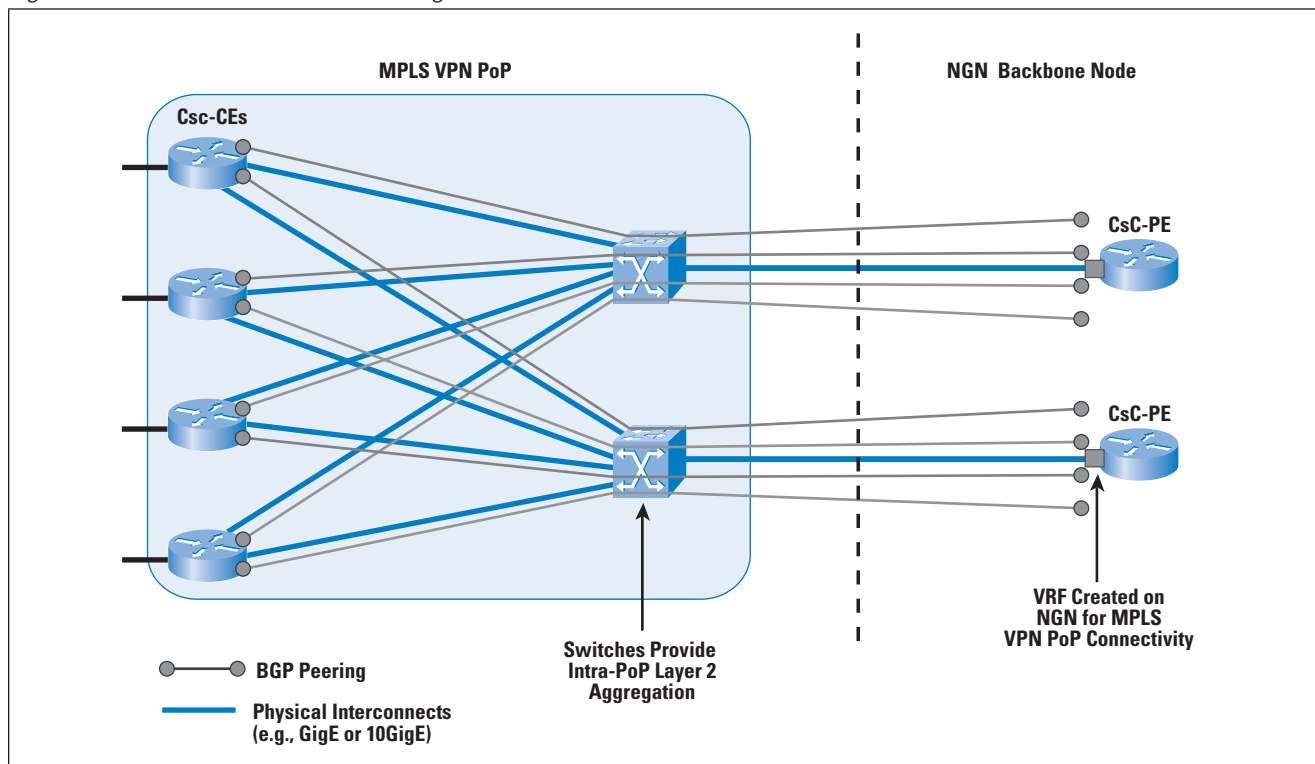
The previous section described the physical connectivity between BGP/MPLS VPN service nodes and the target NGN, with creation of a specific VRF route on the CsC-PEs. This section addresses the way in which the internal routes of the CsC-CEs (that is, the PE routers belonging to the customer carrier BGP/MPLS VPN network) are advertised into this VRF table.

Optional routing protocols include the use of an IGP such as OSPF, or *Exterior Gateway Protocols* (EGPs) such as BGP. With an IGP like OSPF^[5], the routing protocol itself is used for route exchange between the CsC-CEs and CsC-PEs, and must be used in conjunction with an LDP^[6] for MPLS label exchange between the CsC-CEs and CsC-PEs.

Separating the IP prefix and label allocation protocols between an IGP and LDP can introduce complexities with potential divergence between the two control planes. Such divergence in the extreme case can lead to partial or complete loss in forwarding. Use of an EGP like BGP, however, can be used to implement CsC as a single IP prefix and Label Allocation control-plane protocol between CsC-CE and CsC-PE. Piggybacking MPLS label-mapping information in the BGP update messages helps ensure that an IP prefix and its associated MPLS label are always synchronised in their delivery. The way in which this synchronisation is achieved is documented in RFC 3107^[7]. BGP has the benefit of being a mature protocol for use either within the same network organisation or between networks belonging to different operators. Furthermore, BGP employs mechanisms for loop avoidance and control over the number and type of routes advertised and accepted.

Figure 4 shows an example scenario whereby two BGP peerings are established (for resiliency) between each of the four CsC-CEs (which are actually PE routers of the BGP/MPLS VPN customer carrier network) and a pair of target CsC-PE routers (which are the PE routers of the NGN backbone carrier network).

Figure 4: BGP Plus Labels as the Routing Protocol Between CsC-CEs and CsC-PEs

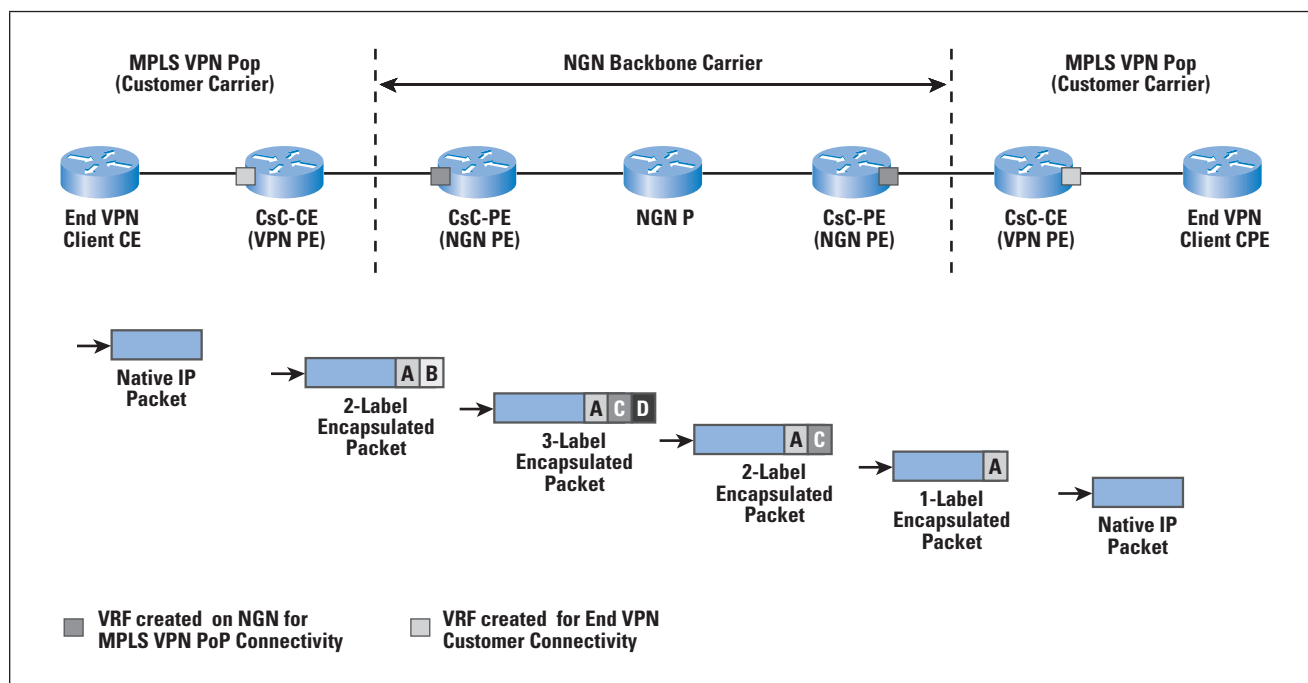


Label Switching of Customer Packets

As shown in Figure 5, viewing packet flow from left to right, a unicast packet originates as a native IP packet when presented from the end client CE router to the MPLS VPN PE router, which is behaving as a CsC-CE in this context. Upon traversal between CsC-CEs in different MPLS VPN POP locations connected by an NGN backbone carrier using CsC, the packet ultimately undergoes three levels of label encapsulation:

- The innermost label corresponds to the *End Customer VRF*. This label is transparent to the NGN backbone carrier (that is, it is not operated upon in lookup and forwarding tables with the NGN). It is label “A” in Figure 5.
- The middle label is the “outer label” as far as the CsC-CE is concerned, swapped at the CsC-PE, and becomes the “inner label” as far as the NGN backbone carrier is concerned. In Figure 5, this label is assigned as label “B” by the CsC-CE as instructed by the CsC-PE through the BGP plus labels (RFC 3107-compliant) peering. At the CsC-PE itself, the label is swapped (to become label “C” in Figure 5) and is used to associate the packet with the CsC VRF. The packet is then identifiable at the destination CsC-PE at the far end of the backbone carrier network; it allows forwarding to the correct interface.
- The outermost label (shown as label “D” in Figure 5) is assigned by the backbone carrier LDP process at the CsC-PE router, and is present only to allow transport across the backbone carrier CsC core. Thus when a packet leaves the CsC-PE for transport across the backbone carrier core it has three levels of labels on each packet.

Figure 5: Label Encapsulation and End-to-End Packet Flow Across a CsC Core Network



As shown in Figure 5, the last P router in the backbone carrier path has “popped” the outermost label (label “D”) using penultimate-hop label forwarding. The destination CsC-PE uses and removes the middle label (label “C”) to indicate the correct outgoing interface, leaving only the innermost label on presentation to the CsC-CE (label “A”). This CsC-CE, which is the PE router in relation to the end VPN services, uses the last remaining label to determine the VRF route and interface on which to send the native IP packet so that it reaches the required client CE router.

Design and Operational Challenges

The previous section outlined the architectural framework of using CsC to integrate one BGP/MPLS core network with another. This section addresses the important design and operational challenges that such a network transformation brings about.

Maintaining Performance Levels

Many existing operators of “carrier-class” BGP/MPLS networks exploit IP QoS mechanisms to allow different IP-based traffic types to be treated in different ways in terms of how the packets are conveyed across the core network. This treatment relates chiefly to prioritisation of delay, jitter, and/or loss-sensitive traffic, against traffic types that are less sensitive to loss or delay. Customers of VPN services supported on such networks generally demand support of a range of traffic types, including corporate intranet, transactional applications, mail services, data backup, video, and VoIP telephony.

To deal with the range of traffic types, BGP/MPLS VPN service providers have developed the means of supporting IP QoS defining different transport classes with associated service levels. One such example may map, for instance, six service classes based on IETF “Per-Hop Behaviours” as defined by the *Differentiated Services* (DiffServ) working group^[8, 9] and the recommended *DiffServ Code Point* (DSCP) values for them. The classes in this example could be broadly described as follows:

- *Expedited Forwarding* (EF), designed and optimised for the delivery of jitter and delay-sensitive applications such as VoIP
- *Assured Forwarding* (AF), intended to support priority data applications; the AF class is split into four equivalent sub-classes (AF1–AF4) used to segregate data or video traffic applications, with priority being maintained over the Default class
- *Default* (DE), to support “best-effort” (that is, unprioritized) data traffic

The DSCP markings dictate the way in which such traffic is placed into queues and conveyed across the core network. At the edge of the MPLS core, the PE maps the incoming DSCP value into the MPLS *Class-of-Service* (CoS) bits (formerly known as EXP bits).

The details of the mapping relate to the specific implementation and policy of the service provider. Under heavy traffic load and congestion situations, such policies dictate how packets are treated in terms of scheduling, queuing, and discard eligibility.

Both the existing BGP/MPLS “customer carrier” and the target NGN “backbone carrier” networks already have their own implementation of QoS classes to allow management and prioritisation of multiple traffic types carried across their respective core infrastructures. A significant design challenge that arises with integrating the networks is that a suitable mapping of the QoS schema present on the PE routers of the customer carrier network (the CsC-CEs in earlier diagrams) to the QoS schema supported on the PE routers of the NGN (the CsC-PEs in earlier diagrams) is necessary.

It is imperative that such a mapping not compromise the existing customer experience for VPN services in terms of packet loss, packet delay, and packet jitter (that is, delay variance). Careful design, mapping of the required service levels, and ultimately end-to-end testing of the QoS mappings is therefore necessary to assure the maintenance of performance levels after the networks are integrated with CsC.

Network Resiliency

As described earlier in the article and shown in Figure 2, an existing standalone BGP/MPLS network platform has interconnected POP locations using underlying core transmission infrastructures such as SONET/SDH/*Dense Wavelength-Division Multiplexing* (DWDM). The actual number of WAN circuits deployed, the use of transmission-layer protection mechanisms, and the overall topological connectivity between POPs determine overall levels of network resiliency. In turn, this aspect of the network architecture significantly affects the overall level of service availability to end customers of VPN services.

When the standalone BGP/MPLS network has its existing core topology replaced with that of the NGN backbone carrier, it is very important to consider the levels of resiliency delivered with the new integrated core architecture, compared with the existing standalone arrangement. Critical considerations include:

- The physical connectivity between the serving nodes of the customer carrier and the backbone carrier should avoid single points of failure where possible.
- If the physical connectivity between the customer carrier and backbone carrier requires the use of WAN transmission links because locations are geographically separate, then suitable levels of circuit protection should be employed
- Because the backbone carrier effectively replaces the existing core topology of the customer carrier, the actual way in which backbone carrier nodes are interconnected and levels of WAN transmission protection etc., should be analysed.

All these aspects should be assessed and incorporated into the actual design process such that there is no detrimental effect on overall levels of service availability to the end customer. Service levels can be verified by reliability modeling of the new network topology, and by comparing the results with the reliability data for the existing topology.

Fault Management

There are many facets of monitoring and managing a core BGP/MPLS network in terms of assurance of service, alarm detection and filtering, customer notification of faults, and so on. In a standalone network environment, it is generally the responsibility of a particular operational team to manage faults on the network and provide service continuity during various types of failure scenarios. As shown in Figure 6, this operational function usually covers all core network elements, including PE and P (core) routers, as well as the WAN topology interconnecting the service nodes or “POPs.”

In an integrated core network scenario, however, part of the customer carrier network—the P (core) routers and WAN transmission links, for example—are replaced by the NGN backbone carrier. The NGN backbone carrier has its own operational team with specific processes and systems for carrying out monitoring and management of fault events. A crucial challenge arises in terms of how to realise end-to-end fault management holistically and transparently between customer carrier and backbone carrier networks (Figure 6). Important considerations include:

- The requirement for a clear and unambiguous demarcation between customer carrier and backbone carrier core platforms must be addressed in terms of operational responsibility for specific faults and the hand-over procedures between operational domains.
- The use of existing monitoring tools and systems in both the customer carrier and backbone carrier domains must be assessed to determine whether new interfaces between such systems need to be developed to facilitate the hand-over procedures.

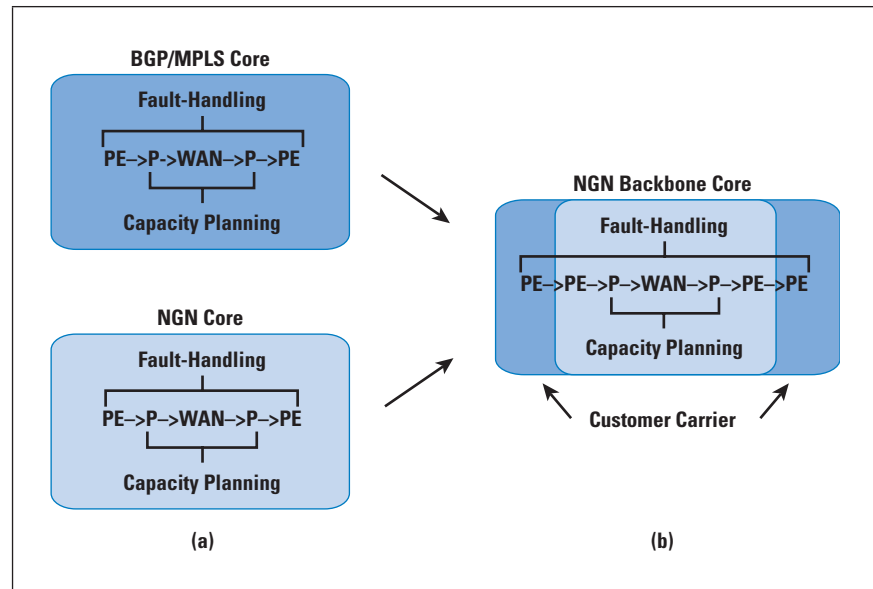
These topics must be factored in to determine the optimal solution for realising smooth and transparent fault-management procedures in an integrated core BGP/MPLS network environment.

Capacity Planning

As shown in Figure 6, in a standalone BGP/MPLS VPN network environment, a particular operational function exists for ongoing core capacity planning to ensure P router and WAN link capacity are suitably dimensioned to cope with current and future traffic demands. When an existing BGP/MPLS VPN network becomes a customer carrier network that is integrated with a target NGN backbone using CsC, there will be a corresponding shift in responsibility for certain aspects of core capacity planning.

VPN service traffic that would have been confined to its own dedicated core network will now be offered onto the NGN backbone carrier core network. As such, the capacity-management function for the NGN backbone carrier must use traffic planning information pertaining to the VPN services in addition to all the other service types supported on the NGN. This aggregated view of traffic demands will accelerate the core capacity dimensioning on the NGN backbone carrier network.

Figure 6: Fault-Management and Capacity-Planning Functions
(a) Before Core Integration
(b) After Core Integration with CsC



Conclusions

The MPLS-based Carrier-supporting-Carrier (CsC) framework provides network operators with a potential solution for integrating an existing BGP/MPLS VPN network, with a target all-IP based NGN. This solution should enable both capital and operational cost reduction by collapsing multiple core networks into a single NGN core domain. The article emphasised that as well as understanding the critical network architectural building blocks required to implement CsC, there are numerous critical design and operational challenges that an integrated core network presents. These challenges include how to maintain service levels and performance metrics for existing VPN customers, resiliency, fault management, and capacity planning. It is important to note, however, that in addition to the broad topic areas covered in this article, many specific additional challenges will present themselves to network operators who have implemented BGP/MPLS VPN networks, and/or NGN networks in their own specific way.

References

- [1] P. Knight and C. Lewis, "Layer 2 and 3 Virtual Private Networks: Taxonomy, Technology, and Standardization Efforts," *IEEE Communications Magazine*, June 2004, pp. 124–131.
- [2] E. Rosen and Y. Rekhter, "BGP/MPLS IP Virtual Private Networks (VPNs)," RFC 4364, February 2006.
- [3] M. Mahmoud, "Carrier-supporting-Carrier: The Whole Story (1)," *Networkers Online*, December 2008.
<http://networkers-online.com/blog/2008/12/carrier-supporting-carrier-the-whole-story-1/>
- [4] M. Mahmoud, "Carrier-supporting-Carrier: The Whole Story (2)," *Networkers Online*, December 2008.
<http://networkers-online.com/blog/2008/12/carrier-supporting-carrier-the-whole-story-2/>
- [5] J. Moy, "OSPF Version 2," RFC 2328, April 1998.
- [6] L. Andersson et al., "LDP Specification," RFC 5036, October 2007.
- [7] Y. Rekhter and E. Rosen, "Carrying Label Information in BGP-4," RFC 3107, January 2001.
- [8] B. Davy et al., "An Expedited Forwarding PHB," RFC 3246, March 2002.
- [9] J. Heinanen et al., "Assured Forwarding PHB Group," RFC 2597, June 1999.

PAUL VEITCH holds an M.Eng. and a Ph.D. from the University of Strathclyde, Glasgow. He joined BT at Martlesham Heath, Ipswich, UK, in September 1996, and worked on various aspects of broadband transmission architectures, multi-service platforms, and 3G network design. In 2000, he joined MCI-WorldCom in Cambridge, UK, and led a number of projects on IP backbone network design. In 2003, he returned to BT to work on IP VPN infrastructure design. He is currently the design authority for BT Retail's Internet networks. He can be reached at: paul.veitch@bt.com

PAUL HITCHEN holds a B.Eng. in Electrical and Electronic Engineering from the University of Salford. He joined BT at Martlesham Heath in September 1990 and has worked on numerous aspects of BT's data services. From 1990 to 1997 he led the development of BT's multiprotocol router portfolio, developing routing and QoS functions with BT's equipment suppliers and provided consulting to BT's customers on IP and Ethernet networks. During the same period he worked on the introduction of Frame Relay, ATM, and SMDS WAN services for BT. In 1997 he developed BT's first IP VPN service offering, working on the development and standardisation of MPLS and VPN technology. From 1997 to the end of 2006 he led the design of BT's Global MPLS Network and service, expanding the network to provide service to more than 150 countries across the world. He is currently a principal consultant working on BT's 21CN IP/MPLS network, focusing on the integration of BT's networks onto 21CN and introducing content delivery and IPTV into the network. He can be reached at: paul.hitchen@bt.com

MARTIN MITCHELL holds an M.Sci. from the University of Bristol and has worked for BT since 2007. He is currently an IP network designer specializing in service provider core design, Ethernet access, and network migrations. He can be reached at: martin.3.mitchell@bt.com

Letter to the Editor

Hi Ole,

I enjoyed the article entitled “PMIPv6: A Network-Based Localized Mobility Management Solution” in the last issue of *The Internet Protocol Journal* (Volume 13, No. 3, September 2010).

I believe that in the “Security Considerations” section it should be mentioned that the CSI (Cga & Send maintenance) working group in the IETF is also working on updating the *Secure Neighbor Discovery* (SEND) specification (RFC 3971) to include the possibility of authenticating the proxied *Neighbor Discovery* (ND) messages sent between the terminal, the *Mobile Access Gateway* (MAG), and the *Local Mobility Anchor* (LMA). This configuration should work in addition to the proposed *IP Security* (IPsec) tunnel between the MAG and the LMA.

The reference material is available at:

<https://datatracker.ietf.org/doc/draft-ietf-csi-proxy-send/>

<https://datatracker.ietf.org/doc/draft-ietf-csi-send-cert/>

Regards,

—Roque Gagliano, Cisco Systems
rogaglia@cisco.com

One of the authors responds:

Dear Ole and Roque,

Thanks for reading our article and providing these valuable comments. We agree with your point. We just considered the basic security mechanisms in our article, limiting the scope to the protocols already standardized, which cover only the protection of the MAG-LMA signaling. We agree that the efforts being carried out within the CSI working group are worth mentioning with regard to the security aspects of PMIPv6.

Thanks,

—Carlos J. Bernardos, Universidad Carlos III de Madrid
cjb@it.uc3m.es

Book Review

A History of the Internet

A History of the Internet and the Digital Future, by Johnny Ryan, Reaktion Books, ISBN 978 1 86189 777 0, September 2010.

Any attempt to document a 50-year history of people and activities that had such a profound and global effect as the Internet faces some challenges. Sequences are complex; written source materials are sketchy; and the many different memories conflict. Added to this reality, of course, are legitimate disagreements about intents and effects. To evaluate such writing effort means first looking for useful criteria. Here are mine: In terms of basic research, was the effort extensive, looking for multiple, appropriate sources and exploring a wide range of probing and constructive questions? Were the sources and questions interesting? This line of thinking leads to a query about the way the author integrates the resulting massive body of data. Is there an effort to develop critical analyses? Are alternative explanations explored?

Johnny Ryan's ambitious *A History of the Internet and the Digital Future* is a rather modest 246 pages, including 28 pages of references. Overall my feeling is that he does quite an interesting job of satisfying the first half of his title, but a somewhat disappointing job with the second half. His research was extensive throughout, but he takes a more critical view of the history than he does of the social aspects of our digital future. In the first half, he integrates information and reports discrepancies and curiosities. In the second half, he indulges in the common, wide-eyed wonderment that technology futurist efforts inherently risk. (Full disclosure: By way of demonstrating the thoroughness of his research, Ryan even included me as one of his many sources.)

Organization

The book is divided into three parts. Broadly, they cover origins, growth, and social effects. Ryan's use of "centrifugal" is contrasted with "centripetal" and is meant to distinguish paradigmatic tensions between approaches that centralize control versus approaches that distribute it. (Oddly, neither of these pivotal terms is in the index.) On page 8 he sets the stage:

"Three characteristics have asserted themselves throughout the Internet's history and will define the digital age to which we must all adjust: The Internet is a centrifugal force, user-driven and open."

By "centrifugal" he means moving outward, away from centralized control. For me, the terminology proved distracting, because I kept hearing my 8th-grade science teacher condescendingly explaining that there is no physics force called centrifugal. Rather it is a perception of the interaction between inertia and centripetal force.

For those with less compulsive (or effective) science teachers, the analogy might prove more helpful, because the design choice really is central to the history of networking. The tension between centralized versus distributed has marked—and continues to mark—much of the development of networking. In fact, I wish Ryan had explored its continuation as much as he explored its effect on origins.

Early History

In general, Ryan presents a narrative with fine-grained detail of the different players who played a critical role in the creation and pursuit of packet switching and then its evolution to link independent networks and technologies^[1]. Efforts to take credit for the former have often become quite public and unseemly; Ryan dissects the play of actors, the essence of their technical ideas, and the details of their activities with documentation and diligence, and even uncovers some discrepancies. He develops a narrative that I found intriguing, enlightening, and credible. What I especially liked was that he explored the organizational milieu in which the activities took place. So we hear of the origins of groups such as the *Advanced Research Projects Agency* (ARPA), Lincoln Labs, and The Rand Corporation; the social and political forces that created them; and the roles they played.

Narrative Arcs

The following is really the strength of this book: It develops narrative arcs about social, political, and organizational environments and the steps taken within them that moved along the path of the Internet. It explores who, when, how, and what, both overall and in detail. At its best, the book provides comparative perspective to help the reader understand what was risky and truly innovative and thereby understand what was really challenging to develop and get adopted. As a minor example, Ryan deserves credit for his exploration and debunking of the media distortions surrounding Al Gore's role and statements concerning the Internet. Strictly speaking, debunking media excesses would not normally seem relevant to a review of the history of a technology, but Ryan uses this example for some consideration of the role of politics in the development of the Internet. The U.S. government could have chosen to assume more control over the Internet; it might have quickly turned it into a telecommunications monopoly, rather than letting it develop through independent market forces.

As would be expected for a story this sweeping, Ryan is sometimes redundant and sometimes inconsistent. Overall, the book would have benefited from more careful editing. So it has a quick reference to the “invention” of e-mail messaging at Bolt Beranek and Newman, but later has a more accurate, detailed account of Ray Tomlinson's 1971 effort, there, to add networking to the *existing* e-mail mechanism. (E-mail messaging was present on the first time-sharing systems of the 1960s, but these systems were standalone services. Tomlinson got them to talk each other.)

Another touchstone I use for discussions of Internet history is the role of the *Computer Science Network* (CSNet), because I worked on that. CSNet served as the forerunner of the larger and more obviously pivotal *National Science Foundation Network* (NSFNet). With NSFNet the Internet developed the ability to support multiple backbones—essential for a truly competitive Internet—and the market-priming creation of regional operational services, from which the seeds of the commercial Internet were sown. Ryan notes the role of CSNet as a kind of market research that led to NSFnet, and in this observation his discussion is notable. But his account of CSNet details is somewhat skewed, because CSNet is cast as having full packet-level connectivity, with e-mail-only telephone-based linkages as a secondary service. In reality full connectivity came later; the original years of CSNet were e-mail-only. Why this fact is important to note—besides overly personal fault-finding—is as a reminder that the accounting efforts for this sort of history are always noisy; the story signal is never pure, even with a diligent effort.

A further touchstone topic is the *Domain Name System* (DNS) and the development of the *Internet Corporation for Assigned Names and Numbers* (ICANN). The interesting part of this saga is later-stage Internet history, and Ryan is relatively sloppy with the details. For example, he muddles what *generic Top-Level Domains* (gTLD) already existed and what new ones were proposed, such as **.com** versus **.biz**; he also muddles the distinction between gTLDs and national domains, such as **.uk**. On the other hand, he certainly captures the continuing tone of controversy that surrounded the development and operation of ICANN, the organization now managing assignment of IP addresses and domain names.

But the most obvious, later-stage touchstone for a history like this one must be the development of the World Wide Web. Ryan gets mixed marks here. He misses the long history of open document publishing that existed even in the earlier *Advanced Research Projects Agency Network* (ARPANET), with “anonymous” FTP, and he misses that the use of *Gopher* predated the web by several years. He also misses just how complete and useful a “dynamically linked document” system Doug Englebart’s NLS (computer) system provided 20 years before the invention of the web^[2]. Hence, he misses the long, historical arc for publishing on the Internet. On the other hand, he does discuss *Gopher* and explores some of the reasons it lost the competition to the web. He focuses on management and intellectual property issues, whereas I tend to consider *Gopher* as having a much poorer cost/benefit mix. *Gopher* was text-only and required going down a potential long lookup tree—quite a few “clicks”—before getting any content. The web is mixed-media and can provide utility to the reader—that is, content—at each step down a lookup path. So the web is more complex to develop than *Gopher*, but it provides enough additional power and better human factors to be worth it.

Ryan's discussion of the commercial explosive growth of the Internet is a good read, including the Dutch tulip market reference and his introduction to some relevant tidbits of economics theory. However, as the book moves into "Web 2.0" and beyond, it provides reasonable descriptions of who did what to create popular new services, but his critical eye largely stops providing serious analysis. Explanations sound more like exuberance than examination. On the other hand, he certainly provides substance to the view that the Internet enables "long-tail" market opportunities to discover and satisfy specialized segments. His discussion of politicians' inventive use of the Internet is nicely concise and integrated. Again, it provides a narrative arc with substance. But his predictions for the future of users as news consumers or as citizens in political processes have too much tone of certitude and positive outcome than is justifiable in my opinion.

Worth Reading

In sum, the book is certainly worth reading. You will likely learn quite a bit, but make sure you read with glasses that have no hint of rose coloring!

References

- [1] Debating which milestone marks "the beginning of the Internet" is a favorite pastime, including among those around during the period in question. Various definitions are legitimate, as long as one is clear about the choice. For me, the operational demonstration of packet switching was when the world changed, so I choose 1969 and the first four nodes of the ARPANET; or its public demonstration in 1972. TCP/IP built on this, by refining and minimizing the work to be done within the infrastructure and by linking independent networks.
- [2] In the early 1970s, my job at UCLA included technical documentation and supporting online use by the Computer Science Department's secretaries. We did all our editing remotely, on the Engelbart system, because it was so powerful.

—Dave Crocker, *Brandenburg Internet Working*
dcrocker@bbiw.net

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the "networking classics." In some cases, we may be able to get a publisher to send you a book for review if you don't have access to it. Contact us at ipj@cisco.com for more information.



Photo: Matsuzaki Yoshinobu

Bjoern A. Zeeb Receives Second Itojun Service Award

The second Itojun Service Award was presented at the 79th meeting of the *Internet Engineering Task Force* (IETF) in Beijing, China. Bjoern A. Zeeb received the award for his dedicated work to make significant improvements in open source implementations of IPv6.

First awarded last year, the *Itojun Service Award* honours the memory of Dr. Jun-ichiro “Itojun” Hagino, who passed away in 2007, aged just 37. The award, established by the friends of Itojun and administered by the *Internet Society* (ISOC), recognises and commemorates the extraordinary dedication exercised by itojun over the course of IPv6 development.

“For many years, Bjoern has been a committed champion of, and contributor to, implementing IPv6 in open source operating systems used in servers, desktops, and embedded computer platforms, including those used by some of the busiest websites in the world,” said Jun Murai of the Itojun Service Award Committee and Founder of the WIDE Project. “On behalf of the Itojun Service Award Committee, I am extremely pleased to present this award to Bjoern for his outstanding work in support of IPv6 development and deployment.”

The Itojun Service Award is focused on pragmatic contributions to developing and deploying IPv6 in the spirit of serving the Internet. The award, expected to be presented annually, includes a presentation crystal, a US\$3,000 honorarium, and a travel grant.

“This is a great honour, and I would like to thank the people who recommended me for the award and the committee for believing my work was valuable. I never met Itojun but he was one of the people helping me, and I have the highest respect for his massive foundational work,” said Bjoern A. Zeeb. “As the Internet community works to roll out IPv6 to more and more people all around the globe, we also need to help others—developers, businesses, and users—understand and use the new Internet protocols so that the vision Itojun was working so hard for comes true.”

Each Internet-connected device uses an IP address and, with the number of Internet-connected devices growing rapidly, the supply of unallocated IPv4 addresses is expected to be exhausted within the next year. To help ensure the continued rapid growth of the Internet, IPv6 provides a huge increase in the number of available addresses. And, while the technical foundations of IPv6 are well established, significant work remains to expand the deployment and use of IPv6.

For more information about the Itojun Service Award see:
<http://www.isoc.org/itojun/>

Remaining IPv4 Address Space Drops Below 5 percent

The *Number Resource Organization* (NRO) recently announced that less than five percent of the world's IPv4 addresses remain unallocated. APNIC, the Regional Internet Registry for the Asia Pacific region, has been assigned two blocks of IPv4 addresses by the *Internet Assigned Numbers Authority* (IANA). This latest allocation means that the IPv4 free pool dipped below 10% in January 2010. Since then, over 200 million IPv4 addresses have been allocated from IANA to the *Regional Internet Registries* (RIRs).

“This is a major milestone in the life of the Internet, and means that allocation of the last blocks of IPv4 to the RIRs is imminent,” stated Axel Pawlik, Chairman of the NRO, the official representative of the five RIRs. “It is critical that all Internet stakeholders take definitive action now to ensure the timely adoption of IPv6.”

IPv6 is the “next generation” of the Internet Protocol, providing a hugely expanded address space, which will allow the Internet to grow into the future. In 2010, the five RIRs are expected to allocate over 2,000 IPv6 address blocks, representing an increase of over 70% on the number of IPv6 allocations in 2009. In contrast, the number of IPv4 allocations is expected to grow by only 8% in 2010. These statistics indicate an absence of any last minute “rush” on IPv4 addresses, and a strong momentum behind the adoption of IPv6.

“The allocation of Internet number resources by the five RIRs enables every region in the world to benefit from fair and equitable distribution of IPv4 and IPv6 addresses. We are also actively collaborating with stakeholders at the local, regional, and global level to offer training and advice to public and private sector organisations on IPv6 adoption to ensure that everyone is prepared for IPv4 depletion and IPv6 deployment,” added Pawlik.

The IANA assigns IPv4 addresses to the RIRs in blocks that equate to 1/256th of the entire IPv4 address space (each block is referred to as a “/8” or “slash-8”). The most recent assignment means that there are now only 12 of these blocks available, which is less than five percent of the entire IPv4 address pool.

The final five blocks of IPv4 addresses will be distributed simultaneously to the five RIRs, leaving only seven blocks to be handed out under the normal distribution method.

According to current depletion rates, the last five IPv4 address blocks will be allocated to the RIRs in early 2011. The pressure to adopt IPv6 is mounting. Many worry that without adequate preparation and action, there will be a chaotic scramble for IPv6, which could increase Internet costs and threaten the stability and security of the global network.

The NRO exists to protect the pool of unallocated Internet numbers (IP addresses and AS numbers) and serves as a coordinating mechanism for the five RIRs to act collectively on matters relating to the interests of RIRs. For further information, visit <http://www.nro.net>

The RIRs are independent, not-for-profit membership organizations that support the infrastructure of the Internet through technical coordination. There are five RIRs in the world today. Currently, the IANA allocates blocks of IP addresses and ASNs, known collectively as *Internet Number Resources*, to the RIRs, who then distribute them to their members within their own specific service regions. RIR members include *Internet Service Providers* (ISPs), telecommunications organizations, large corporations, governments, academic institutions, and industry stakeholders, including end users.

The RIR model of open, transparent participation has proven successful at responding to the rapidly changing Internet environment. Each RIR holds one to two open meetings per year, as well as facilitating online discussion by the community, to allow the open exchange of ideas from the technical community, the business sector, civil society, and government regulators. Each RIR performs a range of critical functions including: The reliable and stable allocation of Internet number resources (IPv4, IPv6 and *Autonomous System Number* resources); The responsible storage and maintenance of this registration data; The provision of an open, publicly accessible database where this data can be accessed. RIRs also provide a range of technical and coordination services for the Internet community. The five RIRs are:

AfriNIC: <http://www.afrinic.net>

APNIC: <http://www.apnic.net>

ARIN: <http://www.arin.net>

LACNIC: <http://www.lacnic.net>

RIPE NCC: <http://www.ripe.net>

The Facebook logo, consisting of the word "facebook" in white lowercase letters on a blue rectangular background.

Find us on Facebook

In addition to *The Internet Protocol Forum*, available at <http://www.ipjforum.org>, IPJ now has its own Facebook page. Join the discussion and get the latest news and updates:

<http://www.facebook.com/#!/pages/Internet-Protocol-Journal/163288673690055>

This publication is distributed on an "as-is" basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2010 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

March 2011

Volume 14, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editors.....	1
Address Exhaustion.....	2
World IPv6 Day	12
Transitional Myths	14
Transitioning Protocols.....	22
Call for Papers.....	47

FROM THE EDITORS

In 2011 we have already seen some important Internet anniversaries and milestones. We have celebrated 25 years of IETF meetings and 40 years of the FTP protocol, but the most significant milestone took place in February when IANA handed out its final blocks of IPv4 addresses to the RIRs (see page 21). It seems like a good time to publish an edition of IPJ devoted entirely to IPv4/IPv6 transition, and to help me with this task I have invited Geoff Huston as co-editor and author for this issue, so let me hand it over to him:

There is a Chinese proverb that states: 寧為太平犬，不做亂世人 “It’s better to be a dog in a peaceful time than be a man in a chaotic period.” For the Internet, this year is shaping up to be a time that looks more like developing chaos than serenity and peace. The IANA has given out the last /8’s, and demand has already depleted the IPv4 address stocks in the Asia Pacific. Meanwhile, the industry has discovered the mass marketing potential of mobile devices, and expects to sell and connect more than 250 million of them in 2011 alone.

The IETF designed IPv6 in the 1990s for this very reason. Its 128-bit address field is easily capable of accommodating the output of a prolific silicon manufacturing industry for many decades to come. But when we look at today’s Internet, very little IPv6 can be seen. Estimates of the number of clients with functional IPv6 services hover at around 0.2 to 0.4 percent of the total.

The story about IPv6 transition technologies is complex, and there are many ways to undertake this effort. In this issue we will examine the various approaches and their relative strengths and weaknesses.

In order to send out a broad message about the need to shift online content from exclusively using IPv4 into a dual-stack world of both IPv4 and IPv6, ISOC is supporting *World IPv6 Day* on June 8. Phil Roberts explains this initiative and its role in helping the overall transition effort.

This transition is going to be difficult. It involves all parts of this diverse industry, and means combining some well-understood and widely-deployed technologies in some surprising and challenging ways. There is much to do, and we hope that this issue of IPJ provides an insight into just what the transition to IPv6 will entail.

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

—Geoff Huston, gih@apnic.net
Chief Scientist, APNIC

—Ole J. Jacobsen, ole@cisco.com
Editor and Publisher, IPJ

A Rough Guide to Address Exhaustion

by Geoff Huston, APNIC

The level of interest in IPv4 address exhaustion seems to be increasing, so I thought I would share some answers to the most common questions I have been asked on this topic in recent times.

What is the most significant challenge to the Internet today?

What a wonderfully open-ended question! There are so many challenges that I could identify: improving the level of security on the network, eradicating spam and viruses, improving capacity of the network infrastructure, improving the efficiency of high-speed data transfer, improving the accuracy of search engines, building more efficient and high-capacity data centers, and reducing the unit cost of Internet services, to name but a few.

If there is a common factor in many of these challenges, it is *scaling* the network to meet an ever-expanding agenda of more users, more devices, more traffic, more services, and more policies. And with more users and more forms of use come higher levels of diversity of use and greater need to replace implicit mechanisms of trust with explicit forms of trust negotiation and greater levels of demonstrable integrity of operation.

But these topics are all tactical in nature. They reflect the “how” of making the network work tomorrow by studying how to undertake marginal improvements on the network of today. However, it is not clear that the networks not just of tomorrow or next year, but a decade or more hence should reflect the usage patterns and user population of today. Perhaps a more fundamental challenge is to understand what is missing in today’s network that we will need in the future.

This discussion leads to a pretty obvious challenge, at least for me. The basic currency of any network is *identifiers*. Identifiers allow the network to distinguish between clients and ensure that conversations occur between those parties who intended to communicate. In the world of packet-switched networking, such as IP, these endpoint identifiers are synonymous with the concept of an *address*. What is missing in today’s network is an abundant supply of new addresses that will allow the network to scale up in size by a further factor of at least 1 million, and hopefully more than a billion-fold.

In fact, the supply of addresses is not just inadequate for future needs for a decade hence. The stock of addresses is facing imminent depletion, and the question of availability of addresses is best phrased in terms of months rather than years.

Perhaps the term “address” is somewhat of a misnomer in this context, but it may well be too late to change that now. The primary role of an IP address is not to uniquely identify the location of an endpoint of a network in relation to some positional or topographical coordinate set, but to simply uniquely identify an endpoint to distinguish it from all other endpoints. Its location is not an intrinsic property of this so-called *address*. But common convention is to call these endpoint identifiers “addresses,” so I will stick with the same convention here.

So my candidate “most significant challenge for the Internet today” is that we are running out of further supply of IP addresses.

What is an IP address, and why is it so important?

One of the revolutionary changes introduced by the so-called *packet-switched* network architecture of the Internet—as compared to its telephone predecessor that used *circuit switching*—was that a massive amount of “intelligence” was ripped out of the network and placed into the devices that connect at the edge.

IP networks are incredibly simple, and at their most basic level they do very little. They are built of routers and interconnecting conduits. The function of a router is quite simple. As a packet arrives at the router from the connected circuitry (or from a wireless interface), it is divided into a common IP header and a payload. The IP header of the packet contains, among other components, two fixed-length fields: the address of the intended *destination* of the packet, and, like a postal envelope, the address of the packet creator, or the *source*. The router uses the destination address of the packet to make a routing decision as to how to dispose of the packet. For each incoming packet, the router inspects the destination address in the packet and either passes it to a connected computer if there is an address match or otherwise passes it down the *default* path to the next router. And that is a working description of the entirety of today’s Internet. The important aspect here is that every connected device must have a unique address. As long as this condition is satisfied, everything else can be made to work.

In the current version of the Internet Protocol, an “address” is a 32-bit field, which can encompass some 4.4 billion unique values.

Why are we running out of addresses?

Blame silicon. Over the past 50 years, the silicon chip industry has graduated from the humble transistor of the 1950s to an astonishing industry in its own right, and the key to this silicon industry is volume.

Individual processor chips may take hundreds of millions of dollars to design, but if fabricated in sufficient volume, each processor chip may take as little as a few dollars to manufacture and distribute. The larger the production run of the silicon die, the lower the unit price of the resultant chip. We currently produce a huge volume of computers every year. In 2008 alone around 10 billion computer processors were manufactured. Although most of these microprocessors are simple 8-bit processors that are used to open doors or run elevators, a sizable proportion are used in devices that support communications, whether it is in laptop computers, smartphones, or even more basic communication applications. Typically we do not invent a new communications protocol for each new application. We recycle. And these days if we want a communications protocol for a particular application, it is easiest to simply embed the IP protocol engine onto the chip. The protocol is cheap, well tested, and it works across almost any scale we can imagine from a couple of bits per second to a couple of billion bits per second.

So it is not just the entire human population of the planet who may well have a desire to access the Internet in the future, but equally important is the emerging world of “things” that communicate. Whether it is the latest fashion in mobile phones or more mundane consumer electronics devices such as televisions or games consoles, all these devices want to communicate, and to communicate they need to have a unique identification code to present to the network, or, an “address.”

We are presently turning on more than 200 million new Internet services every year, and today we have used up most of the 4.4 billion addresses that are encompassed by the IP protocol.

When will we run out?

As of September 2010, some 151 million addresses were left in the general-use pool of unallocated addresses that are managed by the central pool administration, the *Internet Assigned Numbers Authority* (IANA). The world’s IP address consumption rate peaked earlier this year at a new all-time high of an equivalent rate of 243 million addresses per year.

By early February 2011 IANA handed out its last address blocks to the RIRs.

The five *Regional Internet Registries* (RIRs)^[1] still had pools of addresses available for general use at that time, but from that point, as they further run down their local pools, the IANA is now unable to provide any more addresses to replenish them. The Asia Pacific Regional Registry, APNIC, has been experiencing the highest level of demand in the world, accounting for some two-thirds of all addresses consumed in early 2011. APNIC exhausted its general use IPv4 address pool in April 2011.

Although the current models of address consumption show that the other regions will be able to manage available address pools for a few more months, this prediction does not account for the multinational nature of many of the largest of the service providers, and at this stage it is not known how much address-consumption pressure will shift outward from APNIC to the other RIRs now that APNIC's available address pool is effectively drained. So it may well be that 2011 will see IPv4 addresses cease to be generally available in many parts of the world, and by early 2012 there will be no further generally available IPv4 addresses in Europe, North America and Asia.

What is the plan?

This news of imminent exhaustion of the supply of addresses is not a surprise. Although the exact date of predicted address exhaustion has varied over time, the prospect of address exhaustion was first raised in technical circles in August 1990, and work has been undertaken since that time to understand what might be possible and how that could be achieved.

The 1990s saw an intense burst of engineering activity that was intended to provide a solution for this forthcoming address problem. The most significant outcome of this effort was the specification of a successor IP protocol to that of IPv4, called IP Version 6 or *IPv6*.

Why IPv6 and not IPv5?

It would be reasonable to expect the successor protocol of IP Version 4 to be called IP Version 5, but as it turned out Version 5 of the Internet Protocol Family was already taken. In the late 1980s the Internet Protocol itself was the topic of a considerable level of research, as researchers experimented with different forms of network behavior. Version 5 of the Internet Protocol was reserved for use with an experimental IP protocol, the *Internet Stream Protocol, Version 2* (ST-II), written up as RFC 1190 in 1990. When it came time to assign a protocol number of the “next generation” of IPv4, the next available version number was 6, hence IPv6.

The outcome of this process was a relatively conservative change to the IP protocol. The major shift was to enlarge the address fields from 32 bits to 128 bits in length. Other changes were made that were thought to be minor improvements at the time, although hindsight has managed to raise some doubts about that!

The design intent of IPv6 is a usable lifetime of more than 50 years, as compared with a “mainstream” deployment lifetime of IPv4 of 15 years, assuming that you are prepared to draw a line at around 1995 and claim that at that time the protocol moved from an interesting academic and research project to a mainstream pillar of the global communications industry.

That 50 years of usable life for IPv6 is admittedly very ambitious, because it is intended to encompass a growth of the ubiquity of silicon from the current industry volumes of hundreds of millions of new connected devices every year to a future level of activity that may encompass in the order of hundreds of billions to possibly some trillions of new connected devices every year.

So the technical plan to address the address-exhaustion problem was to perform an upgrade of the Internet and convert the Internet from IP Version 4 to IP Version 6.

Nothing else needs to be changed. This change is not intended to be radical or revolutionary. The change from circuit switching to packet switching was a revolutionary change for both the communications industry itself and for you and me as enthusiastic communicators. The change from IPv4 to IPv6 is intended to be a polar opposite, and at best it is intended to be a transparent and largely invisible transition. E-mail will still be e-mail. The web should still look just as it always did, and anything that works on IPv4 is expected to work on IPv6. IPv6 is not inherently any faster, nor any cheaper, nor is it even all that much better. The major change in IPv6 is that it supports a much larger address field.

How many addresses are in IPv6?

In theory, there are 2 to the power 128 unique addresses in IPv6—a very large number. If each IPv6 address were a single grain of sand, the entire IPv6 address space would construct 300 million planets, each the size of the earth!

But theory and practice align only in theory. In practice the IPv6 address plan creates a usable span of addresses that encompasses between 2 to the power 50 and 2 to the power 60 devices. Although this number is nowhere near 2 to the power 128, it is still a range of numbers that are between 1 million to 1 billion times the size of the IPv4 address space.

How do we transition to IPv6?

Unfortunately IPv6 is not “backward-compatible” with IPv4. Backward compatibility would allow for a piecemeal transition, where IPv6 could be regarded as a fully functional substitute for IPv4, so that the existing network base would keep using IPv4 forever, while the most recent devices would use IPv6 and all devices could communicate with each other. The lack of such backward compatibility implies that this communication is simply not possible. IPv4 and IPv6 are distinct and different communications protocols, in the same way that English and, say, German are distinct and different languages.

Attempts have been made to design various forms of automated protocol translator units that can take an incoming IPv4 packet and emit a corresponding IPv6 packet in the same manner as a language interpreter. However, this approach also has some major limitations, so it is usable only in very carefully constrained contexts.

The implication of this lack of backward compatibility and inability to perform automated translation within the network is that if we want to preserve comprehensive any-to-any connectivity during the transition, we have to equip each device that is performing a transition with both protocol stacks, or, in effect, allow the device to become “bilingual,” and conduct a conversation in either IPv4 or IPv6, as required. This transition has been termed a *dual-stack* transition.

When my computer supports IPv6, can I return my IPv4 address?

Each device needs to maintain its capability to converse using IPv4 while there are still other devices out there that remain IPv4-only. So a device that becomes IPv6-capable cannot immediately give up its IPv4 address. It will need to keep this IPv4 capability and operate in dual-stack mode for as long as there are other devices and services out there that are reachable only using IPv4.

The implication of this constraint is that we will need to add dual-stack devices to the Internet and consume both IPv4 and IPv6 addresses during this transition.

So, no, you will need to keep your IPv4 address for as long as there are folk out there with whom you want to communicate who have not also migrated to be a dual-stack IPv4- and IPv6-capable entity.

What needs to be done to transition the network to IPv6?

What is encompassed in “transition?” Do all *Internet Service Providers* (ISPs) have to decide when and how to reprogram their systems and reconfigure their routers, switches, and middleware? Will they need to replace all their customers’ modems with ones that support IPv6? What is the agenda?

This level of uncertainty about the transition to IPv6 is evidently widespread in today’s Internet. Most of the actors in the Internet are unsure about what needs to be done, from the largest of the service providers down to individual end users. Yes, it appears to be a simple matter of reprogramming devices from being just IPv4-capable to being capable of supporting both IPv4 and IPv6, but it is not quite so simple. Dual-stack operation is not easy, nor will it just happen without any form of applied impetus. Imagine that this transition is from everyone on the planet speaking Latin to each other to everyone speaking Esperanto. If this situation were a simple matter of everyone stopping using one language and being rebooted to use the other language one by one, then imagine the plight of the first people to undertake this transition—from being connected and being able to communicate with everyone else using Latin, these first users would find themselves speaking exclusively Esperanto to ... nobody! They would in effect have been disconnected from the network.

So the transition is a little trickier than just turning a big switch from IPv4 to IPv6. Because this transition is a piecemeal and fragmentary one, each device, each router, each firewall, each load server, and all those other components of the network service platform need to be programmed with an additional protocol, and become, in effect, bilingual. And in this case there are no magic interpreters that can “translate” between IPv4 and IPv6. So it is only when the entire network is bilingual in a dual-stack mode that we can turn off IPv4 and consider the transition to be complete.

For an extended period of time the Internet is going to have to operate as two Internets. We have never tried that type of operation before, at least not on a grand scale as this one; in fact, it has often been likened to replacing the jet engines of an airplane while the plane is in flight. Somehow we now have to not only sustain a growth rate of at least some 250 million new connections per year, but at the same time retrofit IPv6 to the existing installed base while continuing to support IPv4. The complexity of this operation is significant, and there is considerable confusion about what to do, when to do it, how much it will all cost, and who will pay. So yes, we are all unsure about what needs to be done.

How long do we expect this dual-stack transition to take?

If only we knew! The Internet today encompasses some 1.7 billion users, and hundreds of millions of devices out there are configured to “talk” only IPv4. Some of these devices will surely die in the coming years, and others may be upgraded or reprogrammed, but others will persist in operation for many years to come while continuing to speak only IPv4. Even looking at what is being sold today, although many general-purpose computers (or at least their operating systems) are now configured to operate in dual-stack mode, when you look at embedded devices such as *Digital Subscriber Line* (DSL) or cable modems, or firewalls, or a myriad of other devices that are integral to the operation of today’s Internet, many of these devices are still configured in firmware to operate exclusively using IPv4.

Some modeling of the transition process has projected an 80-year transition process. That projection is heading into the realms of the absurd, given that our expectations for the operational lifespan of IPv6 have a lower bound of just 50 years or so. However, given the sheer scope of the conversion task and the current level of penetration of IPv6 to levels of between 2 and 5 percent of today’s Internet, and given that a deadline of 2 years from now implies a conversion rate of in excess of 1 million devices every day in that 2-year span. It seems that an expectation that this transition could be substantially completed in as little as 2 years also strikes an unrealistic note.

So a more realistic assumption is that we will probably take around 5 years to complete this transition, and we will need to operate the Internet in dual-stack mode with both IPv4 and IPv6 across this entire period.

But at the current level of Internet growth, the IPv4 address pool cannot sustain a further 5 years of growth—at least not with the current amount of unallocated addresses remaining in the allocation pools. The current address-consumption rate is some 250 million addresses per year. The depleted IPv4 address pool simply cannot withstand the pressures of a 5-year transition without a radical change to the model of the IPv4 network. And if we need to rework the model of the IPv4 network simply to sustain a transition to IPv6, then can't we simply get going with IPv6 a little more quickly instead?

However, “fully depleted” or even “run out” is perhaps not the most appropriate way to describe what will happen to IPv4 addresses in the coming months. It is probably more accurate to say “unobtainable at the current prices.” When the current orderly process of allocation of IPv4 addresses comes to an end, that does not mean that IPv4 addresses will be completely unobtainable. In this world many things that are scarce are still obtainable—for a price. It is quite reasonable to anticipate that for as long as there is still a demand for IPv4 addresses there will be some form of “aftermarket” where addresses are traded for money. However, as with many markets, what is not possible to predict is the price for addresses that will be established by such a market-based address-trading regime.

What about “address sharing” in IPv4?

Why do we need IPv6, given that we could simply share addresses in IPv4?

Yes, of course address sharing^[2] is an option, and we have been doing it for many years already in IPv4. But is it a viable substitute for IPv6?

As part of the engineering effort to develop a successor protocol to IPv4 in the mid 1990s, the IETF published a novel approach of *address sharing*, which we call today *Network Address Translation*, or NAT.^[3] These days almost every DSL modem, and other forms of customer connection equipment, comes equipped with NAT functions. Today most Internet Service Providers give their subscribers a single IPv4 address. At home I have a single IPv4 address, and you probably do too. But in my home I have about 20 connected devices of various sorts (I am counting TiVo units, game consoles, televisions, printers, and such, because they are all in essence Internet-connected devices, and I believe that my situation is not unusual). All these devices “share” the single external IP connection, so all of them “share” this single IPv4 address.

But address sharing has its limitations. When a single household shares a single address, nothing unusual happens. But If I were to try to do the same address-sharing trick of using a single IP address to share across, say 2,000 customers, I would cross over into a world of pain. Many applications today gain speed through parallelism, and they support parallelism through consuming port addresses.

Each IP address can support the parallel operation of 65,535 sessions, using a 16-bit *port identifier* as the distinguishing identifier. But when address sharing is used, these ports are shared across the number of devices that are all sharing this common address. When 2,000 customers are sharing a single address and each customer has some 20 or so devices, then the average number of port addresses per device is 1.5. Common applications that exploit parallel operation include such favorites as *Gmail*, *Google Maps*, and *iTunes*. With a sufficiently constrained number of available ports to use, these applications would cease to work. Indeed, many network applications would fail, and at a level of a single address shared across 2,000 households, I would guess that up to half of these 2,000 customers would not have a working Internet at any single point in time.

Our experience suggests that address sharing works only up to a point, and then it breaks everything badly. We are already address sharing at the level of sharing a single address per household, and households are these days buying more connected devices of various sorts, not *fewer*. So attempting to share that single address across more than one household is at best a temporary solution, and is not a sustainable option that is an alternative to IPv6.

So we need to transition to IPv6, and we need to do so within an impossibly short time.

This discussion all sounds like a terrible problem.

Was this global “experiment” with the Internet all one big mistake?

Should we have looked elsewhere for a networking technology back in the 1990s?

The IP address problem is—for me at any rate—a fascinating one. At the time when researchers were working on the specifications for the Internet Protocol in the 1970s, they decided to use fixed-length 32-bit fields of the interface identifier addresses in the protocol. This decision was a radical one at the time. Contemporary network protocols, such as *DECnet Phase III*, used 16-bit address lengths, and 8-bit addresses were also very common at the time. After all, computers were so big and expensive, who could possibly afford more than 256 unique devices in a single network? Eight bits for addresses was surely enough! Using 32 bits in the address field was not an easy decision to make, because there was constant pressure to reduce the packet headers in order to leave more room for the data payload, so to reserve such a massive amount of space in the address fields of the protocol header to allow two 32-bit address fields was a very bold decision.

However, it was a decision that has proved to be very robust. TCP/IP has sustained the Internet from a mere handful of warehouse-sized computers running at mere kilobits per second to today, where probably more than 3 billion devices connect to the Internet in one way or another, at speeds that range from a few hundred bits per second to a massive 100 Gbps—all talking one single protocol that was invented more than 30 years ago.

IP has demonstrated a scale factor 1 billion! In my mind that achievement demonstrates a level of engineering foresight that is truly phenomenal. So in some sense the underlying observation here is not that IPv4 is running out of addresses today, but that it has been able to get to today at all!

Given that IPv4 has been able to scale by a factor of 1 billion, then if we can make IPv6 scale by a further factor of 1 billion from today we will have done well.

Disclaimer

The views expressed in this article do not necessarily represent the views or positions of the *Asia Pacific Network Information Centre* (APNIC).

References

- [1] Daniel Karrenberg, Gerard Ross, Paul Wilson, and Leslie Nobile, "Development of the Regional Internet Registry System," *The Internet Protocol Journal*, Volume 4, No. 4, December 2001.
- [2] Geoff Huston, "NAT++: Address Sharing in IPv4," *The Internet Protocol Journal*, Volume 13, No. 2, June 2010.
- [3] Geoff Huston, "Anatomy: A Look inside Network Address Translators," *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005; he served on the Board of Trustees of the Internet Society from 1992 until 2001. E-mail: gih@apnic.net

World IPv6 Day

by Phil Roberts, ISOC

On June 8, 2011, websites including *Google*, *Facebook*, *Yahoo!*, and *Bing* will make their main webpages reachable over IPv6 for a 24-hour period from 00:00 to 23:59 *Coordinated Universal Time* (UTC). This activity, *World IPv6 Day*, a “test flight” of IPv6, is motivating organizations across the Internet industry to prepare their services for IPv6, the next generation of the Internet Protocol. Internet Service Providers, hardware makers, operation system and application vendors, and other websites are indeed working to make this activity of testing IPv6 on an Internet scale successful.

The Internet is a never-ending exercise in collaboration. Making a successful transition to IPv6 is one of the major challenges facing the Internet today. Although IPv6 is used extensively in many large networks today, the World IPv6 Day activity is acting as a focal point to bring together all parts of the Internet industry to accelerate deployment of IPv6 in all parts of the Internet.

For some time the deployment of IPv6 has faced a “chicken-and-egg problem.” Website owners have been reluctant to deploy IPv6 because there were not many end users to view their webpages over IPv6. Network operators have been hesitant to deploy IPv6 for many end users because there were few places for those users to view content over IPv6. That the most popular websites in the world according to Alexa rankings are deploying IPv6 on their main webpages is a clear indication that the Internet industry is moving beyond this long-standing impasse. Although June 8 is a 24-hour test, it is clear that this is a move toward regular operation of IPv6, and network operators can confidently roll out IPv6 to end users knowing that the Internet industry is making a concerted effort to make IPv6 an operational reality.

Today, IPv6 connectivity concerns provide another disincentive for a major website to enable IPv6 for regular operation. Badly configured or poorly behaving implementations may prevent end users from reaching a major website that enables IPv6 on its main page. It is currently estimated that this problem will affect only a minor percentage of end users—at the time of the announcement of World IPv6 Day, the estimate was that only 0.05 percent of end users would experience difficulties.

Although this percentage is small, it is potentially a very large number of end users for a website that has visitors numbering in the tens of millions (or more). It is simply impossible from a business point of view for a website of this magnitude to deploy IPv6 alone when this many users could be affected. The users who would not be able to get to that website will simply go to another website in search of similar services.

However, because several such websites have agreed to do this testing at the same time, and for the same duration, individual end users who experience disruption of their connectivity by IPv6 may be able to determine that the problem they are experiencing is indeed not a problem with a set of major websites but may, in fact, be a problem in their own host or network, and will provide an incentive for them to take steps to determine the source of the problem and repair it.

Website owners, network operators, and hardware and software vendors are collaborating to minimize these effects leading up to World IPv6 Day. All of these organizations are working to provide tools to detect these problems and offer suggested fixes in advance of June 8. The test site <http://test-ipv6.com/> allows end users today to test their connectivity and determine whether their connectivity to websites will be affected when those websites enable IPv6.

Some websites have already performed a similar 24-hour test. Last year, the German online news site Heise (<http://www.heise.de>) conducted a similar experiment. The site enabled IPv6 on its main page for 24 hours, turned it off, examined the effects of the experiment, and then permanently enabled IPv6 on its main page. Two major websites in Norway did a similar test, and they also have enabled IPv6 permanently. An activity like this for many websites is clearly a step toward regular and normal IPv6 operations. Website owners will, of course, determine when it makes sense for their business to make IPv6 operations available permanently.

Since the announcement of World IPv6 Day, many other websites from around the world have indicated that they are deploying IPv6, and many of those have decided to join in the global IPv6 test on June 8. The list of websites includes major websites such as *Google*, *Facebook*, and *Yahoo!* and very small websites with small numbers of visitors. It is exciting that websites from every inhabited continent plan to participate. Major websites from the Czech Republic, Portugal, Brazil, and Japan, for example, are joining this test, with more websites joining every day.

For further information about World IPv6 Day, please visit:
<http://www.isoc.org/wp/worldipv6day>

There you will find details about the websites that will be turning on IPv6 on June 8, how to join, and information for networks and individuals, including an FAQ.

PHIL ROBERTS joined the Internet Society (ISOC) in 2008. Prior to that he spent several years with Motorola in research and product development, all in the area of mobile broadband systems. He has been active in the IETF for more than a decade. He can be reached at: roberts@isoc.org

Transitional Myths

by Geoff Huston, APNIC

Last October, I attended the *Réseaux IP Européens* (RIPE)^[1] meeting in Rome, and—not unexpectedly for a group that has some interest in IP addresses—the topic of IPv4 address exhaustion, and the related topic of the transition of the network to IPv6, captured a lot of attention throughout the meeting. One session I found particularly interesting was on the transition to IPv6, where people related their experiences and perspectives on the forthcoming transition to IPv6.

I found the session interesting, because it exposed some commonly held beliefs about the transition to IPv6, so I will share them here, and discuss a little about why I find them somewhat fanciful.

Myth 1: “We have many years for this transition.”

No, I don’t think we do!

The Internet is currently growing at a rate that consumes some 200 million IPv4 addresses every year, or 5 percent of the entire address IPv4 pool. This growth rate reflects an underlying growth of service deployment by the same order of magnitude of some hundreds of millions of new services activated per year. Throughout a dual-stack transition, all existing services will continue to require IPv4 addresses, and all new services will also require access to IPv4 addresses. The pool of unallocated addresses was exhausted in February 2011, and the *Regional Internet Registries* (RIRs)^[2] will exhaust their local pools commencing early 2011 and through 2012. When those pools exhaust, then all new Internet services will need access to IPv4 addresses as part of the IPv4 part of the dual-stack environment, but at that point there will be no more freely available addresses from the registries. Service providers have some local stocks of IPv4 addresses, but even those stocks will not last for long.

As the network continues to grow, the pressure to find the equivalent of a further 200 million or more IPv4 addresses each year will become acute—and at some point will be unsustainable. Even with the widespread use of *Network Address Translators* (NATs)^[3] and further incentives to recover all unused public address space, the inexorable pressure of growth will cause unsustainable pressures on the supply of addresses.

It is unlikely that we can sustain 10 more years of network growth using dual stack, so transition will need to happen faster than that. How about 5 years? Even then, at the higher level of growth forecasts, we will still need to flush out the equivalent of 1.5 billion IPv4 addresses from the existing user base to sustain a 5-year transition, and this number seems to be a stretch target. A more realistic estimate of transition time, in terms of accessible IPv4 addresses from recovery operations, is in the 3–4 year timeframe, and no longer.

So no, we do not have many years for this transition. If we are careful—and a bit lucky—we will have about 4 years.

Myth 2: “It is just a change of a protocol code. Users will not see any difference in the transition.”

If only that were true!

In an open market environment, scarcity is invariably reflected in price. For as long as this transition lasts, this industry is going to have to equip new networks and new services with IPv4 addresses, and the greater the scarcity pressure on IPv4 addresses, the greater the scarcity price of IPv4 addresses. Such a price escalation of an essential good is never a desirable outcome, and although numerous possible measures can be taken to mitigate the problem, to some extent or other, the scarcity pressure and the attendant price escalation suggest a reasonable expectation of some level of price pressure on IPv4 addresses.

In addition, an *Internet Service Provider* (ISP) may not be able to rely solely on customer-owned and-operated NATs to locally mask out some of the incremental costs of IPv4 address scarcity. It is likely—and increasingly so the longer the transition takes—that the ISP will also have to operate NATs. The attendant capital and operational costs of such additional network functions will ultimately be borne by the service provider’s customer base during the transition.

But it is not just price that is affected by this transition—network performance may also be affected. Today a connection across the Internet is typically made by using the *Domain Name System* (DNS) to translate a name to an equivalent IP address, and then launching a connection-establishment packet (or the entire query in the case of the *User Datagram Protocol* [UDP]) to the address in question. But such an operation assumes a uniform single protocol. In a transition world you can no longer simply assume that everything is contactable with a single protocol, and it is necessary to extend the DNS query to two queries, one for IPv4 and one for IPv6. The client then needs to select which protocol to use if the DNS returns addresses in both protocols. Then there is the tricky problem of failover. If the initial packet fails to elicit a response within some parameter of retries and timeouts, then the client will attempt to connect using the other protocol with the same set of retries and timeouts. In a dual-stack transitional world, not only does failure take more time to recognize, but even partial failure may take time.

So users may see some changes in the Internet. They may be exposed to higher prices that reflect the higher costs of operating the service, and they may see some instances where the network simply starts to appear “sluggish” in response.

Myth 3: “NAT upon NAT upon NAT will work.”

Maybe. But maybe not all the time, and maybe not in ways that match what happens today.

The Internet has been operating for more than a decade now with a very prevalent model of a single level of address translation in the path. Application designers now assume its existence, and also make some other rather critical assumptions, notably that the NAT is close to the client in a client-server world, and that there is a single NAT in the path, and that its particular form of address translation behavior can be determined with numerous probe tests. There is even a client-to-NAT protocol to assist certain applications to communicate port-binding preferences to the local NAT. In a multilevel NAT world, such assumptions do not directly translate, but it is not necessarily the case that the application is aware of the added NATs in the end-to-end path.

However, it is not just the added complexity of the multipart NAT that presents challenges to applications. The NAT layering is intended to create an environment where a single IP address is dynamically shared across multiple clients, rather than being assigned to a single client at a time. Applications that use parallelism extensively by undertaking concurrent sessions require access to a large pool of available port addresses. Modern web browsers are a classic example of this form of behavior. The multiple NAT model effectively shares a single address across multiple clients by using the port address, effectively placing the pool of port addresses under contention. The higher the density of port contention, the greater the risk that this multiple layering of NATs will have a visible effect on the operation of the application.

There is also a considerable investment in the area of logging and accountability, where individual users of the network are recorded in the various log functions through their public-side address. Sharing these public addresses across multiple clients at the same time—as is the intended outcome of a multilayer NAT environment—implies that the log function is now forced to record operations at the level of port usage and individual transactions. Not only does this reality have implications in terms of the load and volume of logged information, there is also a tangible increase in the level of potential back tracing of individual users’ online activities if full port usage logging were to be instituted, with the attendant concerns that this back tracing represents an inappropriate balance between accountability and traceability and personal privacy. It is also unclear whether there will be opportunity to have any public debate on such a topic, given that the pressure to deploy multilevel NAT is already visible.

Myth 4: “Changing the Customer Premises Equipment (CPE) is easy.”

No, not necessarily.

I think we have all seen many transition plans, including multilevel Version 4 NATs, NATs that perform protocol translation between IPv4 and IPv6, NATs plus tunneling, as in *Dual-Stack Lite*, the *IVI Bi-direction Mapping Gateway*, *6to4*, *6RD*, and *Teredo*, to call up but a few of the various transitional technologies that have been proposed in recent times. (See the article “Transitioning Protocols” starting on page 22.)

All approaches to dual-stack transition necessarily make changes to some part of the network fabric, whether it is changes to the end systems to include an IPv6 protocol stack in addition to an IPv4 stack, or the addition of more NATs, or gateways into the network infrastructure. Of course, within a particular transitional model there is a selective choice as to what elements of the infrastructure are susceptible to change and what elements are resistant to change. Some models of transition, such as *6RD* and *Dual-Stack Lite*, assume that changing the CPE is easy and straightforward, or at least that such a broad set of upgrades to customer equipment is logistically and economically feasible. *6RD* contains an implicit assumption that the network operator has no economic motivation to alter the network elements, and wishes to retain a single protocol infrastructure that uses IPv4.

Where the CPE is owned, operated, and remotely maintained by the service provider, upgrading the image on the CPE might present fewer obstacles than upgrading other elements of the network infrastructure, such as broadband remote-access servers that operate in a single protocol mode, but sweeping generalizations in this industry are unreliable. Service providers tend to operate customized cost models, and appear to be operating with specialized mixes of vendor equipment and operational support systems. For this reason operators tend to have differing perspectives on what component of their network is more malleable, and correspondingly have differing perspectives on which particular transition technology suits their particular environment.

This industry is volume-based, where an underlying homogeneity of the deployed technology—and economies of scale and precision of process—are critical components of reliable and cost-efficient rollouts. It is somewhat unexpected to see this transition expose a relative high degree of customization and diversity in network service environments.

Myth 5: “My ISP has enough IPv4 addresses to last for years, so it does not have a problem.”

Well, not necessarily.

The assumption behind this statement is that everyone else is also able to persist with IPv4, and everyone you wish to reach, and every service point you wish to access, will maintain some form of connectivity in IPv4 indefinitely.

But this assumption is not necessarily valid. At the point in time when a significant number of clients or services cannot be adequately supported on IPv4, then irrespective of how many IPv4 addresses ISPs have, they will need to provide their clients with IPv6 in order to reach these IPv6-only services. On a network, the actions of others directly affect your own local actions. So if you believe that you need do nothing, and you can use an IPv4 service for years into the future, then this position will be inadequate at the point in time when a significant number of others encounter critical levels of scarcity such that they are incapable of sustaining the IPv4 side of a dual-stack deployment, and are forced to deploy an IPv6-only service. The greater the level of address hoarding, the greater the level of pressure to deploy IPv6-only services on the part of those service providers who are badly placed in terms of access to IPv4 addresses.

Myth 6: “We will always have to run IPv4 protocols.”

Probably not.

Or at least not in terms and volumes that are significant to the industry over the forthcoming decades. Protocols do die. DECnet and *Systems Network Architecture* (SNA) no longer exist as widely deployed networking protocols. In particular, networking in the public space is all about any-to-any connectivity, and to support this connectivity we need a common protocol foundation. In terms of the dynamics of transition, this situation is more about tipping points of the mass of the market than it is about sustained coexistence of diverse protocols. When a new technology—or in this case, protocol—achieves a critical level of adoption, the momentum switches from resisting the change to embracing it.

The aftermath of such transitions does not leave a legacy of enduring demand for the superseded technology. As difficult as it is to foresee today, when the industry acknowledges that the new technology achieves this critical mass of adoption, the dynamics of the networking effect propels the industry into a tipping point where the remainder of the transition is likely to be both inevitable and comprehensive. The likely outcome of this situation is that there is no residual significant level of demand for IPv4.

Myth 7: “There is a technology that will translate between IPv4 to IPv6.”

Yes, but...

Such a technology effectively maps between IPv4 and IPv6 addresses. One approach, the *IVI Bi-direction Mapping Gateway*, provides a 1:1 mapping by embedding fields of one address in the other. Another approach, originally termed *Network Address Translator – Protocol Translator* (NAT-PT), uses a mapping table in a fashion similar to a conventional NAT unit. The common constraint here is that if there are no IPv4 addresses, then such a bidirectional mapping cannot be sustained in each approach. Ultimately, if every packet that traverses the public Internet requires public address values in the source and destination fields, and the ISP must provide a protocol bridge between IPv4 and IPv6, then public IPv4 addresses are required.

But it is not just the requirement for continued access to addresses that is the critical concern here. A reading of RFC 4966^[4], “Reasons to Move the Network Address Translator – Protocol Translator (NAT-PT) to Historic Status” should curb any untoward enthusiasm that this approach is capable of sustaining the entire load of this dual-stack transition without any further implications or problems.

Myth 8: “We do not necessarily have to transition to IPv6. There are substitutes.”

Nothing is visible from here!

If we want to continue to operate a network at the price, performance, and functional flexibility that is offered by packet-switched networks, then the search for alternatives to IPv6 is necessarily constrained to a set of technologies that offer approaches that are—at a suitably abstract level—isomorphic to IP. But going from abstract observations to a specific protocol design is never a fast or easy process, and the lessons from the genesis of both IPv4 and IPv6 point to a period of many years of design and progressive refinement to develop a viable approach. In our current context any such redesign is not a viable alternative to IPv6, given the timeframe of IPv4 address exhaustion. It is unlikely that such an effort would elicit a substitute to IPv6, and it is more likely that such an effort may lead toward an inevitable successor to IPv6, if we dare to contemplate networking technologies further into the future.

Other approaches exist, based on application-level gateways and similar forms of mapping of services from one network domain. We have been there before in the chaotic jumble of networks and services that defined much of the 1980s, and it is a past that I for one find easier to forget! Such an outcome is of considerably higher complexity, considerably less secure, harder to use, more expensive to operate, and more resistant to scaling.

Like it or not, the pragmatic observation of today’s situation is that we do not have a viable choice here. No viable substitutes exist.

Myth 9: “We know what is happening.”

I am not sure that is universally true! The comments I have heard about the current situation lead me to the observation that there are many different perspectives on the situation. Individuals perceive the transition in terms that relate to their own circumstances and their own limitations, and a more encompassing perspective of the entire Internet and this transition is harder to assemble. So, from the perspective of the Internet as a whole, no, we are not really aware of what is happening.

Myth 10: “We know what we are doing.”

Individually this statement is, hopefully, true. But at the level of the entirety of the Internet, no, we do not really have a clear perspective of this transition.

Myth 11: “We have a plan!”

See the comment for myth 10.

Myth 12: “The Internet will be fine!”

I am unsure about this one.

The worrying observation is that the Internet has so far thrived on diversity and competition. We have seen constant innovation and evolution on the Internet, and the entrance of new services and new service providers.

But if we rely solely on IPv4 for the future Internet, then this level of competition and diversity will be extremely challenging to sustain. If we lose that impetus of competitive pressure from innovation and creativity, then the Internet will likely stagnate under the oppression of brutal volume economics. The risks of monopoly formation under such conditions are relatively high.

I hope one observation I heard at the RIPE session will be a myth as this transition gets underway:

*“The incumbents will have all the IPv4 space.
Thanks for playing!”*

If that is *not* a myth, then we are going to be in serious trouble!

Disclaimer

The views expressed in this article do not necessarily represent the views or positions of the *Asia Pacific Network Information Centre* (APNIC).

References

- [1] <http://www.ripe.net/ripe/meetings/ripe-61/>
- [2] Daniel Karrenberg, Gerard Ross, Paul Wilson, and Leslie Nobile, “Development of the Regional Internet Registry System,” *The Internet Protocol Journal*, Volume 4, No. 4, December 2001.
- [3] Geoff Huston, “Anatomy: A Look inside Network Address Translators,” *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [4] Cedric Aoun and Elwyn Davies, “Reasons to Move the Network Address Translator – Protocol Translator (NAT-PT) to Historic Status,” RFC 4966, July 2007.

Pool of Unallocated IPv4 Addresses Now Completely Emptied

On February 3, 2011 a critical point in the history of the Internet was reached with the allocation of the last remaining IPv4 Internet addresses from a central pool. It means the future expansion of the Internet is now dependant on the successful global deployment of the next generation of Internet protocol, called IPv6.

The announcement was made by four international non-profit groups, which work collaboratively to coordinate the world’s Internet addressing system and its technical standards. At a news conference in Miami, Florida, the *Internet Corporation for Assigned Names and Numbers* (ICANN) joined the *Number Resources Organization* (NRO), the *Internet Architecture Board* (IAB) and the *Internet Society* (ISOC) in announcing that the pool of first generation Internet addresses has now been completely emptied. The final allocation of Internet addresses was administered by the *Internet Assigned Numbers Authority* (IANA), which is a function of ICANN.

“This is a major turning point in the on-going development of the Internet,” said Rod Beckstrom, ICANN’s President and Chief Executive Officer. “No one was caught off guard by this. The Internet technical community has been planning for IPv4 depletion for some time. But it means the adoption of IPv6 is now of paramount importance, since it will allow the Internet to continue its amazing growth and foster the global innovation we’ve all come to expect.”

Two “blocks” of the dwindling number of IPv4 addresses—about 33 million of them—were allocated in late January to APNIC, the *Regional Internet Registry* (RIR) for the Asia Pacific region. When that happened, it meant the pool of IPv4 addresses had been depleted to a point where a global policy was triggered to immediately allocate the remaining small pool of addresses *equally* among the five global RIRs.

“It’s only a matter of time before the RIRs and *Internet Service Providers* (ISPs) must start denying requests for IPv4 address space,” said Raúl Echeberría, Chairman of the NRO, the umbrella organization of the five RIRs. “Deploying IPv6 is now a requirement, not an option.”

Transitioning Protocols

by Geoff Huston, APNIC

In the previous article, I looked at some common myths associated with the transition to IPv6. In this article I would like to look behind the various opinions and perspectives about this transition, and examine in a little more detail the nature of the technologies being proposed to support the transition to IPv6.

After some time of hearing dire warnings about the imminent exhaustion of the stocks of available IPv4 address space, we have now achieved the first milestone of address exhaustion, the depletion of the central pool of *Internet Assigned Numbers Authority* (IANA)-managed address space. The last five /8s were handed out from IANA to the *Regional Internet Registries* (RIRs) on February 3, 2011. After some years of industrywide general inattention and inaction with IPv6, perhaps it is not unexpected to now see a panicked response along the lines of “Maybe we should do something now!”

But what exactly should be done? It is one thing to decide to “support” IPv6 in a network, but quite another to develop a specific plan, complete with specific technologies, timelines, costs, vendors, and a realistic assessment of the incremental risks and opportunities. Although working through some of this detail has the normal levels of uncertainty that you would expect to see in any environment that is undergoing constant change and evolution, an additional level of uncertainty here is a by-product of the technology itself.

There is not just *one* approach to adding support for IPv6 in your network, but *many*. And it is not just one major objective you need to address—incremental deployment of IPv6 as a second protocol into your operational network without causing undue disruption to existing services—but two, because the second challenging objective is how to fuel continued growth in your network service platform when the current supply lines of readily available IPv4 addresses are effectively exhausted.

When?

The most common question I have heard recently is: “How long do we have?”

The remaining pools of IPv4 address space continue to be drawn down. At the start of February 2011, the IANA pool was fully depleted, with the final allocation to the RIRs^[1] of IPv4 addresses.

Using a model based on monthly address demands now predicts that the next 18 months or so will see the first three RIRs depleted of IPv4 addresses.

The *Asia Pacific Network Information Centre* (APNIC) was the first RIR to exhaust its available pool of IPv4 addresses in April 2011, with the *RIPE Network Coordination Centre* (RIPE NCC) predicted to follow in late 2011 and the *American Registry for Internet Numbers* (ARIN) in early 2012. The *Latin American and Caribbean Internet Addresses Registry* (LACNIC) is predicted to follow in 2014, and the *African Network Information Centre* (AFRINIC) in 2016.

The good news is that many people have been busy thinking about these intertwined objectives of extending the useful lifetime of IPv4 in the Internet and simultaneously undertaking the IPv6 transition, and there is a wealth of possible measures you can take, and a broad collection of technologies you can use. Fortunately, we are indeed spoiled with choices here!

The not-so-good news is that there is no simple single path to follow. Each individual network needs to carefully consider the transition and select an approach that matches their particular circumstances. For an industry used to playing “follow the leader” for many years, a variety of choice is not always appreciated. And, unfortunately, we are spoiled for choices here.

Let’s look at each of the major transitional technologies that are currently in vogue, and examine their respective strengths and weaknesses and their intended area of applicability. We will look at these technologies first from the perspective of the end user and then from the other side, examining options for *Internet Service Providers* (ISPs).

The Dual-Stack ISP Client

If your service provider provides a dual-stack service with both IPv6 and IPv4, then your task should be relatively straightforward. If you configure your modem or router with IPv6 in addition to IPv4, you are finished, assuming of course that your local modem or router unit actually supports IPv6—an assumption that may not be valid in many of the older and, unfortunately, many of the currently available devices.

The conventional approach in this form of environment is to use *IPv6 Prefix Delegation*, where the ISP provides the client with an IPv6 prefix, usually a /48 or a /56 IPv6 address prefix, which is then passed into the client network through an *IPv6 Router Advertisement*. Local hosts should be constructed to configure their IPv6 stack automatically, and your system should be connected as a dual-protocol system.

You probably do, however, need to be aware of some caveats, of which the most important is likely to relate to the probable absence of a *Network Address Translation* (NAT)^[2] function in IPv6. Currently most commercial IPv4 Internet services assign a single IP address to each client.

To allow this address to be shared within the client's network, most IPv4 "edge" devices autoconfigure themselves as NAT devices, permitting outgoing connections using the *Transmission Control Protocol* (TCP) or *User Datagram Protocol* (UDP), and allowing some *Internet Control Message Protocol* (ICMP) message types to traverse the NAT, but not much else. For many clients this NAT configuration becomes the default local security framework, because it permits outbound connections through TCP and UDP to be made, but not much else, and permits initiation of no sessions as incoming sessions. With IPv6 the local network is generally configured with an entire subnet, and instead of a NAT, this subnet is directly connected to the Internet.

The local network is then in a mixed situation of being behind a NAT in IPv4, but directly connected to the Internet using IPv6. This asymmetric configuration with respect to IPv4 and IPv6 raises some questions about the effect on the security of your local network. You need to think about adding appropriate filter rules to the gateway IPv6 configuration that performs the same level of access control to your local site that you have already set up with IPv4 and the NAT. The best advice here is to configure some filter rules for IPv6 that limit the extent of exposure of your internal network to the broader Internet to be directly comparable to the configuration you are using with IPv4.

The IPv4-Only ISP Client

Even today, when the IPv4 pools are rapidly depleting, it is really not very common to have an ISP offering dual-stack IPv4 and IPv6 services. Let's look at the more common situation, when your ISP is still offering only IPv4. As an end user, can you still set up some form of IPv6 access?

The answer is "Yes," but you must use tunnels, and the story can get somewhat ugly.

6to4 Tunnels

If you have public IPv4 addresses on your local network, you may elect to configure your local system to use the *6to4 Tunneling Protocol*.

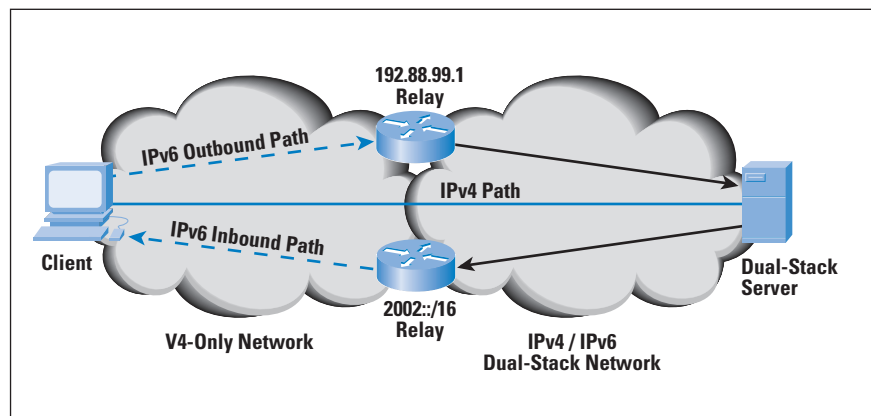
6to4 is an autotunneling protocol coupled with an addressing structure. The IPv6 address of a 6to4-reachable host begins with the IPv6 prefix **2002::/16**. The address architecture embeds a 32-bit IPv4 address of the end host into the next 32 bits. That way the IPv6 address carries the "equivalent" IPv4 address within the IPv6 address.

To send an IPv6 packet, the local host must first tunnel through the local IPv4 network. To perform this tunneling, the local host encapsulates the IPv6 packet in an outer IPv4 packet header. The IP protocol used is neither TCP nor UDP, but protocol 41, an IP protocol number reserved for tunneling IPv6 packets (RFC 2473)^[3].

The IPv4 packet is addressed to an IPv4-to-IPv6 relay. To avoid manual configuration of each client, all these relays share the same *anycast* address, **192.88.99.1**. These relays strip the outer IPv4 packet header off the packet and forward the IPv6 packet into the IPv6 network. The IPv6 destination treats the packet normally, and generates a packet in response without any special processing.

The reverse path to a 6to4 host uses an IPv6-to-IPv4 relay. The IPv6 address of the 6to4 local host started with the IPv6 address prefix **2002::/16**, so the IPv6 packet that is being sent back to this host has a destination address that uses the **2002::/16** 6to4 prefix. This prefix is interpreted as an anycast relay address. A route to the IPv6 **2002::/16** prefix is advertised by IPv6-to-IPv4 relays. When a relay receives a packet destined to a **2002::/16** address, it lifts the IPv4 address from inside the IPv6 address. It then wraps the IPv6 packet in an IPv4 packet header, using as a destination address this extracted IPv4 address, and using protocol 41 as the IP protocol. The resultant IPv4 packet is then passed to the 6to4 host in the IPv4 network (Figure 1).

Figure 1: 6to4 Tunneling Architecture



If the local network has public IPv4 addresses on the local network, then individual hosts on the local network may use 6to4 directly. Of course then the local gateway needs to be configured to accept incoming IP packets that use protocol 41.

An alternative is to configure the gateway device of the local network as a 6to4 gateway, and use the IPv4 address on the ISP side of the gateway as a common 6to4 address for the local network. The gateway then advertises this synthetic 48-bit IPv6 prefix to the interior network with a conventional IPv6 Router Advertisement. The gateway can couple this advertisement with a NAT function and provide native IPv6 to interior hosts that are configured on RFC 1918^[4] local IPv4 addresses.

In general, 6to4 is a relatively poor approach to provisioning IPv6, and you really should avoid it if at all possible. Indeed, your experience will probably be better overall if you continue running IPv4 and avoid accessing IPv6 with 6to4!

The major concern here is that a successful connection relies on the assistance of both an outbound and an inbound 6to4 third-party relay. On the IPv4 side a 6to4 connection relies on the presence of a usable route to a IPv4-to-IPv6 relay, and preferably one that is as close as possible to the IPv4 endpoint. On the IPv6 side a 6to4 connection relies on a usable relay advertising a route to **2002::/16**. Again, to avoid extended path overheads, this relay should be as close as possible to the IPv6 endpoint. This path asymmetry can cause connection “black holes,” where one party can deliver packets to the other but not the reverse.

Also, such configurations have problems if the IPv4 host is configured with stateful filters that insist that the IPv4 source address in incoming packets match the destination address of outgoing packets, not necessarily true in a 6to4 connection.

Finally, it seems that many sites operate with firewall filters that disallow incoming packets other than TCP and UDP (and possibly some forms of ICMP). The 6to4 packets use protocol 41, and there appears to be widespread use of filter rules that block such packets.

Tunneling also adds an additional packet header to a packet, inflating the size of the packet. Such an expansion of the packet on certain path elements of the network may cause path packet size problems, increasing the risk of encountering Path *Maximum Transmission Unit* (MTU) “black holes” due to the increase of the packet size by 20 bytes when the IPv4 packet header is attached to the packet.

Teredo Tunnels

If the local network is behind an IPv4 NAT and the NAT gateway does not support 6to4, then all is not lost, because another form of tunneling could possibly be an answer. *Teredo* is described in RFC 4380^[5].

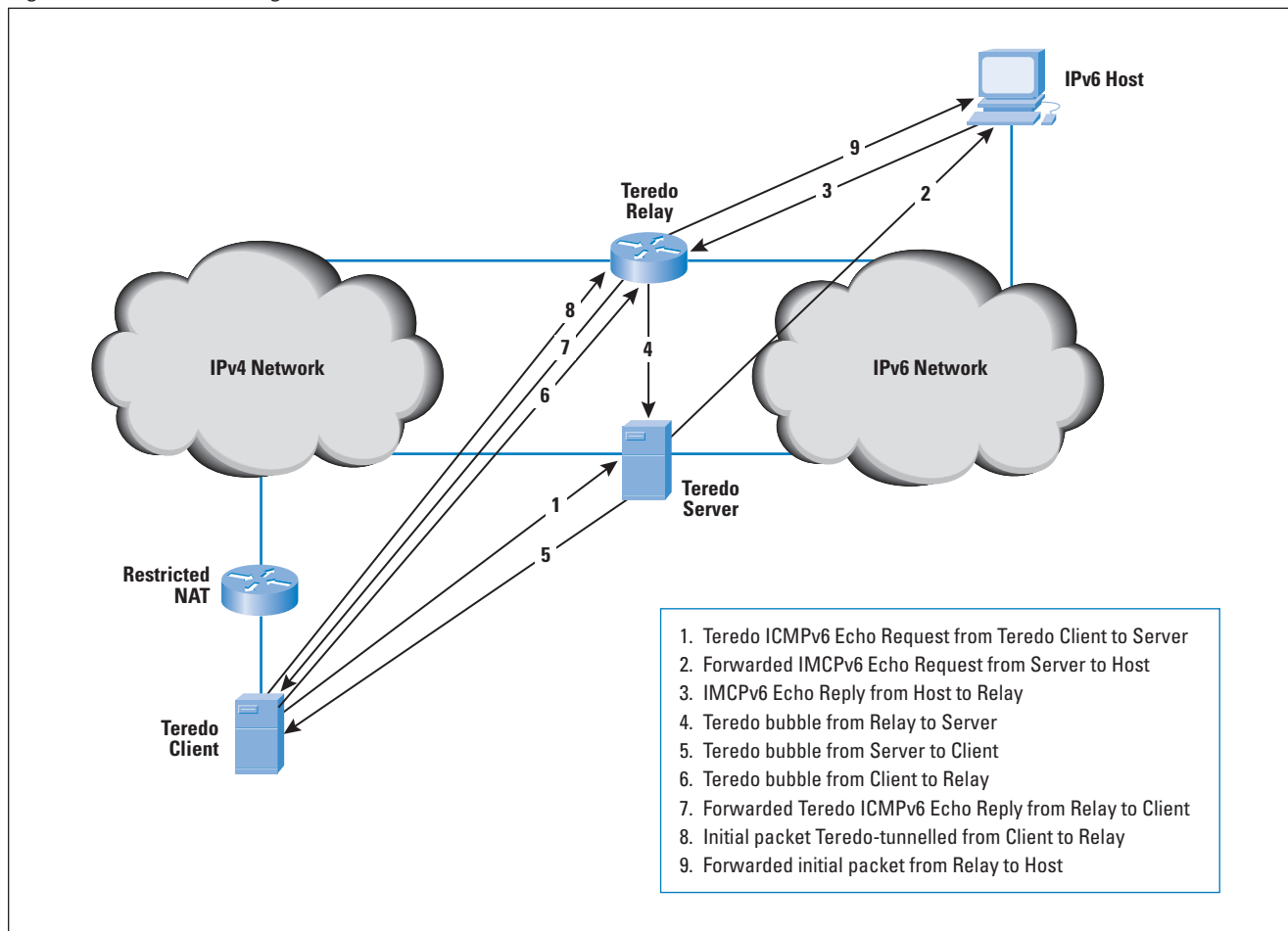
Teredo, like 6to4, is an autotunneling protocol coupled with an addressing structure. Like 6to4, Teredo uses its own address prefix, and all Teredo addresses share a common IPv6 /32 address prefix, namely **2001:0000::/32**. The next 32 bits are the IPv4 address of the Teredo server. The IPv6 interface identifier field is used to support NAT traversal, and it is encoded with the triplet of a field describing the NAT type, the view of the relay of the UDP port number used to reach the client (the external UDP port number used by the NAT binding for the client), and the view of the relay of the IPv4 address used to reach the client (the external IPv4 address used by the NAT binding for the client).

Teredo uses what has become a relatively conventional approach to NAT traversal, using a simplified version of the *Session Traversal Utilities for NAT* (STUN)^[6] active probing approach to determine the type of NAT; it uses concepts of “clients,” “servers,” and “relays.”

A Teredo *client* is a dual-stack host that is located in the IPv4 world, assumed to be located behind a NAT. A Teredo *server* is an address and reachability broker that is located in the public IPv4 Internet, and a Teredo *relay* is a Teredo tunnel endpoint that connects Teredo clients to the IPv6 network. The tunneling protocol used by Teredo is not the simple IPv6-in-IPv4 protocol 41 used by 6to4. NAT devices are sensitive to the transport protocol and generally pass only TCP and UDP transport protocols. In the Teredo case the tunneling is UDP, so all IPv6 Teredo packets are composed of an IPv4 packet header and a UDP transport header, followed by the IPv6 packet as the UDP payload. Teredo uses a combination of ICMPv6^[7] message exchanges to set up a connection and tunneled packets encapsulated using an outer IPv4 header and a UDP header, and it contains the IPv6 packet as a UDP payload.

It should be noted that this reliance on ICMPv6 to complete an initial protocol exchange and confirm that the appropriate NAT bindings have been set up is not a conventional feature of IPv4 or even IPv6, and IPv6 firewalls that routinely discard ICMP messages will disrupt communications with Teredo clients.

Figure 2: Teredo Tunneling



The exact nature of the packet exchange in setting up a Teredo connection depends on the nature of the NAT device that sits in front of the Teredo client. Figure 2 shows an example packet exchange that Teredo uses when the client is behind a Restricted NAT.

Teredo represents a different set of design trade-offs as compared to 6to4. In its desire to be useful in an environment that includes NAT functions in the IPv4 path, Teredo is a per-host connectivity approach, as compared to the 6to4 approach, which can support both individual hosts and entire end sites within the same technology. Also, Teredo is a host-centric multiparty rendezvous application, and Teredo clients require the existence of dual-stack Teredo servers and relays that exist in both the public IPv4 and IPv6 networks. Teredo is more of a connectivity tool than a service solution, and one that is prone to many forms of operational failure.

On the other hand, if you are an isolated IPv6 host behind an IPv4 NAT and you want to access the IPv6 network, then 6to4 is not an option, and you either have to set up static tunnels across the NAT to make it all work or turn on Teredo in your dual-stack host; if everything goes according to theory, you should be able to establish IPv6 connectivity. It is highly likely that the IPv6 Teredo connection will fail in strange ways, and, like 6to4, this is a technology best avoided!

Tunnel Brokers

In contrast to these autotunnel approaches, the simplest form of tunneling IPv6 packets over an IPv4 network is the manually configured IPv6-in-IPv4 tunnel.

Here an IPv6 packet is simply prefixed by a 20-octet IPv4 packet header. In the outer IPv4 packet header, the source address is the IPv4 address of the tunnel ingress, the destination address is the IPv4 address of the tunnel egress, and the IP protocol field uses value 41, indicating that the payload is an IPv6 packet. The packet is passed across the IPv4 network from tunnel ingress to egress using conventional IPv4 packet forwarding, and at the egress point the IPv4 IP packet header is removed and the inner IPv6 packet is routed in an IPv6 network as before. From the IPv6 perspective the transit across the IPv4 network is a single logical hop.

Alternatively, like *Virtual Private Network* (VPN) tunnels, the tunnel can be configured using UDP or TCP, and with some care, the tunnel can be configured through NAT functions in the same way as VPN tunnels can be configured through NAT functions.

The advantage of this approach is that the need to manually configure the tunnel endpoints ensures that the tunnel relay function is not provided, intentionally or unintentionally, by third parties through some well-intentioned, but ultimately random, act of goodwill. The need to perform a manual configuration also reduces the chances that the tunnel will be broken through local firewall filters.

Of course the need to perform a manual configuration does not lend itself to a “plug-and-play” environment, nor is this approach a viable one for a larger mass market of consumer devices and services.

Client Conclusions

None of these approaches to offer IPv6 connectivity to end hosts behind an IPv4-only service provider offers the same level of robustness and performance as native IPv4 services. All of these approaches require a significant degree of local expertise to set up and maintain, and they often require a solid understanding of other aspects of the local environment, such as firewall and filter conditions and Path MTU behavior to maintain. With the exception of the tunnel broker approach, they also require third-party assistance to support the connection, further adding to the set of potential performance and reliability concerns.

It appears that the most robust and reliable way to provision IPv6 to end hosts is for the service provider to provision IPv6 as an integral part of its service offering, and offer clients a dual-stack service in both IPv4 and IPv6.

IPv6 for Internet Service Providers

Although the “self-help” autotunneling approaches for clients outlined earlier in this article are a possible answer, their utility is appropriately restricted to a very small number of end clients who have the necessary technical expertise and who are willing to debug some rather strange resultant potential problems relating to asymmetric paths, third-party relays, potential MTU mismatches, and interactions with filters. This approach is not a reasonable one for the larger Internet.

From the perspective of the mass market for Internet Services, we cannot assume that clients have the motivation, expertise, and means to bypass their ISP and set up IPv6 access on their own, either through autotunneling or manually configured tunnels. The inference from this observation is that for as long as the mass-market ISPs do not commit to IPv6 services, and for as long as they continue to stall in deploying services supporting dual access for their clients, the entire IPv6 transition story remains effectively stalled.

How can ISPs support IPv6 access for their clients?

The Dual-Stack Service Network

Perhaps it is obvious, but the most direct response here is for the ISP to operate a *Dual-Stack Network*.

And the most direct way to achieve this operation is for the ISP's infrastructure to also support IPv6 wherever there is IPv4, so that the delivery of services to the ISP's clients in IPv6 faithfully replicates the service offered in IPv4.

This solution implies that the network needs to support IPv6 in the ISP's routing infrastructure, in the network data plane, in the load-management systems, in the operational support infrastructure, in access and accounting, and in peering and in transit. In short, wherever there is IPv4 there needs to be IPv6.

The infrastructure elements that require dual-stack service at the next level include the routing and switching elements, including the internal and external routing protocols. The task includes negotiating peering and transit services in IPv6 to complement those in IPv4. Network infrastructure also includes VPN support and other forms of tunnels, as well as data center front-end units, including load balancers, filters and firewalls, and various virtualized forms of service provision. The task also includes integration of IPv6 in the network management subsystem and the related network measurement and reporting system. Even a comprehensive audit of the supported *Management Information Bases* (MIBs) in the active elements of the network to ensure that the relevant IPv6 MIBs are supported is an essential task. A similar task is associated with equipping the server infrastructure with IPv6 support, and at the higher levels of the protocol stack are the various applications, including web services, mail, *Domain Name System* (DNS), authentication and accounting, *Voice over IP* (VoIP) servers, Load Balancers, Cloud Servers, and similar applications.

And those are just the common elements of most ISPs' infrastructures. Every ISP also has more specialized elements in its service portfolio, and each one of these elements also requires a comprehensive audit to ensure that there is an IPv6 solution for each of these elements that leads to a comprehensive dual-stack outcome.

As obvious as this approach might appear, it has two significant problems. First, it requires a comprehensive overhaul of every element in the ISP's service network. Even for small-scale ISPs this overhaul is not trivial, and for larger service provider platforms it is an exercise that may take months if not years and make considerable inroads into the operating budgets of the ISPs. Secondly, it still does not account for the inevitable fact that in the coming months the current supply lines of IPv4 addresses will end and any continued expansion of the service platform will require some different approaches to the way in which IPv4 addresses are deployed in the service platform.

Although the approach of simply provisioning IPv6 alongside IPv4 in a simple dual-protocol service infrastructure may appear to be the most obvious response to the need to transition to IPv6, it may not necessarily be the most appropriate response for many ISPs to the dual factors of IPv6 transition and IPv4 address exhaustion.

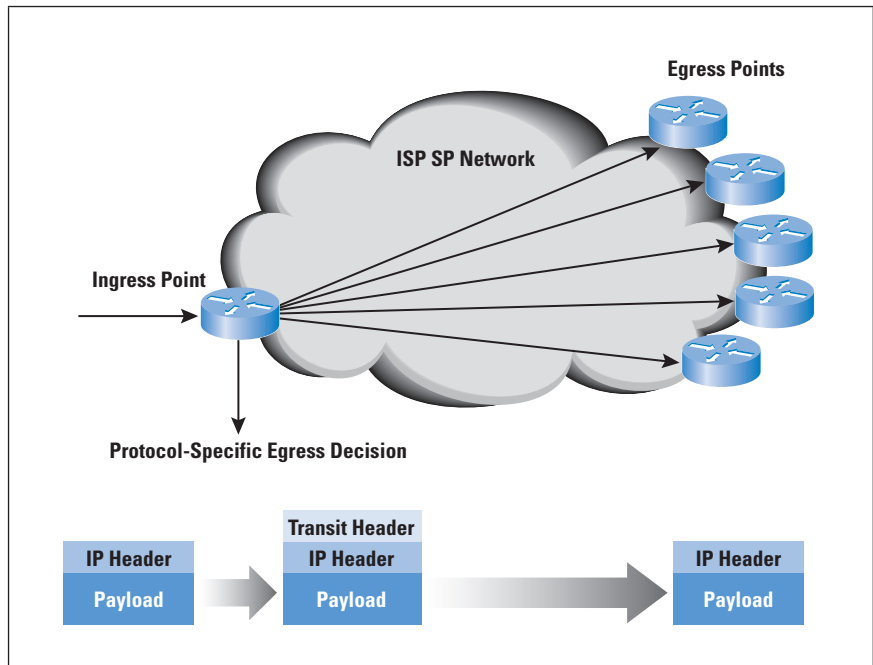
Are there alternative approaches for ISPs? Of course.

Hybrid Approaches

Saying that an ISP must deploy IPv6 across all of its infrastructure and actually doing it are often quite different. The cost of converting all parts of an ISP's operation to run in dual-stack mode can be quite high, and the benefit of running every aspect of an ISP's service offering in dual-stack mode is dubious at best.

Are there middle positions here? Is it possible for an ISP to deliver robust IPv6 services to clients while still operating an IPv4-only internal network? One way to look at an ISP's network is as a transit conduit (Figure 3).

Figure 3: Generic ISP Packet Transit Architecture



The ISP needs to be able to accept packets from an external interface, determine the appropriate egress point for the packet within the context of the local network, and then ensure that the packet is passed out this egress interface. The internal network need not operate in the same protocol context as the protocol of the packets the network is handling. Viewed at a level of the minimal essentials, the network needs to be able to have some protocol-specific capability at its ingress points in order to determine the appropriate egress point of each incoming packet, and thereafter during the transit of the service provider's network, the minimum necessary association to maintain the identity of this preselected egress point with the packet. Now if the network uniformly supports the same protocol as the packet, then the same egress decision can be made at each forwarding point within the network.

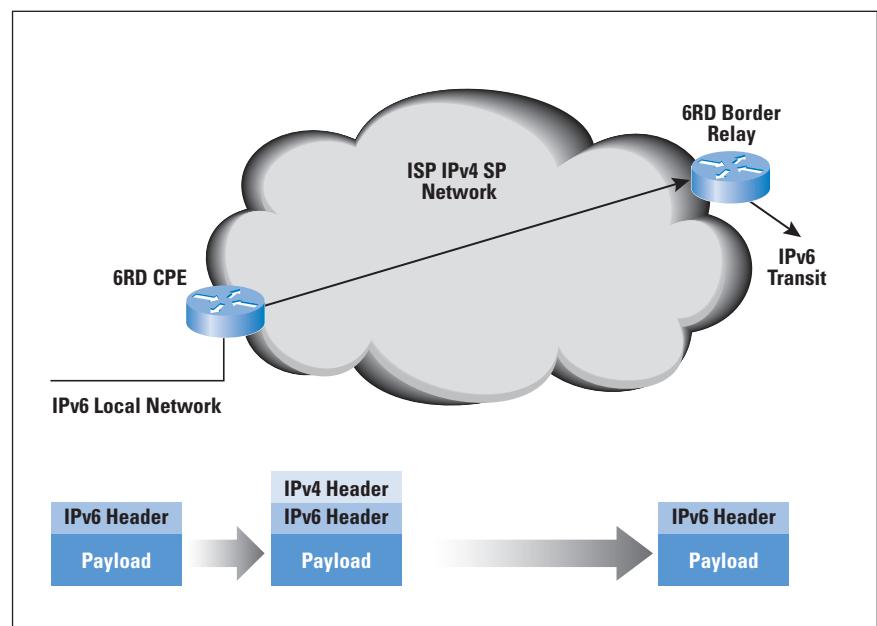
Alternatively, the packet can be encapsulated with an outer wrapper that identifies the egress point using the same protocol context as that used by the service provider's internal switching elements, and the packet can be passed through the service provider's transit network using only this temporary wrapper to determine the sequence of forwarding decisions. *Multiprotocol Label Switching* (MPLS) networks are an excellent example of this form of approach, as are other forms of IP-in-IP encapsulation. The advantage of this approach is that the internal infrastructure of the service provider network need not be altered to support additional carriage protocols: the changes to specifically support IPv6 are required only at the network ingress elements, and a basic encapsulation stripping function is used at all egress points.

With this information in mind, let's look at some of these hybrid approaches to supporting IPv6 in a service provider network.

6RD

6RD, described in RFC 5969^[8], is an interesting refinement of the 6to4 approach. It shares the same basic encapsulation protocol and the same address structure of embedding of the IPv4 tunnel endpoint into the IPv6 address. However, it has removed the concept of third-party relays and the use of the common `2002::/16` IPv6 prefix, and instead uses the provider's IPv6 prefix. The effect of these changes is to limit the scope of the tunneling mechanism to that of tunneling across the network infrastructure of a single provider, and the intended function is to tunnel from the *Customer Premises Equipment* (CPE) to IPv6 *Border Relays* operated by the customer's ISP (Figure 4).

Figure 4: 6RD Tunneling



If 6to4 is not recommended for use because of high failure rates of connections and suboptimal performance, then why would 6RD be any better?

The most compelling reason to believe that 6RD will perform more reliably than 6to4 is that 6RD removes the wild-card third-party relay element from the picture. For outbound traffic the CPE provides the tunnel encapsulation, which is, hopefully, under the ISP's operational control. The IPv6-in-IPv4 tunnel is directed to the ISP's own 6RD Border Relay rather than the 6to4 relay anycast address. Because this process is also under the ISP's direct operational control, it eliminates the outbound third-party relay function. For the reverse path, the use of the provider's own IPv6 prefix in 6RD, instead of the generic `2002::/16` prefix, ensures that the inbound packets are sent through IPv6 directly to the ISP, and the IPv6-in-IPv4 tunnel is again limited to a hop across the ISP's own internal infrastructure.

As long as the ISP effectively manages all CPE devices, and as long as the CPE itself is capable of supporting the configuration of additional functional modules that can deliver unicast IPv6 to the client and 6RD tunnels inward to the ISP, then 6RD is a viable option for the ISP. At the cost of upgrading the CPE set to include 6RD support, and the cost of deployment of 6RD Border Relays that terminate these CPE tunnels, together with IPv6 transit from these Border Relays, the ISP is in a position to provide dual-stack support to its client base from an internal network platform that remains an IPv4 service platform, thereby deferring the process of conversion of its entire network infrastructure base to support IPv6.

For ISPs seeking to defray the internal infrastructure IPv6 conversion costs over a number of years, or for ISPs seeking an incremental path to IPv6 support that allows the existing infrastructure to remain in place temporarily, 6RD can be an interesting and cost-effective alternative to a comprehensive dual-stack deployment, as long as the ISP has some mechanism to load the CPE with IPv6 support and 6RD relay functions.

MPLS and 6PE

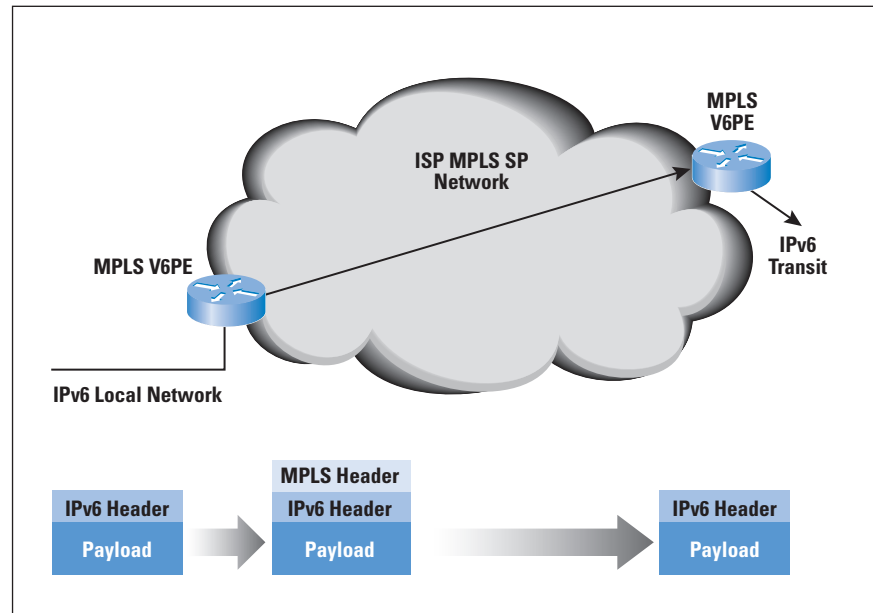
The 6RD approach has many similarities to MPLS, in that an additional header is added to incoming packets at the network boundary, and the encapsulation effectively directs the packet to the appropriate network egress point (as identified by ingress), where the encapsulation is stripped and the original packet is passed out.

Rather than using an IPv4 header to direct a packet from ingress to egress, if the network is already using MPLS, why not simply support IPv6 on an existing MPLS network as a PE-to-PE MPLS path set and bypass the IPv4 step?

Why not, indeed, and RFC 4659^[9] describes how this bypass can be achieved.

If you are running an MPLS network, then the role of the interior routing protocol and label distribution function is to maintain viable paths between all network ingress and egress points. The protocol-specific function in such networks is not the interior network topology management function, but the maintenance of the mapping of egress to protocol-specific destination addresses (Figure 5).

Figure 5: MPLS and 6PE



As with 6RD, if the local problem is some form of prohibitive barrier to the immediate deployment of IPv6 in a dual-stack configuration across the network infrastructure, then this approach allows an IPv4 MPLS network to set up paths across the network IPv4 MPLS infrastructure from provider edge to provider edge. These paths may be used to tunnel IPv6 packets across the network by associating the IPv6 destination address of the incoming packet with the IPv4 address of the egress router, using the *interior Border Gateway Protocol* (iBGP) *Next-Hop* address, for example.

The incremental changes to support IPv6 are constrained to adding IPv6 to the service provider's iBGP routing infrastructure, and to the provider-edge devices in the MPLS network, while all other parts of the service provider's service platform can continue to operate as an MPLS IPv4 network for now.

IPv4 Address Compression

It is not just the challenge of adding a new protocol to the existing IPv4 network infrastructure that confronts ISPs. The entire reason for this activity is the prospect of exhaustion of supply of IPv4 addresses. When this prospect was first aired, in 1990, it was assumed that the Internet would be supported by industry players that acted rationally in terms of common interests.

One of the more critical assumptions made in the development of transitional tools was that transition activity would be undertaken well in advance of IPv4 address exhaustion. Competitive interest would see each actor making the necessary investments in new technologies to mitigate the risks of attempting to operate a network in an environment of acute general scarcity of addresses. As much fun as the debate as to whom the “last” IPv4 address should be given might be, it was assumed that this event was, in fact, never going to happen. The assumption was that industry actors would anticipate this situation and take the necessary steps to avoid it. The transition to IPv6 would be effectively complete well before the stocks of IPv4 addresses had been exhausted, and IPv4 addresses would be an historical artefact well before we needed to use the last one!

Obviously, this scenario has not happened.

This industry is going to exhaust the available supplies of IPv4 addresses well before the transition to IPv6 is complete—and in some cases well before the transition process has even commenced! This situation creates an additional challenge for ISPs and the Internet, and raises a further question as well. The challenge is to fold into this dual-stack transition the additional factor of having to work with fewer and fewer IPv4 addresses as the transition process continues. This situation implies that the necessary steps that the ISP must take include ones that increase the intensity of use of each IPv4 address, and wherever possible substitute a private-use IPv4 address for public IPv4 addresses.

The question that this scenario raises is one of guessing how long this hybrid model of an Internet where a significant proportion of network services and network clients remains entrenched in an IPv4-only world will persist. For as long as such IPv4-only network domains persist, and for as long as these IPv4-only network domains encompass significant service and customer populations, all the other parts of the Internet are forced to maintain residual IPv4 capability and cannot transition their customers and services to an IPv6-only environment. Students of economic game theory may see some rich areas of study in this developing situation.

More practically, for an ISP the question becomes one of attempting to understand how long this hybrid period of attempting to operate a dual-stack network with continuing postexhaustion demand for further IPv4 addresses will last. Will an after-market for the redistribution of addresses emerge? How will the increasing scarcity pressure affect pricing in such a market? How long will demand persist for IPv4 addresses in the face of escalating prices? Will the industry turn to IPv6 in a rapid surge in response to cost escalation for additional IPv4 addresses, or will a dual-stack transition lumber on for many years? In such a large, diverse, heterogeneous environment of today’s Internet, the one constant factor is that the immediate future of the Internet is clouded with extremely high levels of uncertainty.

The cumulative effect of the individual decisions made by service providers, enterprises, carriers, vendors, policy makers, and consumers has created a somewhat chaotic environment that adds a significant level of uncertainty and associated investment risk into the current planning process for ISPs.

Carrier-Grade NATs

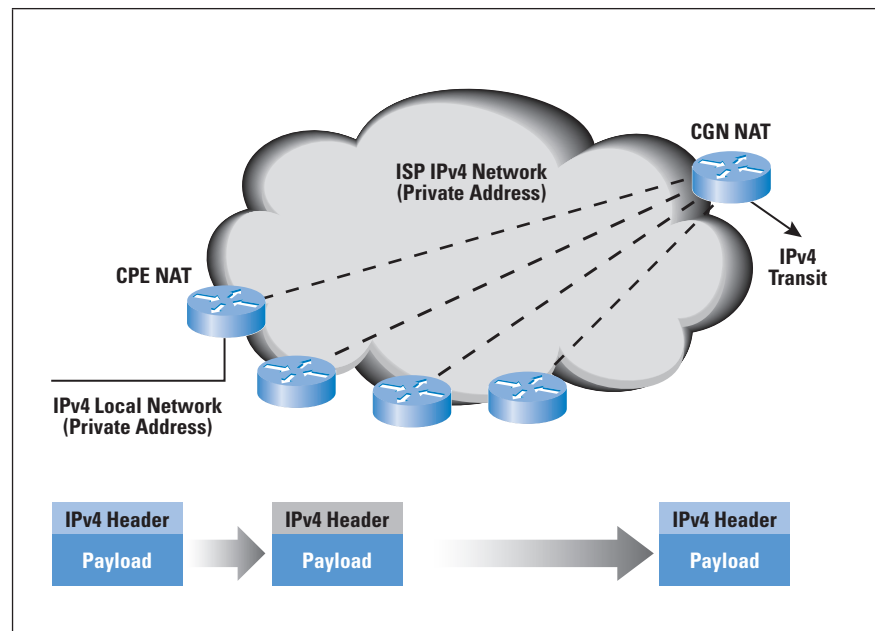
I have often heard it said that address scarcity in IPv4 is nothing new, and it first occurred when the first NAT device that supported port mapping was deployed. At this point the concept of *address sharing* was introduced to the Internet, and, from the perspective of the NAT industry, we have not looked back since.

In today's world NATs are extremely commonplace. Most clients are provisioned with a single address from their ISP, which they then share across their local network using a NAT. Whether it is well advised or not, NATs typically form part of a client's network security framework, and they often are an integral part of a customer's multihoming configuration if the client uses multiple providers.

But in this model of NATs as the CPE, the ISP uses one IPv4 address for each client. If the ISP wants to achieve greater levels of address compression, then it is necessary to share a single IPv4 address across multiple customers.

The most direct way to achieve this scenario is for ISPs to operate their own NAT, variously termed a *Carrier-Grade NAT* (CGN) or a *Large-Scale NAT* (LSN), or *NAT444*. This approach is the simplest, and, in essence, is a case of "more of the same" (Figure 6).

Figure 6: Carrier-Grade NATs



The Carrier-Grade NAT allows a single public address to be shared across multiple clients, who, in turn, further share this address across the end systems in their local networks.

From behind the CPE in the client edge network not much has changed with the addition of the CGN in terms of application behavior. It still requires an outbound packet to trigger a binding that would allow a return packet through to the internal destination, so nothing has changed there. Other aspects of NAT behavior, notably the NAT binding lifetime and the form of NAT “cone behavior” for UDP, take on the more restrictive of the two NAT functions in sequence. The binding times are potentially problematic in that the two NATs are not synchronized in terms of binding behavior. If the CGN has a shorter binding time, it is possible for the CGN to misdirect packets and cause application-level problems. However, this situation is not overly different from a single-level NAT environment where aggressively short NAT binding times also run the risk of causing application-level problems when the NAT drops the binding for an active session that has been quiet for an extended period of time.

However, one major assumption is broken in this structure, namely that an IP address is associated with a single customer. In the CGN model a single public IP address may be used simultaneously by many customers at once, albeit on different port numbers. This scenario has obvious implications in terms of some current practices in filters, firewalls, “black” and “white” lists, and some forms of application-level security and credentials where the application makes an inference about the identity and associated level of trust in the remote party based on the remote party’s IP address.

This approach is not without its potential operational problems as well. For the service provider, service resiliency becomes a critical concern in so far as moving traffic from one NAT-connected external service to another will cause all the current sessions to be dropped. Another concern is one of resource management in the face of potentially hostile applications. For example, an end host infected with a virus may generate a large amount of probe packets to a large range of addresses. In the case of a single edge NAT, the large volumes of bindings generated by this behavior become a local resource-management problem because the customer’s network is the only affected site. In the case where a CGN is deployed, the same behavior will consume port-binding space on the CGN and, potentially, can starve the CGN of external address port bindings. If this problem is seen to be significant, the CGN would need to have some form of external address rationing per internal client in order to ensure that the entire external address pool is not consumed by a single errant customer application.

The other concern here is one of *scalability*. Whereas the most effective use of the CGN in terms of efficiency of usage of external addresses occurs when the greatest numbers of internal edge NATed clients are connected, there are some real limitations in terms of NAT performance and address availability when a service provider wants to apply this approach to networks where the customer population is in the millions or larger. In this case the service provider must use an IPv4 private address pool to number every client. But if network 10 is already used by each customer as its “internal” network, then what address pool can be used for the service provider’s private address space? One of the few answers that come to mind is to deliberately partition the network into numerous discrete networks, each of which can be privately numbered from **172.16.0.0/12**, allowing for some 600,000 or so customers per network partition, and then use a transit network to “glue” together the partitioned elements.

The advantage of the CGN approach is that nothing changes for the customer. There is no need for any customers to upgrade their NAT equipment or change it in any way, and for many service providers this motivation is probably sufficient to choose this path. The disadvantages of this approach lie in the scaling properties when looking at very large deployments, and the concerns of application-level translation, where the NAT attempts to be “helpful” by performing *Deep Packet Inspection* and rewriting what it thinks are IP addresses found in packet payloads. Having one NAT do this process is bad enough, but loading them up in sequence is a recipe for trouble.

Are there alternatives?

The Address-plus-Port Approach

One NAT in the path is certainly worse than none from the perspective of application agility and functions. And two NAT functions do not make it any better! Inevitably, that second NAT device adds some additional levels of complexity and fragility into the process.

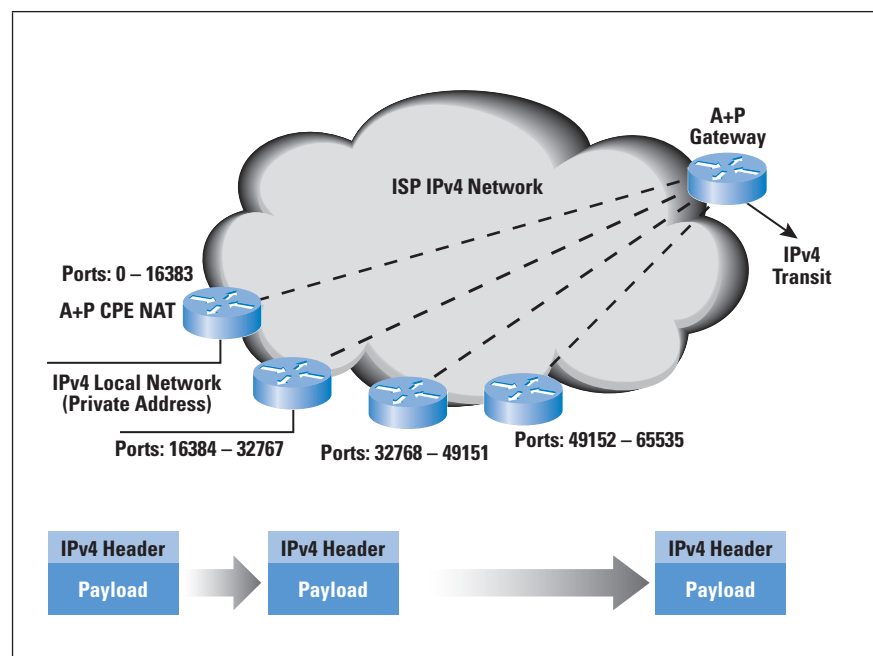
The question is, can these two NAT functions be collapsed back into a single NAT, yet still allow sharing of public IPv4 addresses across multiple end clients? CPE NAT devices currently map connections into the 16-bit *port* field of the single external address. If the CPE NAT could be coerced into performing this mapping into, say, 15 bits of the port field, then the external address could be shared between two edge CPEs, with the leading bit of the port field denoting which CPE. Obviously, moving the bit marker further across the port field will allow more CPE devices to share the one address, but it will reduce the number of available ports for each CPE in the process.

The theory is again quite simple. The CPE NAT is dynamically configured with an external address, as happens today, and a port range, which is the additional constraint. The CPE NAT performs the same function as before, but it is now limited in terms of the range of external port values it can use in its NAT bindings to those that lie within the provided port range. Other CPE devices are concurrently using the same external IP address, but with a different port range.

For outgoing packets this scenario implies only a minor change to the network architecture, in that the RADIUS exchange to configure the CPE now must also provide a port range to the CPE device. The CPE is then constrained such that as it maps private addresses and TCP or UDP port values to the external address and port values, the mapped port value must fall within the configured range.

The handling of incoming packets is more challenging. Here the service provider must forward the packet based not only on the destination IP address, but also on the port value in the TCP or UDP header, because there are now multiple CPE egress points that share the same IP address. A convenient way to perform forwarding is to take the Dual-Stack Lite approach and use an IPv4-in-IPv6 tunnel between the CPE and the external address-plus-port (A+P) gateway. This address-plus-port gateway needs to be able to associate each address and port range with the IPv6 address of a CPE (which it can learn dynamically as it decapsulates outgoing packets that are similarly tunneled from the CPE to the address-plus-port gateway). Incoming packets are encapsulated in IPv6 using the IPv6 destination address that it has learned previously. In this manner the NAT function is performed just once, at the edge, much as it is today, and the interior device is a more conventional form of tunnel server (Figure 7).

Figure 7: Address-plus-Port-Approach



This approach relies on every CPE device being able to operate using a restricted port range, to perform IPv4-in-IPv6 tunnel ingress and egress functions, and act as an IPv6 provisioned endpoint for the service provider network. This set of constraints is perhaps unrealistic for many service provider networks. Further modifications to this model propose the use of an accompanying CGN operated by the service provider to handle those CPE devices that cannot support this address-plus-port function.

This approach has some positive aspects. Pushing the NAT function back to the network edge has some considerable advantage over the approach of moving the NAT to the interior of the network. The packet rates are lower at the edge, allowing for commodity computing to process the NAT functions across the offered packet load without undue stress. The ability to control the NAT behavior with the *Internet Gateway Device* protocol as part of the *Universal Plug and Play* (uPnP) framework will still function in an environment of restricted port ranges. Aside from the initial provisioning process to equip the CPE NAT with a port range, the CPE and the edge environment are largely the same as that of today's CPE NAT model.

That is not to say that this approach is without its negative aspects, and it is unclear as to whether the perceived benefits of a "local" NAT function outweigh the problems in this particular model of address sharing. The concept of port "rationing" is a very suboptimal means of address sharing, given that when a CPE is assigned a port range, those port addresses are unusable by any other CPE. The prudent service provider would assign to each CPE a port address pool equal to some estimate of peak demand, so that, for example, each CPE would be assigned some 1024 ports, allowing a single external IP address to be shared across only some 60 such CPE clients. The Carrier-Grade NAT and Dual-Stack Lite approaches do not attempt this form of rationed allocation, allowing the port address pool to be treated as a common resource, with far higher levels of usage efficiency. The leverage obtained in terms of efficiently using these additional 16 bits of address space is reduced by the imposition of a fixed boundary between customer and service provider use. The central NAT model effectively pools the port address range and would result in more efficient sharing of this common pool across a larger client base.

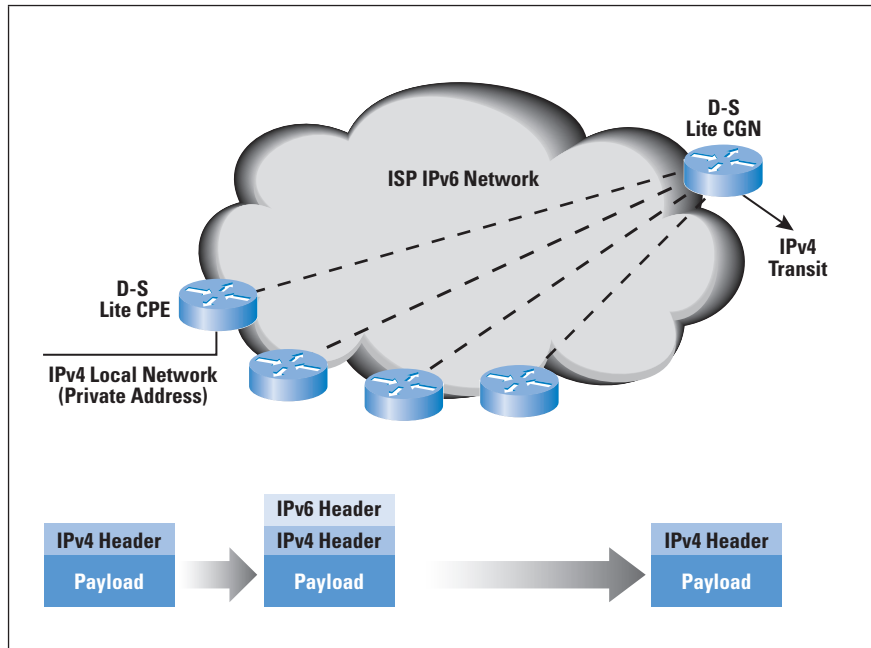
The other consideration here is that this approach means a higher overhead for the service provider, in that the service provider would have to support both "conventional" CPE equipment and address-plus-port equipment. In other words, the service provider will have to deploy a CGN and support customer CPE using a two-level NAT environment in addition to operating the address-plus-port infrastructure. Unless customers would be willing to pay a significant price premium for such address-plus-port service, it is unlikely that this option would be attractive for the service provider as an additional cost above the CGN cost.

Dual-Stack Lite

The concept behind the *Dual-Stack Lite* approach is that the service provider's network infrastructure will need to support IPv6 running in native mode in any case, so is there a way in which the service provider can continue to support IPv4 customers without running IPv4 internally?

Here the customer NAT is effectively replaced by a tunnel ingress-egress function in the Dual-Stack Lite home gateway. Outgoing IPv4 packets are not translated, but are encapsulated in an IPv6 packet header, which contains a source address of the carrier side of the home gateway unit, and a destination address of the ISP's gateway unit. From the service provider's perspective, each customer is no longer uniquely addressed with an IPv4 address, but instead is addressed with a unique IPv6 address, and provided with the IPv6 address of the provider's combined IPv6 tunnel egress point and IPv4 NAT unit (Figure 8).

Figure 8: Dual-Stack Lite



The service provider's Dual-Stack Lite gateway unit will perform the IPv6 tunnel termination and a NAT translation using an extended local binding table. The NAT "interior" address is now a 4-tuple of the IPv4 source address, protocol ID, and port, plus the IPv6 address of the home gateway unit, while the external address remains the triplet of the public IPv4 address, protocol ID, and port. In this way the NAT binding table contains a mapping between interior "addresses" that consist of IPv4 address and port plus a tunnel identifier, and public IPv4 exterior addresses. This way the NAT can handle a multitude of net 10 addresses, because they can be distinguished by different tunnel identifiers.

The resultant output packet following the stripping of the IPv6 encapsulation and the application of the NAT function is an IPv4 packet with public source and destination addresses. Incoming IPv4 packets are similarly transformed, where the IPv4 packet header is used to perform a lookup in the Dual-Stack Lite gateway unit, and the resultant 4-tuple is used to create the NAT-translated IPv4 packet header plus the destination address of the IPv6 encapsulation header.

The advantage of this approach is that there now needs to be only a single NAT in the end-to-end path, because the functions of the customer NAT are now subsumed by the carrier NAT. This scenario has some advantages in terms of those messy “value-added” NAT functions that attempt to perform deep packet inspection and rewrite IP addresses found in data payloads. There is also no need to provide each customer with a unique IPv4 address, public or private, so the scaling limitations of the dual-NAT approach are also eliminated. The disadvantages of this approach lie in the need to use a different CPE device—or at least one that is reprogrammed. The device now requires an external IPv6 interface and at the minimum an IPv4/IPv6 tunnel gateway function. The device can also include a NAT if so desired, but it is not required in terms of the basic Dual-Stack Lite architecture.

This approach pushes the translation into the interior of the network, where the greatest benefit can be derived from port multiplexing, but it also creates a critical hotspot for the service itself. If the Dual-Stack Lite NAT fails in any way, the entire customer base is disrupted. It seems somewhat counterintuitive to create a resilient end-to-end network with stateless switching environments and then place a critical stateful unit right in the middle!

Protocol Translation

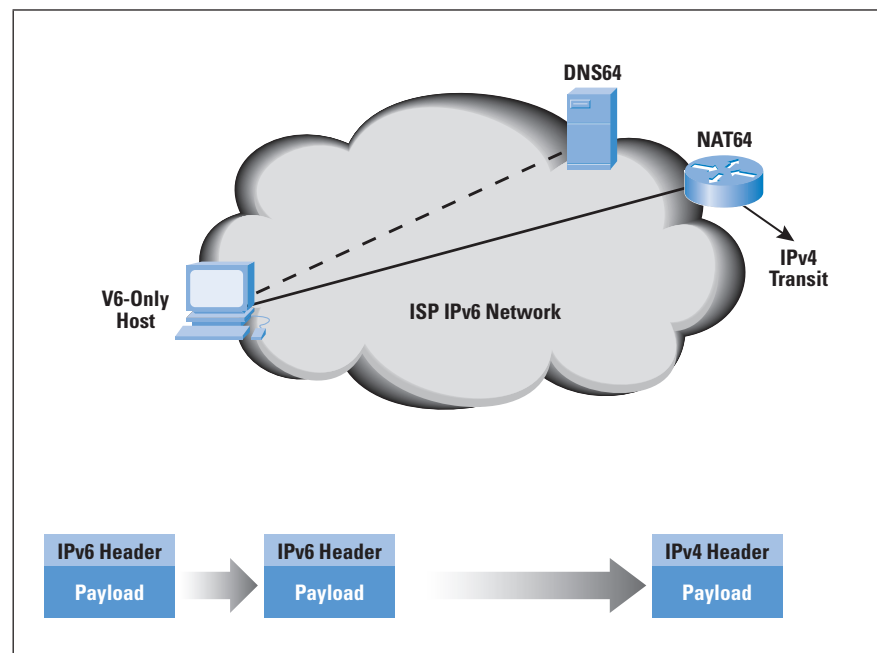
So far we have looked at two general forms of approach to hybrid networks that are intended to support both IPv6 transition and greater levels of address usage in IPv4, namely address mapping and tunneling. A third approach lies in the area of protocol translation.

RFC 2765^[10] contains the details of a relatively simple protocol-translation mechanism. The approach relies on the basic observation that IPv6 did not make any radical changes to the basic IP architecture of IPv4, and that it was therefore possible to define a stateless mapping algorithm that could translate between certain IPv4 and IPv6 packets. Of course the one major problem here is that there are far more addresses in IPv6 than in IPv4, so the approach used was to map IPv4 addresses into the trailing 32 bits of the IPv6 address prefix `::FFFF:0:0/96`. The approach assumed that to the IPv6-only end host the entire IPv4 network was visible in this mapped IPv6 prefix, and that when the IPv6-only end host wished to communicate with a remote host who was addressed using this IPv4-mapped prefix it would use a source address also drawn from the same IPv4-mapped prefix. In other words, it assumed that all IPv6-only hosts were also assigned a unique IPv4 address.

The *NAT-Protocol Translation* (NAT-PT) approach attempted to relax this constraint, allowing IPv6-only hosts to use a dynamic mapping to a public IPv4 address through the NAT-PT function, in the same way as NAT functions work in an all-IPv4 domain (Figure 9). The proposed approach assumed that the local host was located behind a modified DNS environment where the IPv4 “A” record of an IPv4-only remote service is translated by the DNS gateway into a local IPv6 address where the initial 96 bits of the IPv6 address identify the internal address of the NAT-PT gateway and the trailing 32 bits are the IPv4 address of the remote service. When the local host then uses this address as an IPv6 destination address, the packet is directed by the local routing environment to the NAT-PT device. This device can construct an “equivalent” IPv4 packet by using the local IPv4 address as the source address and the last 32 bits of the IPv6 address as the destination address, and bind the IPv6 source port to a free local port value. These sets of transforms can be locally stored as an active NAT binding. Return IPv4 packets can be mapped back into their “equivalent” IPv6 form by using the values in the binding to perform a reverse set of transforms on the IP address and port fields of the packet.

This approach was published as RFC 2766^[11] in February 2000. Some 7 years later in July 2007, the IETF published RFC 4966^[12], deprecating NAT-PT to “historic,” with an associated list of applications that would not operate correctly through such a device. This negative judgement of NAT-PT seems rather curious to me, given that conventional CPE NAT functions in IPv4 appear to share most, if not all, of the same shortfalls that are listed in RFC 4966. Given the extensive set of compromises that are required in the environment that is partially crippled by IPv4 address exhaustion, it seems rather contradictory to insist upon extremely high levels of functions and robustness from these hybrid translation approaches.

Figure 9: NAT Protocol Translation – NAT64



Not unsurprisingly, NAT-PT is undergoing a revival, this time under the name “NAT64.” Not much has changed from the basic approach outlined in NAT-PT. The IPv6-only client performs a DNS lookup through a modified DNS server that is configured with DNS64. If the queried name contains only an IPv4 address, the DNS64 server synthesises an IPv6 response by merging the prefix address of the NAT64 gateway with the IPv4 address. When the client uses this address, the IPv6 packet is directed to the NAT64 gateway, and the same transform as described previously for NAT-PT takes place.

This setup is similar to the CGN model, in so far as the service provider operates a common NAT that shares an IPv4 address pool across a set of end clients.

ISP Conclusions

There really is no single clear path forward from this point. Different ISPs will see some advantages in pursuing different approaches to this dual problem of introducing IPv6 into their service portfolio and at the same time introducing additional measures that allow more efficient use of IPv4 addresses.

However, one common theme is becoming clear. So far ISPs have been able to “externalize” many of these problems by pushing much of the complexity and fragility of NAT functions out to the customer and loading up the CPE with these functions. This approach of externalizing much of the complexity of address compression in NAT functions over to the customer’s network cannot be sustained with the IPv6 transition, and no matter which approach is used, whether it is a CGN, NAT64, Dual-Stack Lite, 6RD, or MPLS with 6PE, the ISP now has to actively participate in the delivery of IPv6 and in increasing the efficiency of the use of IPv4.

So for the ISP it is time to start making some technical choices as to how to address the combination of these two rather unique challenges of transition and exhaustion.

References

- [1] Daniel Karrenberg, Gerard Ross, Paul Wilson, and Leslie Nobile, “Development of the Regional Internet Registry System,” *The Internet Protocol Journal*, Volume 4, No. 4, December 2001.
- [2] Geoff Huston, “Anatomy: A Look inside Network Address Translators,” *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [3] Alex Conta and Stephen Deering, “Generic Packet Tunneling in IPv6 Specification,” RFC 2473, December 1998.
- [4] Yakov Rekhter, Bob Moskowitz, Daniel Karrenberg, Geert Jan de Groot, and Eliot Lear, “Address Allocation for Private Internets,” RFC 1918, February 1996.
- [5] Christian Huitema, “Teredo: Tunneling IPv6 over UDP through Network Address Translations (NATs),” RFC 4380, February 2006.

- [6] Jonathan Rosenberg, Rohan Mahy, Philip Matthews, and Dan Wing, “Session Traversal Utilities for NAT (STUN),” RFC 5389, October 2008.
- [7] Alex Conta, Stephen Deering, and Mukesh Gupta, “Internet Control Message Protocol (ICMPv6) for the Internet Protocol Version 6 (IPv6) Specification,” RFC 4443, March 2006.
- [8] Mark Townsley and Ole Troan, “IPv6 Rapid Deployment on IPv4 Infrastructures (6rd) – Protocol Specification,” RFC 5969, August 2010.
- [9] Jeremy De Clercq, Dirk Ooms, Marco Carugi, and Francois Le Faucheur, “BGP-MPLS IP Virtual Private Network (VPN) Extension for IPv6 VPN,” RFC 4659, September 2006.
- [10] Erik Nordmark, “Stateless IP/ICMP Translation Algorithm (SIIT),” RFC 2765, February 2000.
- [11] George Tsirtsis and Pyda Srisuresh, “Network Address Translation – Protocol Translation (NAT-PT),” RFC 2766, February 2000.
- [12] Cedric Aoun and Elwyn Davies, “Reasons to Move the Network Address Translator - Protocol Translator (NAT-PT) to Historic Status,” RFC 4966, July 2007.

Further Reading

The IETF has been working on the issues related to the transition to IPv6 for the past 18 years, and in the intervening period has generated many hundreds of documents. In selecting the following documents as a helpful reading list, I have tried to select only from the more recent documents and those that are overviews of transition technologies rather than reference specifications for individual technologies.

- [1] Jari Arkko and Fred Baker, “Guidelines for Using IPv6 Transition Mechanisms during IPv6 Deployment,” Internet Draft, Work in Progress, December 2010.

*The document discusses the IPv6 deployment models and migration tools, and considers what appears to be effective in networks to date. This Internet Draft, **draft-arkko-ipv6-transition-guidelines-14.txt**, is about to be published as an Informational RFC.*

- [2] Brian Carpenter and Sheng Jian, “Emerging Service Provider Scenarios for IPv6 Deployment,” RFC 6036, October 2010.

This document describes practices and plans that are emerging among Internet Service Providers for the deployment of IPv6 services, using data collected in a survey of numerous ISPs carried out in early 2010.

- [3] Reinaldo Penno, Tarun Saxena, Mohamed Boucadair, and Senthil Sivakumar, “Analysis of 64 Translation,” Internet Draft, Work in Progress, **draft-ietf-behave-64-analysis-01**, January 2011.

This paper is a working document of the IETF’s BEHAVE Working Group. The document notes that because of specific problems, NAT-PT was deprecated by the IETF as a mechanism to perform IPv6-IPv4 translation. Since then, new efforts have been undertaken within IETF to standardize alternative mechanisms to perform IPv6-IPv4 translation. This document evaluates how the new translation mechanisms avoid the problems that caused the IETF to deprecate NAT-PT.

- [4] Fred Baker, Xing Li, and Kevin Yin, “Framework for IPv4/IPv6 Translation,” Internet Draft, Work in Progress, August 2010.

*It is common in the IETF these days to generate a “framework” document as part of the process of developing technical specifications. This draft is a framework document for the general IPv4/IPv6 translation technology. This Internet Draft, **draft-ietf-behave-v6v4-framework-10.txt**, will soon be published as an Informational RFC.*

- [5] Elwyn Davies, Suresh Krishnan, and Pekka Savola, “IPv6 Transition/Coexistence Security Considerations,” RFC 4942, September 2007.

The transition into a dual-stack environment, while attempting to preserve the integrity of a single service regime, presents numerous security concerns. This document is a good overview of such concerns.

- [6] Dan Wing and Andrew Yourtchenko, “Improving User Experience with IPv6 and SCTP,” *The Internet Protocol Journal*, Volume 13, No. 3, September 2010.

Building efficient applications in a dual-stack world can be very challenging. It is often the case that poor management of a dual-stack system can make the user experience far slower than just continuing in the IPv4 world. One way to redress this problem is to exchange sequential testing of IPv6 and IPv4 connectivity into a parallel operation—both protocols at once. This article explains the concept.

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005; he served on the Board of Trustees of the Internet Society from 1992 until 2001. E-mail: **gih@apnic.net**

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ contains standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2011 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

June 2011

Volume 14, Number 2

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
Securing BGP	2
IPv6 Site Multihoming.....	14
Reflecting on World IPv6 Day	23
Letters to the Editor	25
Call for Papers	29
Fragments	30

The process of adding security to various components of Internet architecture reminds me a little bit of the extensive seismic retrofitting that has been going on in California for decades. The process is slow, expensive, and occasionally intensified by a strong earthquake after which new lessons are learned. Over the past 13 years this journal has carried many articles about network security enhancements: *IP Security* (IPSec), *Secure Sockets Layer* (SSL), *Domain Name System Security Extensions* (DNSSEC), *Wireless Network Security*, and *E-mail Security*, to name but a few. In this issue we look at routing security again, specifically the efforts underway in the *Secure Inter-Domain Routing* (SIDR) Working Group of the IETF to provide a secure mechanism for route propagation in the *Border Gateway Protocol* (BGP). The article is by Geoff Huston and Randy Bush.

Our second article discusses *Site Multihoming* in IPv6. Multihoming is a fairly common technique in the IPv4 world, but as part of the development and deployment of IPv6, several new and improved solutions have been proposed. Fred Baker gives an overview of these solutions and discusses the implications of each proposal.

By all accounts, *World IPv6 Day* was a successful demonstration and an important step toward deployment of IPv6 in the global Internet. Several major sites left IPv6 connectivity in place after the event, an encouraging sign. Discussions are already underway for another similar event, this time perhaps lasting for as long as a week. Phil Roberts gives an overview of what happened on June 8 and provides pointers to some of the important lessons learned from this experiment.

I want to take a moment to mention the IPJ subscription renewal campaign. As you know, each subscriber is issued a unique subscription ID that, coupled with an e-mail address, gives access to the subscription database by means of a “magic URL.” Unfortunately, sometimes the e-mail containing this URL may not arrive in the subscriber’s mailbox, perhaps because of spam filtering. Additionally, readers change e-mail addresses as well as postal addresses. If your subscription has expired or you have changed e-mail, postal mail, or delivery preference, send an e-mail to ipj@cisco.com with the updated information and we will make sure your subscription is re-instated. The purpose of the renewal campaign is to ensure that we are sending copies of IPJ to the correct addresses and only to those who prefer paper copies. IPJ is always available via our website at <http://cisco.com/ipj>

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

Securing BGP with BGPsec

by Geoff Huston, APNIC and Randy Bush, IIR

For many years the fundamental elements of the Internet: *names* and *addresses*, were the source of basic structural vulnerabilities in the network. With the increasing momentum behind the deployment of *Domain Name System Security Extensions* (DNSSEC)^[0], there is some cause for optimism that we have the elements of securing the name space now in hand, but what about addresses and routing? In this article we will look at current efforts within the *Internet Engineering Task Force* (IETF) to secure the use of addresses within the routing infrastructure of the Internet, and the status of current work of the *Secure Inter-Domain Routing* (SIDR) Working Group.

We will look at the approach the SIDR Working Group has taken, and examine the architecture and mechanisms that it has adopted as part of this study. This work was undertaken in three stages: the first concentrated on the mechanisms to support attestations relating to addresses and their use; the second looked at how to secure origination of routing announcements; and the third looked at how to secure the transitive part of *Border Gateway Protocol* (BGP) route propagation.

Supporting Attestations About Addresses Through the RPKI

Prior work in the area of securing the Internet routing system has focused on the operation of BGP in an effort to secure the operation of the protocol and validate, as far as is possible, the contents of *BGP Update* messages. Some notable contributions in more than a decade of study include *Secure-BGP* (S-BGP)^[1, 16], *Secure Origin BGP* (soBGP)^[2], *Pretty Secure BGP* (psBGP)^[3], IRR^[4], and the use of an *Autonomous System* (AS) *Resource Record* (RR) in the *Domain Name System* (DNS), signed by DNSSEC^[5].

The common factor in this prior work was that they all required, as a primary input, a means of validating basic assertions relating to origination of a route into the interdomain routing system: that the IP address block and the AS numbers being used are valid and that the parties using these IP addresses and AS numbers in the context of routing advertisement are properly authorized to so do.

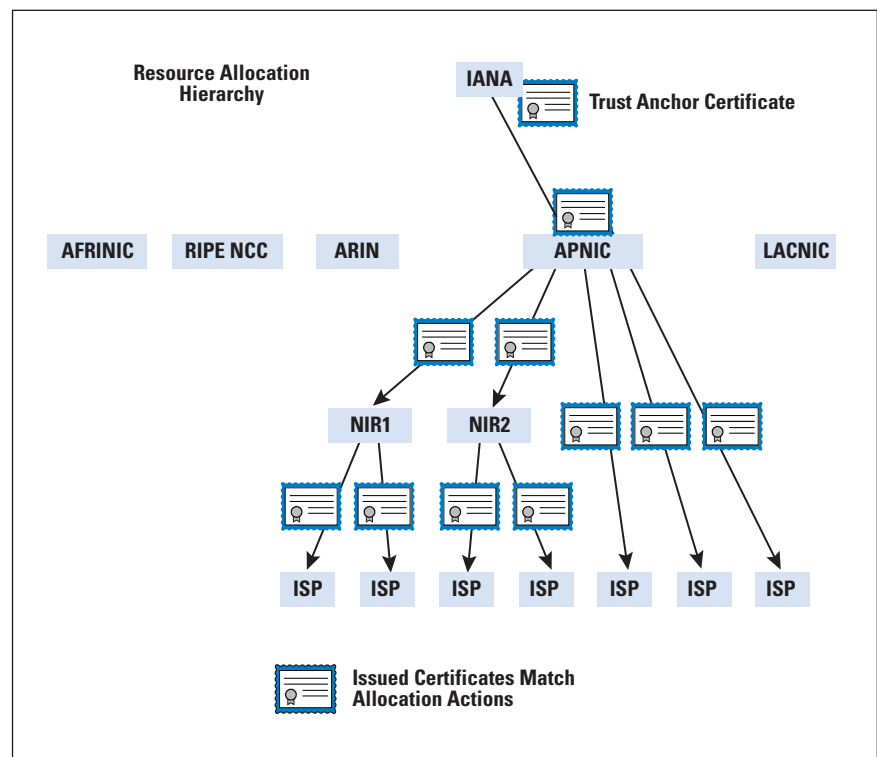
The approach adopted by SIDR for the way in which trust is formalized in the routing environment is through the use of *Resource Certificates*. These certificates are X.509 certificates that conform to the *Public-Key Infrastructure X.509* (PKIX) profile^[6]. They also contain an extension field that lists a collection of IP resources (IPv4 addresses, IPv6 addresses, and AS Numbers)^[7]. These certificates attest that the certificate issuer has granted to the certificate subject a unique “right-of-use” for the associated set of IP resources, by virtue of a resource allocation action.

This concept mirrors the resource allocation framework of the Internet Assigned Numbers Authority (IANA), the *Regional Internet Registries* (RIRs), operators, and others, and the certificate provides a means for a third party (relying party) to formally validate assertions related to resource allocations^[8].

The hierarchy of the *Resource Public Key Infrastructure* (RPKI) is based on the administrative resource allocation hierarchy, where resources are distributed from the IANA to the RIRs, *Local Internet Registries* (LIRs), *National Internet Registries* (NIRs), and end users. The RPKI mirrors this allocation hierarchy with certificates that match current resource allocations (Figure 1).

The *Certification Authorities* (CAs) in this RPKI correspond to entities that have been allocated resources. Those entities are able to sign authorities and attestations, and to do so they use specific-purpose *End Entity* (EE) certificates. This additional level of indirection allows the entity to customize each issued authority for specific subsets of number resources that are administered by this entity. Through the use of single-use EE certificates, the issuer can control the validity of the signed authority through the ability to revoke the EE certificate used to sign the authority. As is often the case, a level of indirection comes in handy.

Figure 1: Hierarchy of the RPKI



Signed attestations relating to addresses and their use in routing are generated by selecting a subset of resources that will be the subject of the attestation, by generating an EE certificate that lists these resources, and by specifying validity dates in the EE certificate that correspond to the validity dates of the authority. The authority is published in the RPKI repository publication point of the entity. The RPKI makes conventional use of *Certificate Revocation Lists* (CRLs) to revoke certificates that have not expired but are no longer valid. Every Certification Authority in the RPKI regularly issues a CRL according to the declared CRL update cycle of the Certification Authority. A Certification Authority certificate may be revoked by an issuing authority for numerous reasons, including key rollover, the reduction in the resource set associated with the certificate subject, or termination of the resource allocation. To invalidate an object that can be verified by a given EE certificate, the Certification Authority that issued the EE certificate can revoke the corresponding EE certificate.

The RPKI uses a distributed publication framework, wherein each Certification Authority publishes its products (including EE certificates, CRLs, and signed objects) at a location of its choosing. The set of all such repositories forms a complete information space, and it is fundamental to the model of securing BGP in the public Internet that the entire RPKI information space be available to every *Relying Party* (RP). It is the role of each RP to maintain a local cache of the entire distributed repository collection by regularly synchronizing each element in the local cache against the original repository publication point. To assist RPs in the synchronization task, each RPKI publication point uses a *manifest*, a signed object that lists the names (and hash values) of all the objects published at that publication point. It is used to assist RPs to ensure that they have managed to synchronize against a complete copy of the material published at the Certification Authority publication point.

The utility of the RPKI lies in its ability to validate digitally signed information and, therefore, give relying parties some confidence in the validity of signed attestations about addresses and their use. The particular utility of the RPKI is not as a means of validation of attestations of an individual's identity or that individual's role, but as a means of validating that person's authority to use IP address resources. Although it is possible to digitally sign any digital object, it has been suggested that the RPKI system uses a very small number of standard signed objects that have particular meaning in the context of routing security.

Securing Route Origination

The approach adopted by SIDR to secure origination of routing information is one that uses a particular signed authority, a *Route Origination Authorization* (ROA)^[10]. An ROA is an authority created by a prefix holder that authorizes an AS to originate one or more specific route advertisements into the interdomain routing system.

An ROA is a digital object formatted according to the *Cryptographic Message Syntax Specification* (CMS)^[11] that contains a list of address prefixes and one AS number. The AS is the specific AS being authorized to originate route advertisements for one or more of the address prefixes in the ROA. The CMS object also includes the EE resource certificate for the key used to verify the ROA. The IP Address extension in this EE certificate must encompass the IP address prefixes listed in the ROA contents.

The ROA conveys a simple authority. It does not convey any further routing policy information, nor does it convey whether or not the AS holder has even consented to actually announce the prefix(es) into the routing system. The associated EE certificate is used to control the validity of the ROA, and the CMS wrapper is used to securely bind the ROA and the EE certificate within a single signed structure.

There is one special ROA, one that authorizes AS 0 to originate a route. Because AS 0 is a reserved AS that should never be used by a BGP speaker, this ROA is a “negative” authority, used to indicate that no AS has authority to originate a route for the address prefix(es) listed in the ROA.

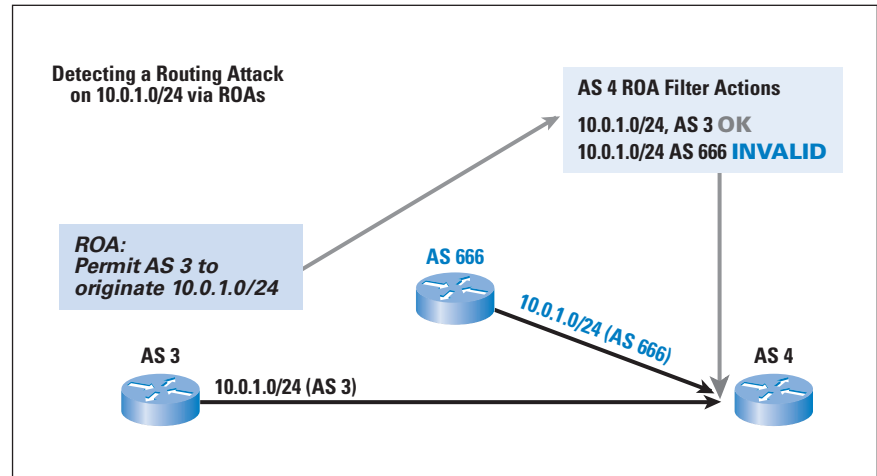
If the entire routing system were to be populated with ROAs, then identification of an invalid route advertisement would be directly related to detection of an invalid ROA or a missing ROA. However, in a more likely scenario of partial use of ROAs (such as when only some legitimate route originations are authorized in a ROA), the absence of an ROA cannot be interpreted simply as an unauthorized use of an address prefix. This scenario leads to the use of a tri-state validation process for routes, as follows.

If a given route matches exactly the information contained in an ROA whose EE certificate can be validated in the RPKI (a “valid” ROA), then the route can be regarded as a “valid” origination. Where the address prefix matches that in a valid ROA but the origination AS does not match the AS number in the ROA, and there are no other valid ROAs that explicitly validate the announcing AS, then the route can be considered to be “invalid.” Also, where the address prefix is more specific than that of a valid ROA, and there are no other valid ROAs that match the prefix, then the route can also be considered “invalid.” Where the prefix in a route is not described in any ROA and is not a more specific prefix of any ROA, the route has an “unknown” validation outcome.

These three potential outcomes can be considered a set of relative local preferences. Routes whose origins can be considered “valid” are generally proposed to be preferred over routes whose origins are unknown, which, in turn, can generally be preferred over routes whose origins are considered invalid. However, such relative preferences are a matter to be determined by local routing policy. Local policies may choose to adopt a stricter policy and, for example, discard routes with an invalid validation outcome^[12].

The way in which ROAs are used to validate the origin of routes in BGP differs from many previous proposals for securing BGP. In this framework the ROAs are published in the RPKI distributed repository framework. Each RP can use the locally cached collection of valid ROAs to create a validation filter collection, with each element of the set containing an address, prefix size constraints, and an originating AS. It is this filter set—rather than the ROAs themselves—that are fed to the local routers^[13]. An example of the way in which ROAs can be used to detect prefix hijack attempts is shown in Figure 2.

Figure 2: Use of ROAs to detect Unauthorized Route Origination



The model of injecting validation of origination into the BGP domain is an example of a highly modular and piecemeal deployment. There are no changes to the BGP protocol for this origin validation part of the secure routing framework.

The process of securing origination starts with the address holder, who generates local keys and requests certification of their address space from the entity from whom their addresses were allocated or assigned. With this Certification Authority resource certificate, the address holder is then in a position to generate an EE certificate and a ROA that assigns an authority for a nominated AS to advertise a route for an address prefix drawn from its address holdings. The one condition here is that if an address holder issues a ROA for an address prefix providing an authority for one AS to originate a route for this prefix, then the address holder is required to issue ROAs for all the ASes that have been similarly authorized to originate a route for this address prefix. The address holder publishes this ROA in its publication point in the distributed RPKI repository structure.

Relying parties can configure a locally managed cache of the distributed RPKI repository and collect the set of valid ROAs. They can then, with the dedicated RPKI cache-to-router protocol^[13], maintain, on a set of “client” routers, the set of address prefix/originating AS authorities that are described in valid ROAs. The BGP-speaking router can use this information as an input to the local route decision process.

This model of operation supports piecemeal incremental deployment, wherein individual address holders may issue ROAs to authorized routing advertisements independent of the actions of other address holders. Also, ASs may deploy local validation of route origination independently of the actions of other ASs. And given that there are no changes to the operation of BGP, then there are no complex inter-dependencies that hinder piecemeal incremental deployment of this particular aspect of securing routing.

Securing Route Propagation: BGPsec

Origin validation as described earlier does not provide cryptographic assurance that the origin AS in a received BGP route was indeed the originating AS of this route. A malicious BGP speaker can synthesize a route as if it came from the authorized AS. Thus, it is very useful in detecting accidental misannouncements, but origination validation does little to prevent malicious routing attacks from a determined attacker.

In looking at the operation of the BGP protocol, some parts of the protocol interaction are strictly local between two BGP-speaking peers, such as advising a peer of local attributes. Another part of the BGP protocol is a “chained” interaction, in which each AS adds information to the protocol object. This attribute of a BGP update, the *AS Path*, is not only useful to detect and prevent routing loops, it is also used in the BGP best-path-selection algorithm.

A related routing security question concerns the validity of this “chained” information, namely the AS Path information contained in a route. Within the operation of the BGP protocol, each AS that propagates an update to its AS neighbors is required to add its AS number to the AS Path sequence. The inference is that at any stage in the propagation of a route through the interdomain routing system, the AS Path represents a viable AS transit sequence from the local AS to the AS originating the route. This AS Path attribute of a route is used for loop detection. Locally, the AS Path may also be used as input to a local route policy process, using the length of the AS Path as a route metric.

Attacks on the AS Path can be used to subvert the routing environment. A malicious BGP speaker may manipulate the AS Path to prevent an AS from accepting a route by adding its AS number to the AS Path, or it may attempt to make a particular route more likely to be selected by a remote AS by stripping out ASs from the AS Path. Accordingly, it is important to equip a secure BGP framework with the ability to validate the authenticity of the AS Path presented in a BGP update^[14].

When attempting to validate an AS path, many potential validation questions must be addressed.

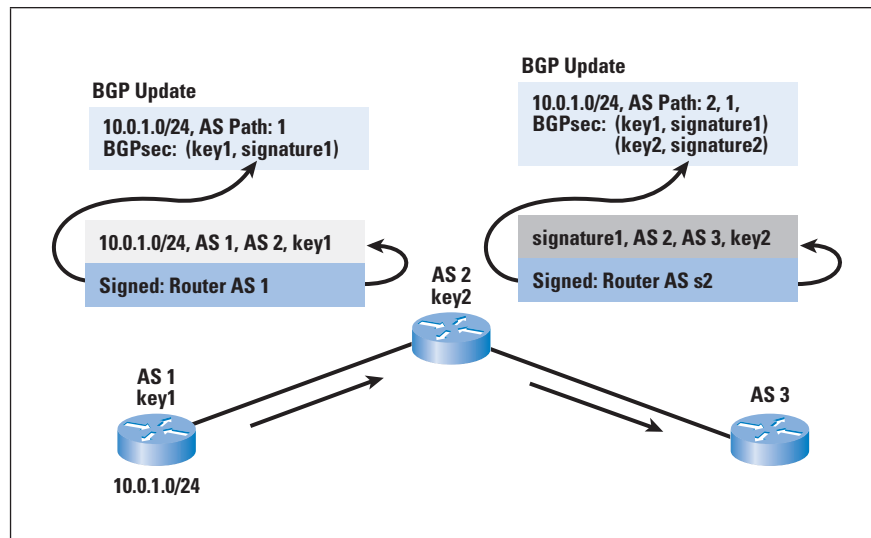
- The first and weakest question is: Are all ASs in the AS Path valid ASs?
- A slightly stronger validation question is: Do all the AS pairs in the AS Path represent valid AS adjacencies (where both ASs in the pair-wise association are willing to attest to their mutual adjacency in BGP)?
- A even stronger question is: Does the sequence of ASs in the AS Path represent the actual propagation path of the BGP route object?

This last question forms the basis for the SIDR activity in defining an AS Path validation framework, BGPsec. This attempt is to assure a BGP speaker that the operation of the BGP protocol is operating correctly and that the content of a BGP update correctly represents the inter-AS propagation path of the update from the point of origination to the receiver of the route. This tool is not the same as a policy validation tool and it does not necessarily assure the receiver of the route that this update conforms to the routing policies of neighboring BGP speakers. This route also does not necessarily reflect the policy intent of the originator of the route. The BGPsec framework proposed for securing the AS Path also uses a local RPKI cache, but it includes an additional element of certification. The additional element of the security credentials used here is an extension to the certification of AS numbers with a set of operational keys and their associated certificates used for signing update messages on *External Border Gateway Protocol* (eBGP) routers in the AS. These “router certificates” can sign BGP update attributes in the routing infrastructure, and the signature can be interpreted as being a signature made “in the name of” an AS number.

In the BGPsec framework, eBGP-speaking routers within the AS have the ability to “sign” a BGP update before sending it. In this case, the added signature “covers” the signature of the received BGP update, the local AS number, the AS number to which the update is being sent, as well as a hash of the public key part of the router key pair used to sign route updates.

The couplet of the public key hash and the signature itself are added to the BGP protocol update as BGPsec update attributes. As the update traverses a sequence of transit ASs, each eBGP speaker at the egress of each AS adds its own public key hash and digital signature to the BGPsec attribute sequence (Figure 3).

Figure 3: BGPsec AS Path Protection



This interlocking of signatures allows a receiver of a BGP update to use the interlocking chain of digital signatures to validate (for each AS in the AS Path) that the corresponding signature was correctly generated “in the name of” that AS in the AS Path, and that the next AS in the path matches the next AS in the signed material. The “forward signing” that includes the AS to which the update is being sent prevents a man-in-the-middle attack of the form of taking a legitimate outbound route announcement destined for one neighbor AS and redirecting it to another AS. But this signing of the AS Path is not quite enough to secure the route update, because the AS Path needs to be coupled to the actual address prefix by the route originator. The route originator needs to sign across not only the local AS and the AS to whom the route update is being sent, but also the address prefix and the expiry time of the route. This action allows the path to be “bound” to the prefix and prevents a man-in-the-middle from splicing a signed path or signed-path fragment against a different prefix.

If the signatures that “span” the AS Path in the BGP update can all be validated, then the receiver of the BGP update can validate, in a cryptographic sense, the currency of the routing update. It can also validate that the route update was propagated across the inter-AS routing space in a manner that is faithfully represented in the AS Path of the route.

The expiry time of the EE certificates used in conjunction with signed route updates introduces a new behavior into BGPsec. In the context of BGP, an announced route remains current until it is explicitly withdrawn or until the peer session that announced the route goes down. This property of BGP introduces the possibility of “ghost-route” attacks in BGP, wherein a BGP speaker fails to propagate a withdrawal in order to divert the consequent misdirected traffic from its peers.

In BGPsec, all route advertisements are given an expiry time by the originator of the route. This expiry time corresponds to the “notAfter” time of the EE certificate used to sign the protocol update, after which time the route is considered invalid. The implication is that a route originator is required to readvertise the route, and refresh the implicit expiry timer of the associated digital signature at regular intervals.

This approach to route-update validation is not quite the “light-touch” of origination validation. In this case the mechanism requires the use of a new BGP attribute and negotiation of a new BGP capability between eBGP peers, in turn meaning that the model of incremental deployment is one that is more “viral” than truly piecemeal. By “viral” we mean that this model is one of incremental deployment in which direct eBGP peers of a BGPsec-speaking AS will be able to speak BGPsec between themselves in a meaningful way. In turn these adjacent ASs can offer to speak BGPsec with their eBGP peers, and so on. This reality does not imply that BGPsec deployment must necessarily start from a single AS, but it does imply that communities of interconnected ASs all speaking BGPsec will be able to provide assurance via BGPsec on those routes originated and propagated within that community of interconnected ASs. It also implies that the greatest level of benefit to adopters of secure BGP will be realized by ASs that adopt BGPsec as a connected community of ASs.

Other changes to the behavior of BGP are implied by this mechanism. BGP conventionally permits “update packing,” where numerous address prefixes can be placed in a single update message if they share a common collection of attributes, including the AS Path. At this stage it appears that such update packing would not be supported in secure BGP, and each update in secure BGP would refer to a single prefix. Obviously this situation would have some effect on the level of BGP traffic, but early experiments suggest not at an unreasonable cost.

There are further effects on BGP that have not been fully quantified in studies to date. The addition of a compound attribute of a signature and a public key identifier for every AS in the AS Path has size implications on the amount of local storage a secure BGP speaker will need to store these additional per-prefix per-peer attributes. It also has broader implications if used in conjunction with current proposals for multipath BGP where multiple paths, in addition to the “best” path, are propagated to eBGP peers. Also, the computational load of validation of signatures in secure BGP is significantly higher in terms of the number of cryptographic operations that are required to validate a BGP update.

However, BGPsec is not intended to “tunnel” across those parts of the interdomain routing space that do not support BGPsec capabilities. When an update leaves a BGPsec realm, the BGPsec signature attributes of the route are stripped out, so the storage overheads of BGPsec are not seen by other BGP speakers.

Similarly, the periodic updates that result from the expiry timer should not propagate beyond the BGPsec realm. If the boundary is prepared to perform BGP update packing to non-BGPsec peers, then even the unpacked update overhead is not carried outside of the BGPsec realm.

It is also noted that the “full” load of BGPsec would only necessarily be carried by “transit” ASs; that is, those ASs that propagate routes on behalf of other ASs. Historically we see some 15 percent of ASs are “transit” ASs, while all other ASs behave as “stub” ASs that only originate routes and do not appear to transit routes for others. Such stub ASs can support a “lightweight” simplex version of BGPsec that can either point a default route to its upstream AS provider or trust its upstream ASs to perform BGPsec validation. In this case the stub AS needs to provide BGPsec signed originated routes to its upstream ASs, but no more.

Conclusion

The work on the specification of the RPKI itself and the specification of origin validation is nearing a point of logical completion of the first phase of standardization within the IETF, and the working draft documents are being passed from the working group into the review process leading to their publication as proposed standard RFCs. The RIRs are in the process of launching their RPKI services based on these specifications, and the initial deployment of working code has been made by numerous parties, who are also working on integration of origination validation in BGP implementations.

The work on securing the AS Path is at an earlier phase in the development process, and the SIDR Working Group is considering the initial design material. It is expected to take a similar path of further review and refinement in light of developing experience and study of the proposed approach.

The RPKI has been designed as a robust and simple framework. As far as possible, existing standards, technologies, and processes have been exploited, reflecting the conservatism of the routing community and the difficulty in securing rapid, widespread adoption of novel technologies.

Acknowledgements

The work described here is the outcome of the efforts of many individuals who have contributed to securing BGP over a period that now spans two decades, and certainly too many to ensure that all the contributors are recognized here. Instead, the authors would like to acknowledge their work and trust that the mechanisms described here are a faithful representation of the cumulative sum of their various contributions.

References

- [0] Miek Gieben, “DNSSEC: The Protocol, Deployment, and a Bit of Development,” *The Internet Protocol Journal*, Volume 7, No. 2, June 2004.
- [1] Stephen Kent, Charlie Lynn, and Karen Seo, “Secure Border Gateway Protocol (S-BGP),” *IEEE Journal on Selected Areas in Communications*, Volume 18, No. 4, pp 582–592, April 2000.
- [2] Russ White, “Securing BGP through secure origin BGP,” *The Internet Protocol Journal*, Volume 6, No. 3, September 2003.
- [3] Paul van Oorschot, Tao Wan, and Evangelos Kranakis, “On Inter-domain Routing Security and Pretty Secure BGP (psBGP),” *ACM Transactions on Information and System Security*, Volume 10, No. 3, July 2007.
- [4] Geoffrey Goodell, William Aiello, Timothy Griffin, John Ioannidis, and Patrick D. McDaniel, “Working Around BGP: An Incremental Approach to Improving Security and Accuracy of Interdomain Routing,” *Proceedings of Internet Society Symposium on Network and Distributed System Security (NDSS '03)*, February 2003.
- [5] Tony Bates, Randy Bush, Tony Li, and Yakov Rekhter, “DNS-based NLRI origin AS verification in BGP,” Internet Draft, Work in Progress, July 1998.
- [6] David Cooper et al., “Internet X.509 Public Key Infrastructure Certificate and Certificate Revocation List (CRL) Profile,” RFC 5280, May 2008.
- [7] Charlie Lynn, Stephen Kent, and Karen Seo, “X.509 Extensions for IP Addresses and AS Identifiers,” RFC 3779, June 2004.
- [8] Matt Lepinski and Stephen Kent, “An Infrastructure to Support Secure Internet Routing,” Internet Draft, Work in Progress, February 2008.
- [9] Geoff Huston, George Michaelson, and Robert Loomans, “A Profile for X.509 PKIX Resource Certificates,” Internet Draft, Work in Progress, September 2008.
- [10] Matt Lepinski, Stephen Kent, and Derrick Kong, “A Profile for Route Origin Authorizations (ROAs),” Internet Draft, Work in Progress, July 2008.
- [11] Russ Housley, “Cryptographic Message Syntax (CMS),” RFC 3852, July 2004.

- [12] Geoff Huston and George Michaelson, "Validation of Route Origination using the Resource Certificate PKI and ROAs," Internet Draft, Work in Progress, November 2010.
- [13] Randy Bush and Rob Austein, "The RPKI/Router Protocol," Internet Draft, Work in Progress, March 2011.
- [14] Kim Zetter, "Revealed: The Internet's Biggest Security Hole," *WIRED*, August 2008, <http://www.wired.com/threatlevel/2008/08/revealed-the-in/>
- [15] Geoff Huston, "Resource Certification," *The Internet Protocol Journal*, Volume 12, No. 1, March 2009.
- [16] Stephen Kent, "Securing the Border Gateway Protocol," *The Internet Protocol Journal*, Volume 6, No. 3, September 2003.

Ed.: A version of this article also appeared in *The IETF Journal*, Volume 7, Issue 1, July 2011. *The IETF Journal* can be obtained from: <http://isoc.org/ietfjournal/>

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005; he served on the Board of Trustees of the Internet Society from 1992 until 2001. E-mail: gih@apnic.net

RANDY BUSH is a Research Fellow and Network Operator at Internet Initiative Japan (IIJ), Japan's first commercial ISP. He specializes in network measurement, especially routing, network security, routing protocols, and IPv6 deployment. Randy has been in computing for 45 years, and has a few decades of Internet operations experience. He was the engineering founder of Verio, which is now NTT/Verio. He has been heavily involved in transferring Internet technologies to developing economies for more than 20 years. E-mail: randy@psg.com

Views of IPv6 Site Multihoming

by Fred Baker, Cisco Systems

In today's Internet, *site multihoming*—an edge network configuration that has more than one service provider but does not provide transit communication between them—is relatively common. Per the statistics at www.potaroo.net, almost 40,000 *Autonomous Systems* are in the network, of which about 5,000 seem to offer transit services to one or more customers. The rest are in terminal positions, possibly meaning three things. They could be access networks, broadband providers offering Internet access to small companies and residential customers; they could be multihomed edge networks; or they might be networks that intend to multihome at some point in the future. The vast majority, on the order of 75 percent, are multihomed or intend to multihome. That is but one measure; you do not have to use *Border Gateway Protocol* (BGP) routing to have multiple upstream networks. Current estimates suggest that there is one multihomed entity per 50,000 people worldwide, and one per 18,000 in the United States.

We also expect site multihoming to become more common. A current proposal in Japan suggests that each home might be multihomed; it would have one upstream connection for Internet TV, and one or more other connections provided by *Internet Service Providers* (ISPs), operating over a common *Digital Subscriber Line* (DSL) or fiber-optic infrastructure. That scenario has one multihomed entity for every four people.

Why do edge networks multihome? Reasons vary. In the Japanese case just propounded, it is a fact of life—users have no other option. In many cases, it is a result of a work arrangement, or a strategy for achieving network reliability through redundancy.

For present purposes, this article considers scaling targets derived from a world of 10 billion people (circa 2050), and a ratio of one multihomed entity per thousand people—on the order of 10,000,000 multihomed entities at the edge of the Internet. Those estimates may not be accurate 40 years from now, but given current trends they seem like reasonable guesses.

RFC 1726^[1], the technical criteria considered in the selection of what at the time was called *IP Next Generation* (IPng), did not mention multihoming per se. Even so, among the requirements are scalable and flexible routing, of which multihoming is a special case. When IPv6 was selected as the “next generation,” multihoming was one of the topics discussed. The Internet community has complained that this particular goal was not fulfilled. Several proposals have been proffered; unfortunately, each has benefits, and each has concerns. No single perfect solution is universally accepted.

In this article, I would like to look at the alternatives proposed and consider the effects they have. In this context, the goals set forth in RFC 3582^[2] are important; many people tried to state what they would like from a multihoming architecture, and the result was a set of goals that solutions only asymptotically approach.

The proposals considered in this article include:

- *Provider Independent Addressing*, also known as *BGP Multihoming*
- *Exchange-Based Addressing*
- *Shim6*, also known as *Level 3 Multihoming*
- *Identifier-Locator Network Protocol* (ILNP)
- *Network Prefix Translation*, also known as *NAT66*

BGP Multihoming

BGP Multihoming involves a mechanism relatively common in the IPv4 Internet; the edge network either becomes a member of a *Regional Internet Registry* (RIR) [APNIC, RIPE, LACNIC, AFRINIC, ARIN] and from that source obtains a *Provider-Independent* (PI) prefix, or obtains a *Provider-Allocated* (PA) prefix from one provider and negotiates contracts with others using the same prefix. In any case, it advertises the prefix in BGP, meaning that all ISPs—including in the PA case—the provider that allocated it, must carry it as a separate route in their routing tables.

The benefit to the edge is easily explained, and in the case of large organizations it is substantial. Consider the case of Cisco Systems, whose internal network rivals medium-sized ISPs for size and complexity. With about 30 *Points of Attachment* (PoAs) to the global Internet, and at least as many service providers, Cisco has an IPv6 /32 PI prefix, and hundreds of offices to interconnect using it. One possible way to enumerate the Cisco network would be to use the next five bits of its address (32 /37 prefixes) at its PoAs, and allocate prefixes to its offices by the rule that if their default route is to a given PoA, their addresses are derived from that PoA. By advertising the PoAs /37 and a backup /32 into the Internet core at each PoA, Cisco could obtain effective global routing. It would also obtain relative simplicity for its internal network—only one subnet is needed on any given *Local-Area Network* (LAN) regardless of provider count or addressing, and routing can be optimized independently from the outside world.

The problem that arises with PI addressing, if taken to its logical extreme, is that the size of the routing table explodes. If every edge network obtains a PI prefix—neglecting for the moment both BGP traffic engineering and the kind of de-aggregation suggested in Cisco’s case—the logical outcome of enumerating the edge is a routing table with on the order of 10^7 routes. The memory required to store the routing table, and in the *Secure Interdomain Routing* (SIDR) case the certificates that secure it, is one of the factors in the cost of equipment. The volume of information also affects the time it takes to advertise a full routing table, and in the end the amount of power that a router uses, the heat it produces, and a switching center’s air conditioning requirements. Thus both the capital cost of equipment used in transit networks and the cost of operations would be affected. In effect, the Internet becomes the “poster child” for the *Tragedy of the Commons*.

Exchange-Based Addressing

Steve Deering proposed the concept of exchange-based addressing at the IETF meeting in Stockholm in 1995, under the name *Metropolitan Addressing*. In this model, prefixes do not map to companies, but to Internet exchange consortia, likely regional. One organizing principle might be to associate an Internet exchange with each commercial airport worldwide, about 4000 total, resulting in a global routing table on the same order of magnitude in size. Edge networks, including residential networks, within that domain obtain their prefix from the exchange, and they are used by any or all ISPs in the region. Routes advertized to other regions, even within the same ISP, are aggregated to the consortium prefix.

The benefits to the edge network in exchange-based addressing are similar to the benefits of PI addressing for a large corporation. In effect, the edge networks served by an exchange consortium behave like the “departments” of a “user consortium,” and they enjoy great independence from their upstream providers. They can multihomed or move between providers without changing their addressing, and on a global scale the routing table is contained to a small multiple of the number of such consortia.

However, the benefit to users is in most cases a detriment to their ISPs; the ISPs are forced to maintain routes to each user network served by the consortium—or at least routes for their own customers and a default route to the exchange. Thus, the complexity of routing is moved from the transit core to the access networks serving regional consortia. In addition, if there is no impediment to a user flitting among ISPs, users can be expected to flit, imposing business costs.

The biggest short-term effect on the ISP might well be the reengineering of its transit contracts. In today’s Internet, a datagram sent by users to their ISPs is quickly shuttled to the destination’s ISPs, which then carry it over the long haul. In an exchange-based network, there is no way to remotely determine which local ISP or ISP instance is serving a given customer.

Hence, the sender's ISP carries the datagram until it reaches the remote consortium, whence it switches to the access network serving the destination. One could argue that a "sender-pays" model might have benefits, but it is very different from the present model.

The edge network has problems, too. If the edge network is sufficiently distributed, it will have services in several exchange consortia, and therefore several prefixes. Although there is nothing inherently bad about that, it may not fit the way a cloud computing environment wants to move virtual hosts around, or miss other requirements.

Level 3 Multihoming: Shim6

The IETF's *shim6* model^[9] starts from the premise that edge networks obtain their prefixes from their upstream ISPs—PA Addressing. If a typical residential or small business does so, there is no question of advertising its individual route everywhere; the ISP can route internally as it needs to, but globally, the number of ISPs directs the size of the routing table. If that is, as **potaroo** suggests, on the order of 10,000, the size of the routing table will be on the same order of magnitude.

The benefit to the ISP should be obvious; it does not have to change its transit contracts, and although there will be other concerns, it does not have the routing table ballooning memory costs or route exchange latencies.

However, as exchange-based addressing moves operational complexity from the transit core to the access network, *shim6* moves such complexities to the edge network itself and to the host in it. If a network has multiple upstream providers, each LAN in it will carry a subnet from each of those providers—not one subnet per LAN, but as many as the providers of the host's LAN will use. At this point, the ingress filtering of RFC 3704^[21] at the provider becomes a problem at the edge; the host must select a reasonable address for any session it opens, and must do so in the absence of specific knowledge of network routing. A wrong guess can have dramatic effects; a session routed to the wrong provider may not work at all, and an unfortunate address choice can change end-to-end latency from tens of milliseconds to hundreds or worse by virtue of backbone routing.

Application layer referrals and other application uses of addresses also have difficulties. Although the address a session is using will work both within and without the network, if a host has more than one address, one of the other addresses may be more appropriate to a given use. Hence, the application that really wants to use addresses is saddled with finding all of the addresses that its own host or a peer host might have.

There is also an opportunity. TCP today associates sessions with their source and destination addresses. The shim6 model, implemented in the *Stream Control Transmission Protocol* (SCTP)^[17] and *Multipath TCP* (MPTCP)^[16], allows a session to change its addresses, meaning that a session can survive a service provider outage. Doing the same in TCP requires the insertion of a shim protocol between IP and TCP; at the Internet layer, the address might change, but the shim tracks the addresses for TCP.

There are, of course, ways to solve the outstanding problems. For simple cases, RFC 3484^[3, 4] describes an address-selection algorithm that has some promise. In the Japanese case, a residential host might use link-local addresses within its own network, addresses appropriate to the television service on its TV and set-top box, and an ISP's prefix for everything else. If there is more than one router in the residential LAN serving more than one ISP, exit routing can be accomplished by having the host send data using an ISP's source address to the router from which it learned the prefix. When the network becomes more complex, though, we are looking at new routing protocols that can route based on a combination of the source and the destination addresses, and we are looking at network management methodologies that make address management simpler than it is today, adding and dropping subnets on LANs—and as a result renumbering networks—without difficulty. It also implies a change to the typical host implementing the shim protocol. Those technologies either do not exist or are not widely implemented today.

Identifier-Locator Network Protocol

The concept of separating a host's identity from its location has been intrinsic to numerous protocol suites, including the *Xerox Network Systems* (XNS), *Internetwork Packet Exchange* (IPX), and *Connectionless Network Service* (CLNS) models. In the IP community, it was first proposed in Saltzer's ruminations on naming and binding, RFC 1498^[5], and in Noel Chiappa's NIMROD routing architecture, RFC 1992^[6]. In short, a host (or a set of applications running on a host, or a set of sessions it participates in) has an identifier independent of its network topology, and sessions can change network paths by simply changing the topological locations of their endpoints. Mike O'Dell, in Internet Drafts in 1996 and 1997 called 8+8 and *GSE*, suggested an implementation of this scenario using the prefix in the IPv6 address as a locator and the interface identifier as an identifier. One implication of the GSE model is the use of a network prefix translation between an edge network and its upstream provider whatever prefix the edge network uses internally, in the transit backbone, the locator appears to be a PA prefix allocated by the ISP in question. As a result, the routing table, as in shim6, enumerates the ISPs in the network—on the order of 10,000.

The *Identifier-Locator Network Protocol* (ILNP) takes the solution to fruition, operating on that basic model and adding a *Domain Name System* (DNS) Resource Record and a random number nonce to mitigate on-path attacks that result from the fact that the *IPv6 Interface Identifier* (IID) is not globally unique.

As compared to the operational complexities and costs of PI Addressing, Exchange-Based Addressing, and shim6, ILNP has the advantage of being operationally simple. Each LAN has one subnet, when adding or changing providers no edge network renumbering is required, and, as noted, the cost of the global routing table does not increase. Additionally, it is trivial to load-share traffic across points of attachment to multiple ISPs, because the locator is irrelevant above the network layer. And unlike *IPv4/IPv4 Network Address Port Translation* (NAPT), the translation is stateless; as a result, sessions using *IP Security* (IPsec) *Encapsulation Security Protocol* (ESP) encryption can cross it.

In this case, the complexities of the network are transferred to the application itself, and to its transport. The application must, in some sense, know all of its “outside” addresses. It can learn them, of course, by using its domain name in referrals and other uses of the address; in some cases however, the application really wants to know the address itself. If it is communicating those addresses to other applications—the usual usage—the assumption that its view of its address is meaningful to its remote peer is, in the words of RFC 3582^[2], *Unilateral Self-Address Fixing* (UNSAF), and the concerns raised in RFC 2993^[7] are the result. To mitigate those concerns, ILNP excludes the locator from the TCP and *User Datagram Protocol* (UDP) pseudo-headers (and as a result from the checksum).

The implication of ILNP is, as a result, that TCP and UDP must be either changed or exchanged for other protocols such as *Stream Control Transmission Protocol* (SCTP) or *Multipath TCP* (MPTCP), and that applications must either use DNS names when referring to themselves or other systems in their network—sharply dividing between the application and network layers—or devise a means by which they can determine the full set of their “outside” addresses.

Network Prefix Translation, Also Known as NAT66

Like ILNP, *Network Prefix Translation* (NPTv6) derives from and can be considered a descendant of the GSE model. It differs from ILNP in that it defines no DNS Resource Record, defines no end-to-end nonce, and requires no change to the host, especially its TCP/UDP stacks. To achieve that, the translator updates the TCP/UDP checksum in the source and destination addresses.

If the ISP prefix is a /48 prefix, this prefix allows for load sharing of sessions across translators leading to multiple ISPs; if the ISP prefix is longer, such as a /56 or /60, the checksum update must be done in the IID, and as a result load sharing can be accomplished only across translators between the same two networks. Like ILNP and unlike IPv4/IPv4 NAT, the translation is stateless; as a result, sessions using IPsec ESP encryption can cross it.

The complexities of the network are again transferred to the application itself, but not to its transport. The application must, in some sense, know all of its “outside” addresses. Using its domain name in referrals and other uses of the address can determine these addresses; in some cases, however, the application really wants to know the address itself. If it is communicating those addresses to other applications—the usual usage—the assumption that its view of its address is meaningful to its remote peer is, again in the words of RFC 3582^[2], “UNSAFE,” and some of the concerns raised in RFC 2993^[7] result.

The implication of NPTv6 is that applications must either use DNS names when referring to themselves or other systems in their network—sharply dividing between the application and network layers—or devise a means by which they can determine the full set of their “outside” addresses. However, the IPv6 goal of enabling any system in the network to communicate with any other given administrative support is retained.

Ways Forward

From the perspective of this author, the choice of multihoming technology will in the end be an operational choice. The practice of multihoming is proliferating and will continue to do so. There is a place for provider-independent addressing; it may not in reality make sense for 40,000 companies, but it probably does for the largest edge networks. At the other extreme, shim6-style multihoming makes sense in residential networks with a single LAN; as described earlier, there are simple approaches to making that work through reasonable policy approaches.

For the vast majority of networks in between, policy suggestions that do not substantially benefit the network or users who implement them do not have a good track record. Hence, while Exchange-Based Addressing materially assists in edge network problems, there is no substantive reason to believe that the transit backbone will implement it. Similarly, although shim6 materially helps with the capital and operational expenses of operating the transit backbone, it is not likely that edge networks will implement it.

We also have a poor track record in changing host software. For example, SCTP is in many respects a superior transport protocol to TCP—it allows for multiple streams, it is divorced from network layer addressing, and it allows endpoints to change their addresses midsession.

In a 2009 “Train Wreck” workshop at Stanford University, in which various researchers argued all day in favor of the development of a new transport with requirements much like those of SCTP, the research community acted as if ignorant of it when the protocol was brought up in conversation.

NPTv6 is not a perfect solution, but this author suspects that it will be operationally simple enough to deploy and manage and close enough to the requirements of edge networks and applications that it will, in fact, address the topic of multihoming.

References

- [1] Craig Partridge and Frank Kastenholz, “Technical Criteria for Choosing IP The Next Generation (IPng),” RFC 1726, December 1994.
- [2] Joe Abley, Benjamin Black, and Vijay Gill, “Goals for IPv6 Site-Multihoming Architectures,” RFC 3582, August 2003.
- [3] Richard Draves, “Default Address Selection for Internet Protocol version 6 (IPv6),” RFC 3484, February 2003.
- [4] Arifumi Matsumoto, Jun-ya Kato, and Tomohiro Fujisaki, “Update to RFC 3484 Default Address Selection for IPv6,” Internet Draft, Work in Progress, March 2011,
<http://tools.ietf.org/html/draft-ietf-6man-rfc3484-revise>
- [5] Jerome Saltzer, “On the Naming and Binding of Network Destinations,” RFC 1498, August 1993.
- [6] Isidro Castineyra, Noel Chiappa, and Martha Steenstrup, “The Nimrod Routing Architecture,” RFC 1992, August 1996.
- [7] Tony Hain, “Architectural Implications of NAT,” RFC 2993, November 2000.
- [8] Leslie Daigle, Ed., IAB “IAB Considerations for UNilateral Self-Address Fixing (UNSAF) Across Network Address Translation,” RFC 3424, November 2002.
- [9] Erik Nordmark and Marcelo Bagnulo, “Shim6: Level 3 Multihoming Shim Protocol for IPv6,” RFC 5533, June 2009.
- [10] Ole Troan, David Miles, Satoru Matsushima, Tadahisa Okimoto, and Dan Wing, “IPv6 Multihoming without Network Address Translation,” Internet Draft, Work in Progress,
<http://tools.ietf.org/html/draft-ietf-v6ops-ipv6-multihoming-without-ipv6nat>

- [11] Margaret Wasserman and Fred Baker, “IPv6-to-IPv6 Network Prefix Translation,”
Internet Draft, Work in Progress, <http://tools.ietf.org/html/draft-mrw-nat66>
- [12] Ran Atkinson and Scott Rose, “DNS Resource Records for ILNP,” Internet Draft, Work in Progress,
<http://tools.ietf.org/html/draft-rja-ilnp-dns>
- [13] Ran Atkinson, “ICMP Locator Update message,” Internet Draft, Work in Progress,
<http://tools.ietf.org/html/draft-rja-ilnp-icmp>
- [14] Ran Atkinson, “ILNP Concept of Operations,” Internet Draft, Work in Progress,
<http://tools.ietf.org/html/draft-rja-ilnp-intro>
- [15] Ran Atkinson, “ILNP Nonce Destination Option,” Internet Draft, Work in Progress,
<http://tools.ietf.org/html/draft-rja-ilnp-nonce>
- [16] Alan Ford, Costin Raiciu, Mark Handley, and Olivier Bonaventure, “TCP Extensions for Multipath Operation with Multiple Addresses,” Internet Draft, Work in Progress,
<http://tools.ietf.org/html/draft-ietf-mptcp-multiaddressed>
- [17] Randall Stewart, Ed., “Stream Control Transmission Protocol,” RFC 4960, September 2007.
- [18] Randall Stewart, Qiaobing Xie, Michael Tuexen, Shin Maruyama, and Masahiro Kozuka, “Stream Control Transmission Protocol (SCTP) Dynamic Address Reconfiguration,” RFC 5061, September 2007.
- [19] Jon Postel, “User Datagram Protocol,” RFC 768, August 1980.
- [20] Jon Postel, “Transmission Control Protocol,” RFC 793, September 1981.
- [21] Fred Baker and Pekka Savola, “Ingress Filtering for Multihomed Networks,” RFC 3704 [BCP 84], March 2004.
- [22] David Meyer, “The Locator Identifier Separation Protocol (LISP),” *The Internet Protocol Journal*, Volume 11, No. 1, March 2008.

FRED BAKER, a Cisco Fellow, has been active in technology development and Internet standardization since the 1980s. He participated in early development of IEEE 802.1d switching and IP routing. In the IETF, he has written or edited RFCs on a variety of topics, and chaired both working groups and the IETF itself. At this time, he is the IETF's Voting Member on the U.S. NIST Smart Grid Interoperability Panel, a member of the SGIP's Architecture Committee, and co-chair of the IETF IPv6 Operations Working Group. At Cisco, his group supports research at universities; he is looked to for research advice and mentorship both within and outside the company. E-mail: fred@cisco.com

Reflecting on World IPv6 Day

by Phil Roberts, ISOC

On June 8, 2011, many websites around the world made their main webpage reachable over IPv6 for 24 hours, and many of those that did this left their sites IPv6-accessible afterward.

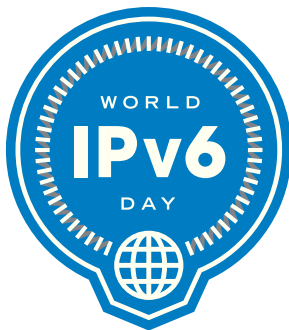
Major worldwide websites enabled IPv6 on their main page. Google enabled not only its main website but also YouTube and Blogger. Facebook and Yahoo! both enabled their main webpages as well. These websites are the five most visited websites in the world according to Alexa rankings. Other major worldwide websites that enabled IPv6 include Yahoo! Japan, Bing, Microsoft, BBC, CNN, and AOL.

Important local websites in countries around the world also joined in. In South Korea both Naver and Daum (the first and fourth most visited sites in South Korea according to Alexa) joined the event. In the Czech Republic four of the top 25 local websites joined. There were also major sites from Brazil, Portugal, and Indonesia.

Purposes

Enabling IPv6 in this way served numerous purposes:

- Network operators clearly saw that content is going to be available on IPv6. Although the major websites may not be quite there yet, it is clear that they are seriously moving in that direction.
- The industry worked to improve problems with IPv6 connectivity. Some immediate improvement resulted, and more fixes are underway to further improve IPv6 connectivity.
- Setting a public date created a deadline that accelerated deployment for many of the organizations that contacted us.
- It was important to be compared with Google, Facebook, and Yahoo!. Participants in this experiment wanted to be seen doing the same thing as the industry giants.
- This event was a clear example of how the Internet industry can work together to deploy technology that is for the good of the Internet, without intervention from outside entities. The multi-stakeholder model of Internet development continues to function well.



More than 1000 organizations contacted the Internet Society. Many of these organizations had already permanently enabled IPv6. Of the 430 or so websites the Internet Society monitored on the day, roughly two-thirds have continued to provide IPv6 access after the day.

In addition, major hosting companies enabled IPv6 for large numbers of domains, including Domain Factory, which, as a result of participating in World IPv6 Day, has made IPv6 “on by default” for all of its more than 800,000 domains. Another hosting company, Stratos, left IPv6 on after June 8 for its more than 4 million domains.

RIPE Labs did extensive measurements of IPv6 leading up to, on, and after the day, and it has published results indicating an increase in IPv6 traffic on the day—and an overall increase in IPv6 traffic also after the day.

References

- [1] Phil Roberts, “World IPv6 Day,” *The Internet Protocol Journal*, Volume 14, No. 1, March 2011.
- [2] RIPE Labs, “Measuring World IPv6 Day—Long-Term Effects,” <http://labs.ripe.net/Members/emileaben/measuring-world-ipv6-day-long-term-effects>
- [3] RIPE Labs, “Measuring World IPv6 Day—Some Glitches And Lessons Learned,” <http://labs.ripe.net/Members/emileaben/measuring-world-ipv6-day-glitches-and-lessons-learned>
- [4] RIPE Labs, “Measuring World IPv6 Day—First Impressions,” <http://labs.ripe.net/Members/mirjam/measuring-world-ipv6-day-first-impressions>

PHIL ROBERTS joined the Internet Society (ISOC) in 2008. Prior to that he spent several years with Motorola in research and product development, all in the area of mobile broadband systems. He has been active in the IETF for more than a decade. He can be reached at: roberts@isoc.org

Letters to the Editor

Hi Geoff,

Thanks you for your contribution to the March 2011 issue of *The Internet Protocol Journal*. Your description in “A Rough Guide to Address Exhaustion” and the article on “Transitional Myths” were very insightful into the whole issue of IPv4 to IPv6, and the issues concerning migration. Some of your thoughts on the migration hit home, as I am speaking to customers about the planning for the transition and I see a lot of “Got You” that I must now incorporate in my discussions with my customer.

If you do have a means of updating the technical community with activities in the area of IPv6 and how to move customers to this protocol platform, can you please point me in that direction? I like your approach and so would like to stay close to what you are doing in this area. Again, thank you for your contribution!

Ole, thanks for getting this type of information out to the technical community. Great work.

—Joel Smith, Verizon Business, Toronto, Ontario, Canada
joel.smith@one.verizon.com

The author responds:

Hi Joel,

Thank you for your comments.

Running IPv6 in a dual-stack configuration certainly presents some issues, some of which are unique to particular networks and configurations, some of which appear to be common to particular roles (such as content delivery platform, Internet Service Provider, Enterprise Provider, and end user), and some of which are common across most, if not all, circumstances.

In assisting to set up some dual-stack services a year ago, I wrote down some of the issues that I found helpful in an article: “Two Simple Hints for Dual Stack Servers” (<http://www.potaroo.net/ispcol/2010-05/v6hints.html>). You may find those hints to be of some value to your work. Some other sites that have a good collection of information are: <http://www.ipv6actnow.org/> and the community site http://www.getipv6.info/index.php/Main_Page, which also contains a wealth of information of a technical nature.

The basic guideline is to approach adding IPv6 to a network like any other engineering project: exercise care and attention to detail, and you will find it to be very straightforward!

Kind regards,

—Geoff Huston, APNIC
gih@apnic.net

Geoff and Ole,

Many thanks for your excellent papers in the March 2011 issue of IPJ. You have brought all the issues together in one place. They are clearly explained. Now I'll do my small part by suggesting to one and all that they read it. My IPv6 service comes from a manually configured tunnel from Hurricane Electric.

—Dan Cotts
dcotts@lisco.com

The author responds:

Thanks, Dan, for this feedback. It's certainly the right time for both users and content providers to act now to ensure that we continue to enjoy an Internet that still operates with a coherent end-to-end architecture into the future. The only way we can ensure that this happens is to act now and insist on IPv6—everywhere!

—Geoff Huston, APNIC
gih@apnic.net

Hello,

I enjoyed the recent IPv6 issue (Volume 14, No. 1, March 2011), but was dismayed by the lack of any frank discussion of the IPv6 “any-to-any” mantra versus the benefits of IPv4 *Network Address Translation* (NAT).

Internet purists don't hide their desire to rid the world of NAT and return to an any-to-any Internet where they could use FTP to/from any host. But for the past 15 years, NAT, RFC 1918, and perimeter security have been great for the Internet and for home and enterprise networking. When dealing with billions of endpoints, the implicit security of NAT far outweighs any alternative. Just think back to the pre-broadband/NAT days when hosts were attacked within seconds of dialing into an ISP.

Of the ~1.7 billion publicly addressed Internet devices, the vast majority would be perfectly happy behind *Carrier-Grade NAT* (CGN). In fact, as ISPs begin introducing NAT offerings, millions will stampede to them for their lower cost. Mobile phone networks are the lowest-hanging fruit, followed by residential broadband. ISPs will still offer public IP products, of course, just at a higher price point.

The IETF needs to stop pussy-footing around the issue. CGN is not just an IPv6 transitional technology; it could very well become the de facto operating standard for the next decade.

The IETF desperately needs to:

- Amend RFC 5382 (“NAT Behavioral Requirements for TCP”) to allow endpoint-independent mapping. This will improve CGN scalability by several orders of magnitude. For example, rather than 2000 hosts per public IP mentioned in Mr. Huston’s “Rough Guide” on address sharing, CGN could support 200,000 or more hosts per public IP.
- Develop an IETF standard for P2P connection establishment. It took 8+ years for the IETF to take an interest in P2P mechanics (RFC 5128). Now it’s time to show leadership. If a CGN-compatible P2P establishment standard were drafted, it would be adopted by P2P libraries overnight. While they’re at it, look at standards for tying *Universal Plug and Play* (uPnP) into CGN.
- Help coordinate a discussion of operational issues with ISP administration, law enforcement, DMCA enforcement, geolocation services, black/white lists, etc. Perhaps it’s time to extol the benefits of millisecond-accurate IPFIX logs with NAT extensions, or develop a new TCP option to embed NAT details?
- Legitimize common ISP self-preservation tactics, such as restricting SMTP, metering connections/sec, and so on.

Most importantly, IPv6 proponents should stop taking CGN as a personal affront. There is no malice; it’s simply the path of least resistance for the IPv4 conundrum.

—Craig Weinhold, Madison, Wisconsin
craig.weinhold@cdw.com

The author responds:

Thank you for your note, Craig.

The discussion of how far the Internet could scale with integration of NATs into the interior of the network as well as the current pattern of NATs at the edge is not a new discussion. The *Realm Specific IP* (RSIP) Working Group was active over a decade ago in the IETF, looking at how a network would operate that consisted of a union of distinct realms, each of which was, in address terms, a discretely addressed IP network. With the benefit of hindsight, the outcomes of that effort in supporting a case for infrastructure NATs as a long-term architectural direction for the Internet were not overly encouraging.

From the perspective of the technology community, it reinforced the conclusion that IPv6 represented the best possible response to the recognized problem of IPv4 address exhaustion. NATs were a poor compromise in so far as, at the most basic level, NATs add state into the interior of the network. This imposition of state into the network infrastructure imposes a cost in terms of service fragility and network robustness that cannot be avoided.

There was an assumption some years ago that the industry would grapple with the transition to IPv6 well before the exhaustion of IPv4 addresses, and we would never have to deal with a dual-stack transition where one-half of the dual stack, the IPv4 part, would need to operate in a mode that included infrastructure NATs. We now appear to be beyond choice here—for the Internet to continue to grow by a further 300 million new services per year at present, and grow by yet more in the coming years, there is no choice but to operate the IPv4 part of the dual-stack environment with infrastructure NATs.

But this is a short-term hack, as distinct from a tenable longer-term position. The address pool of IPv4 is not getting any larger, and as more and more new services are added into a dual-stack network, the growth in the IPv4 part of the network can be absorbed only by progressive reduction of the number of available ports to each client of the infrastructure NAT. Services become more fragile and the network becomes less resilient. The inevitable next step in progressive scarcity of IPv4 addresses in the face of such inexorable growth is to drop the entire notion of end-to-end service and introduce application-level proxies into the IPv4 network. At this point we lose any ability to further sustain an open IPv4 Internet. The only applications that could be supported are those that are supported by the application-level proxies, and all other applications simply fail. The segregation of one Internet into a number of effectively disconnected “walled gardens” of networking is a rapid outcome in such a scenario.

One of the strengths of the Internet is its openness and neutrality. The open architectural model allows novel services to be added into the network by simply equipping clients and services with the service, leaving the interior of the network untouched. The interior of the network is entirely neutral to such innovations, as it is unaware of the content or intent of the packets that are passed through its switching infrastructure.

So the long-term path of greatest common benefit to all in the Internet is a network that, as far as possible, simply vanishes! It is an Internet where content and services can rendezvous with users without having to negotiate with any network elements. It is a network that is free of toll gates. And the network has now grown to such an extent that the only path from here that can sustain that architectural simplicity and sustain yet more growth is one that shifts determinedly and rapidly to IPv6. With the limited time and resources available, attempting to improve upon NATs is, in my opinion, not the best use of the resources we can apply to this problem.

Regards,

—*Geoff Huston, APNIC*
gih@apnic.net

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ contains standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

RFC Series Editor Search Announcement

The *Internet Engineering Task Force* (IETF) is seeking an *RFC Series Editor* (RSE). The RSE has overall responsibility for the quality, continuity, and evolution of the *Request for Comments* (RFC)^[3] Series, the Internet's seminal technical standards and publications series. The position has operational and policy development responsibilities. The overall leadership and supervision of RFC Editor function is the responsibility of the RFC Series Editor. The RSE is a senior professional who must be skilled in leading, managing and enhancing a critical, multi-vendor, global information service. The following qualifications are desired:

- Leadership and management experience. In particular, demonstrated experience in strategic planning and the management of entire operations. Experience that can be applied to fulfill the tasks and responsibilities described in “RFC Editor Model (version 2)”^[1].
- Excellent written and verbal communication skills in English and technical terminology related to the Internet a must; additional languages a plus.
- Experience with editorial processes.
- Familiar with a wide range of Internet technologies.
- An ability to develop a solid understanding of the IETF, its culture and RFC process.
- Ability to work independently, via e-mail and teleconf, with strong time management skills.
- Willingness and ability to travel as required.
- Capable of effectively functioning in a multi-actor and matrixed environment with divided authority and responsibility; ability to work with clarity and flexibility with different constituencies.
- Experience as an RFC author desired.

More information about the position can be found on the RFC Editor Webpage^[2]. The RSE reports to the *RFC Series Oversight Committee* (RSOC). Expressions of interest in the position, Curriculum Vitae (including employment history), compensation requirements, and references should be sent to the RSOC search committee at rse-search@iab.org. Questions are to be addressed to the same e-mail address. Applications will be kept confidential. The RSOC will interview interested parties at the IETF meeting in Quebec City that begins July 24, 2011, but the application period is open until the position is filled.

—Fred Baker, Chair, RFC Series Oversight Committee

References

- [1] <http://www.ietf.org/id/draft-iab-rfc-editor-model-v2-02.txt>
- [2] <http://www.rfc-editor.org/rse/RSE-position.html>
- [3] Leslie Daigle, “RFC Editor in Transition: Past, Present, and Future,” *The Internet Protocol Journal*, Volume 13, No. 1, March 2010.

Global IPv6 Deployment Monitoring Survey 2011

The *Global IPv6 Deployment Monitoring Survey 2011* is now online at: <http://www.surveymonkey.com/s/GlobalIPv6survey2011>

This survey has been designed by GNKS Consult in collaboration with TNO and the RIPE NCC to further understand where the community stands on IPv6 and what needs be done to ensure that the Internet community is ready for the widespread adoption of IPv6.

Anyone can participate in this survey and we hope that the results will establish a comprehensive view of current IPv6 penetration and future plans for IPv6 deployment. The survey comprises 23 questions and can be completed in about 15 minutes. For those without IPv6 allocations or assignments or who have not yet deployed IPv6, there will be fewer questions.

The survey closes July 31, 2011. We thank you for your time and interest in completing this survey. If you have any questions concerning the survey, please e-mail: info@gnksconsult.com

For more information about the survey and links to previous year’s survey results, please see:

<https://www.ripe.net/internet-coordination/news/industry-developments/global-ipv6-deployment-monitoring-survey-2011>

RFC 6127 Published

The topic of IPv4 depletion and IPv6 deployment is covered in the recently published RFC 6127 entitled “IPv4 Run-Out and IPv4-IPv6 Co-Existence.” From the introduction: “When IPv6 was designed, it was expected that the transition from IPv4 to IPv6 would occur more smoothly and expeditiously than experience has revealed. The growth of the IPv4 Internet and predicted depletion of the free pool of IPv4 address blocks on a foreseeable horizon has highlighted an urgent need to revisit IPv6 deployment models. This document provides an overview of deployment scenarios with the goal of helping to understand what types of additional tools the industry needs to assist in IPv4 and IPv6 co-existence and transition.” RFCs can be obtained from the RFC Editor web page, see:

<http://www.rfc-editor.org/rfc.html>



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2011 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

September 2011

Volume 14, Number 3

A Quarterly Technical Publication for
Internet and Intranet Professionals

In This Issue

From the Editor	1
TRILL.....	2
IP Backhaul.....	21
Fragments	30
Call for Papers.....	31

FROM THE EDITOR

I recently attended a conference in Japan where the attendee network offered IPv6 service only. In the past, conferences such as the *Asia Pacific Regional Conference on Operational Technologies* (APRICOT) and meetings of the *Internet Engineering Task Force* (IETF) have conducted IPv6 experiments, but these have all been “opt-in” events. The conference in Japan was different: there was no IPv4 service available. Making this work involved a few manual configuration steps, but for the most part everything worked more or less the same as it did under IPv4. Some applications, including my instant message client and Skype did not work, and all connections to IPv4-only hosts needed to use *Fully Qualified Domain Names* (FQDNs) instead of IP addresses, but overall the experience gave me confidence that IPv6 is becoming a reality. As you might expect, this IPv6-only experiment also uncovered a number of bugs and incompatibilities that were duly reported to developers around the world.

Our first article is an overview of *TRansparent Interconnection of Lots of Links* (TRILL). TRILL uses Layer 3 routing techniques to create a large cloud of links that appear to IP nodes to be a single IP subnet. The protocol has been developed in the IETF and is currently being refined and enhanced in the TRILL working group. The article is by Radia Perlman and Donald Eastlake.

Developments in Internet technologies have lead to changes that go beyond the Internet itself. Not only is *Voice over IP* (VoIP) often used in place of traditional circuit-switched telephony, the telecommunication networks themselves are evolving to incorporate IP routers in place of traditional telephone switches. This evolution also applies to cellular telephone networks, specifically to what is known as *backhaul*—the transportation of voice and data from the cell sites to the mobile operators’ core networks. Jeff Loughridge explains more in “The Case for IP Backhaul.”

Once again I would like to remind you about the IPJ subscription renewal campaign. Each subscriber to this journal is issued a unique subscription ID that, coupled with an e-mail address, gives access to the subscription database by means of a “magic URL.” If your subscription has expired or you have lost your subscription ID, changed e-mail, postal mail, or delivery preference, just send an e-mail to ipj@cisco.com with the updated information and we will take care of the rest.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

Introduction to TRILL

by Radia Perlman, Intel Labs, and Donald Eastlake, Huawei Technologies

T*ransparent Interconnection of Lots of Links (TRILL)*^[1] is an *Internet Engineering Task Force (IETF)* protocol standard that uses Layer 3 routing techniques to create a large cloud of links that appear to IP nodes to be a single IP subnet. It allows a fairly large Layer 2 cloud to be created, with a flat address space, so that nodes can move within the cloud without changing their IP addresses, while using all the Layer 3 routing techniques that have evolved over the years, including shortest paths and multipathing. An early problem and applicability statement for TRILL can be found in [6]. Additionally, TRILL supports Layer 2 features such as *Virtual Local-Area Networks (VLANs)*, the ability to autoconfigure (while allowing manual configuration if so desired), and multicast/broadcast with no additional protocol.

Additionally, TRILL is evolutionary in the sense that an existing Ethernet deployment, where the links are connected with bridges, can be converted into a TRILL cloud by replacing any subset of the bridges with devices implementing TRILL. Devices implementing TRILL are called *Routing Bridges*, or *RBridges*. As bridges are replaced, nothing changes for the IP nodes connected to the cloud except that the cloud becomes more stable and uses available bandwidth more effectively.

To understand why TRILL was needed, it is helpful to explore the history of Ethernet and IP.

Network protocols are usually described in terms of *layers*. The description usually quoted in textbooks is the *Open Systems Interconnection (OSI) Reference Model*, which describes seven protocol layers^[4]. It is important to realize that the layers are useful primarily as a way to think about networking, but actual network protocols are far more complex. Layers get subdivided or combined, and often a technology usually thought of as belonging to a lower layer (for example, Layer 2) can be layered on top of a higher layer (for example, Layer 3). Most descriptions of network layers agree on the bottom four layers, and vary according to details such as whether syntax (for example, *Extensible Markup Language [XML]*^[7]), which would be a *Presentation Layer* in the OSI model, is a layer or not. Such descriptive choices do not affect how protocols are built, and luckily, for understanding of TRILL, the relevant layers to focus on are just the bottom three:

- Layer 1, *Physical Layer*: Physical, electrical, and optical specification for connectors, bit signaling, etc.
- Layer 2, *Data Link Layer*: The protocol that lets neighbor nodes on a link exchange packets
- Layer 3, *Network Layer*: The protocol that provides routing to create a path from a source node to a destination node

TRILL, as we will see, is a Layer 2 and ½ protocol: It glues links together so that IP nodes see the cloud as a single link. Therefore, TRILL is below Layer 3; but, it is above Layer 2 because it terminates traditional Ethernet clouds, just like IP routers would do.

It is definitely time to be confused. Why are there multiple links at Layer 2? Isn't that the job of Layer 3?

Evolution of Layer 2 from Point-to-Point Links to LANs

In the beginning (the 1970s or so for the purposes of this article), Layer 2 really was a direct link between neighbor nodes. Most links were point-to-point, and Layer 2 protocols primarily created *framing*—a way to signal the beginning and end of packets within the bit stream provided by Layer 1—and *checksums* on packets^[11]. For links with high error rates, Layer 2 protocols such as *High-Level Data Link Control* (HDLC)^[12] provided message numbering, acknowledgements, and retransmissions, so the Layer 2 protocol resembled, in some ways, a reliable protocol such as TCP. HDLC and other Layer 2 technologies sometimes provided an ability to have multiple nodes share a link in a master/slave manner, with one node controlling which node transmits through techniques such as polling.

Then the concept of *Local-Area Networks* (LANs) evolved, the most notable example being Ethernet. Ethernet technology enabled interconnection of (typically) hundreds of nodes on a single link in a peer-to-peer rather than master/slave relationship. Ethernet was based on CSMA/CD, where CS = *Carrier Sense* (listen before talking so you don't interrupt); MA = *Multiple Access*; and CD = *Collision Detect* (listen while you are talking to see if someone starts talking while you are so you are both interfering with each other). Interestingly, although IP had a 4-byte address and was the basis of addressing for the entire Internet, Ethernet had a larger 6-byte address, with aspirations for connecting only a small number of nodes in a fairly small region such as a single building.

The reason for the larger address space for Ethernet was to avoid the need to configure addresses when plugging nodes into a network. Instead, manufacturers of equipment would purchase blocks of Ethernet addresses and embed a unique address for each device in their hardware (the “MAC address”), and an Ethernet node would then be able to use that address in any Ethernet without fear of address collision.

Evolution of Ethernet to Spanning Tree

LANs came onto the scene with such fanfare that people came to believe that LAN technology was a replacement of traditional Layer 3 protocols such as IP. People built applications that were implemented directly on Layer 2 and had no Layer 3. This situation meant that the application would be limited by the artifacts of the Layer 2 technology, because a Layer 3 router cannot forward packets that do not contain the Layer 3 header implemented by the router.

In the case of the original Ethernet, it meant the application would work only within a maximum distance of perhaps a kilometer.

When people using technologies built directly on a LAN realized they wanted networks larger (in distance and total number of nodes) than the LAN technology allowed, the industry invented the concept of “bridges”—packet-forwarding devices that forwarded Layer 2 packets.

Forwarding Ethernet packets might seem easy because the Ethernet header looks similar to a Layer 3 header. It has a source and destination address, and the addresses are actually larger than IP addresses. But Ethernet was not designed to be forwarded. Most notably absent from the Ethernet header is a *hop count* (also sometimes referred to as a “time to live,” or TTL) to detect and discard looping packets. But other features of a typical Layer 3 protocol were also missing in Ethernet, such as an address that reflects where a node is in the topology, node discovery protocols, and routing algorithms. These features were not in Ethernet because the intention of the Ethernet design was that it be a Layer 2 protocol, confined to operation on a single link.

The transparent bridge was invented as a mechanism to forward Ethernet packets emitted by end nodes that did not implement Layer 3. Ethernet at the time had a hard packet size limit, so bridges could not modify the packet in any way.

The transparent bridge design, which met those constraints, consisted of having bridges listen promiscuously, remember the source addresses seen on each port, and forward based on the learned location of the destination address. If the destination was unknown, the packet would be forwarded onto all ports except the one that it was received on.

This simple method worked only if there was only one path between any pair of nodes. So the concept was enhanced with a protocol known as the *Spanning Tree Algorithm*.^[8] The physical topology could be an arbitrary mesh, but bridges, using the spanning-tree algorithm, would prune the topology into a loop-free (tree) topology on which data packets were forwarded. (“Spanning” means that packets can reach all the nodes.)

As Figure 1 shows, the spanning-tree concept is that an arbitrary topology could be built using Ethernet links (horizontal lines) and bridges (circles). Bridges running the spanning-tree algorithm determine a loop-free subset of the topology, and put some ports into standby (the ones that are shown in Figure 2 as dotted lines). Data packets flow on the ports that spanning tree determines should be active. This model does not yield optimal routes, as indicated in Figure 3, where packets between A and X go through the path of bridges 11, 7, 6, 2, 14, 4, and 3.

Figure 1: A Bridged Network

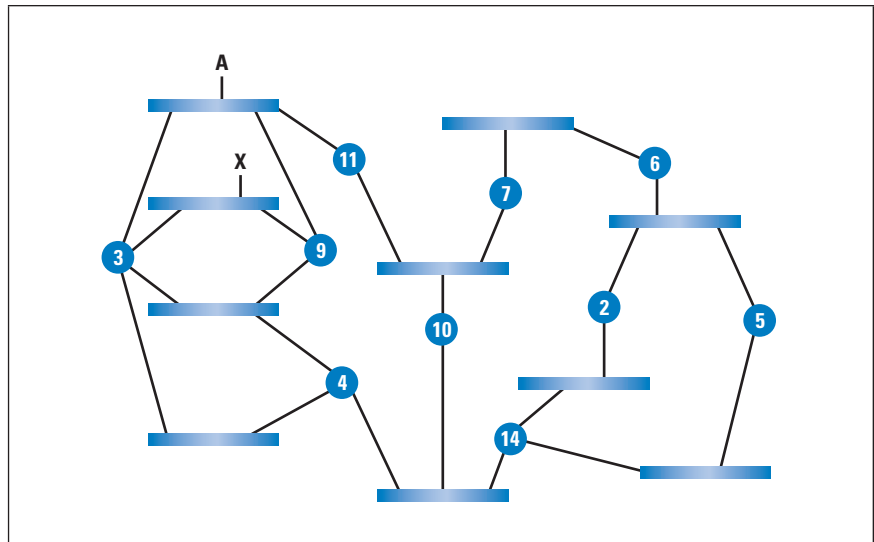


Figure 2: Bridged Network with Spanning Tree

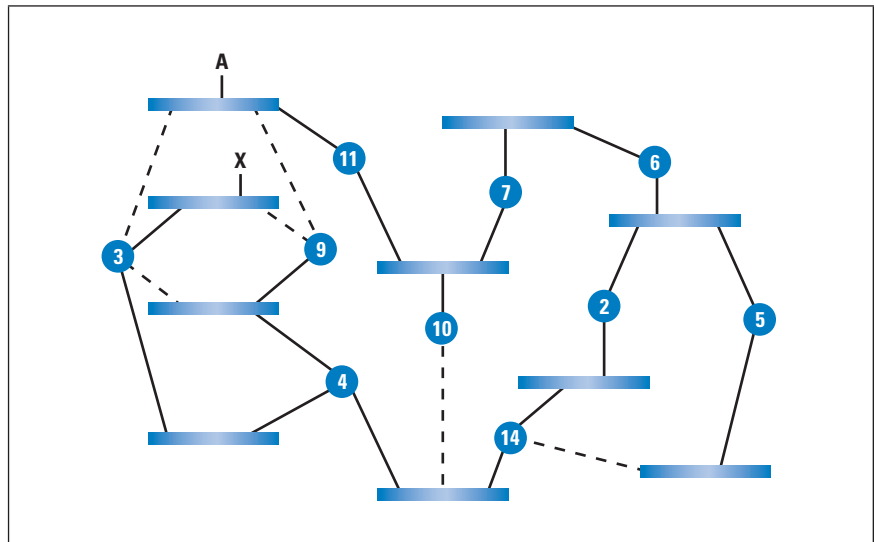
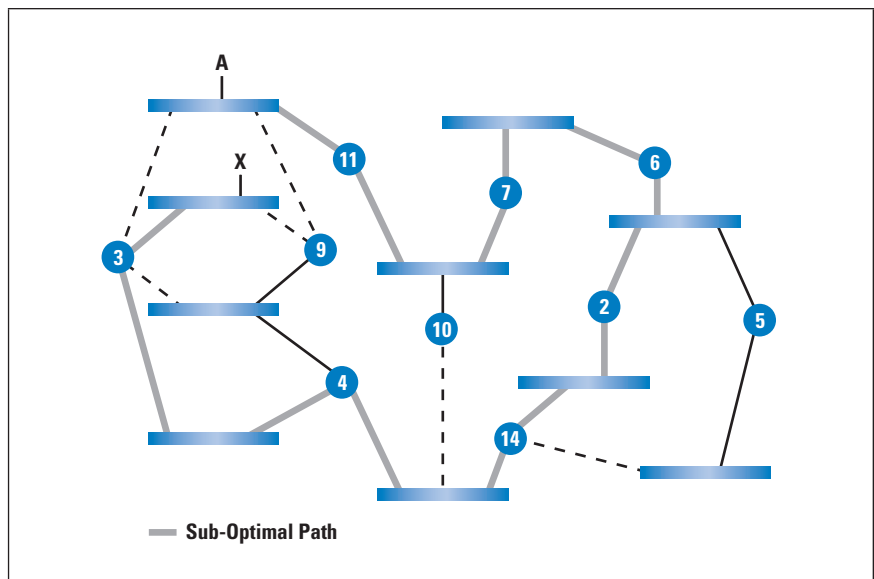


Figure 3: A Sub-Optimal Path



The spanning-tree algorithm is also inherently unstable. It requires bridges to be engineered to be able to examine every incoming packet at wire speed, to determine if the packet is a spanning-tree message, and if so, process it. The spanning-tree algorithm requires a bridge to forward unless there is a “more qualified” neighbor bridge on the link. Details of the spanning-tree algorithm, fascinating as they are, are beyond the scope of this article. If a bridge loses enough spanning-tree messages from its “more qualified” neighbor bridge because congestion overwhelms its ability to process incoming messages, the bridge will conclude that it does not have a more qualified neighbor, and therefore should start forwarding onto the link. This situation is extremely dangerous without a hop count, a field that would naturally be included in a protocol designed to be Layer 3 and forwardable.

The originally invented Ethernet, CSMA/CD, is pretty much nonexistent. Almost all Ethernet today consists of bridges connected with point-to-point links. The header still looks like Ethernet, but new fields have been added, such as VLANs discussed later in this article.

Characteristics of IP

Transparent bridging was necessitated by a quirk of history, in that applications were being built without Layer 3. But today, applications are almost universally built on top of IP. So why not replace all bridges with IP routers?

The reason is an idiosyncrasy of IP. In IP, routing is directed to a *link*, not a *node*. Each link has its own block of addresses. A node connected to multiple links will have multiple IP addresses, and if the node moves from one link to another, it must acquire a new IP address within the block for that link.

This property is not an inherent property of Layer 3, just a characteristic of IP. An alternative technology, proposed in 1992 as a replacement to IPv4, was *Connectionless-mode Network Protocol* (CLNP), an ISO packet format that had 20-byte addresses (actually, variable length). Its address, like IP, was hierarchical, routing to the longest matching address prefix in the forwarding table that matched the destination address. But in IP, the bottom level of routing was to a single link. In CLNP, the bottom level of routing consisted of routing to a cloud known as an “area,” that included lots of links (typically hundreds). Within the area, end nodes announced themselves and routers routed directly to the end node. An end node could move within an area without changing its Layer 3 address. Routers within an area would not need to be configured.

In contrast, with IP, a block of IP addresses needs to be carved up to assign a unique block to each link, IP routers need to be configured with the address block for each of their ports, and nodes that move from one link to another have to change their Layer 3 addresses. Therefore, it is still popular to create large bridged Ethernets, because a bridged set of links looks to IP like a single link.

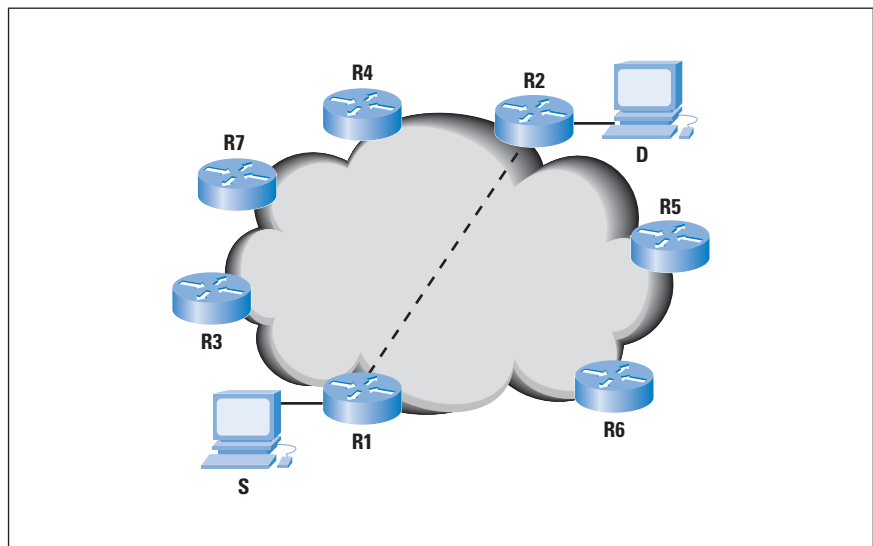
TRILL: Best of Both Worlds

TRILL allows the ease of configuration of Ethernet while benefitting from the routing techniques provided at Layer 3. It also coexists with existing bridges; it is not necessary to replace all the bridges in an Ethernet, but the more bridges replaced by RBridges, the better the bandwidth usage and the more stable the cloud becomes (because the spanning trees get smaller and smaller, and ultimately disappear if all bridges are replaced by RBridges).

Figure 4 shows the basic concepts in TRILL handling a unicast packet where the location of the destination is known:

- RBridges run a link state routing protocol, which gives each of them knowledge of the topology consisting of all the RBridges and all the links between RBridges. Using this protocol, each RBridge calculates shortest paths from itself to each other RBridge, as well as trees for delivering multidestination traffic.
- When an RBridge, R1, receives an Ethernet frame from an end node S, addressed to Ethernet destination D, R1 encapsulates the frame in a TRILL header, addressing the packet to the RBridge R2, to which D is attached. The TRILL header contains an “ingress RBridge” field (R1), an “egress RBridge” field (R2), and a hop count.
- When R2 receives the encapsulated packet, R2 removes the TRILL header and forwards the Ethernet packet on to D.

Figure 4: RBridging



What the TRILL header looks like, how R1 knows that R2 is the correct “egress RBridge,” and some of the concepts in the link state protocol *Intermediate System-to-Intermediate System* (IS-IS) are described in the next section. We also explain how TRILL handles multidestination frames, VLANs, and IP Multicast.

The TRILL Header

The main fields in the TRILL header are: ingress RBridge nickname (16 bits), egress RBridge nickname (16 bits), hop count (6 bits), and a multidestination flag bit (1 bit). A typical Layer 3 header would contain a source, a destination, and a hop count. So TRILL is basically an encapsulation header with flat 16-bit addresses. How RBridges obtain “nicknames” is described later in this article.

This header is very simple for core RBridges to forward, compared with either an IP or an Ethernet header. The destination field is just 16 bits, so it can be a simple table lookup to find the entry in the output port, as opposed to the Ethernet 6-byte destination, which typically requires content-addressable memory or hashing, or the longest prefix matching of IP.

Learning End-Node Locations

How does R1 know that R2 is the correct egress RBridge for some destination D? The default mechanism is learning the correspondence between (ingress RBridge, source MAC address) when the egress RBridge decapsulates a packet. If R1 does not know where the destination MAC is located, R1 encapsulates the packet in a TRILL header with the multidestination flag set, indicating that it should be transmitted through a tree to all the RBridges.

An additional mechanism, which is optional, is known as *End-Station Address Distribution Information* (ESADI). ESADI allows R1 to announce some or all of the end nodes that are attached to R1. Both announcing to and listening to ESADI are optional. This mechanism has advantages over flooding and learning from data packets:

- ESADI packets can have cryptographic protection.
- R1 might have a more definite reason to know that S is attached to R1 than simply seeing a packet with the S address in the header. For instance, R1 might have been configured to lock down a port to the S MAC address. Or there might be a cryptographically protected enrollment protocol when S attaches to R1.
- R1 might be able to have tighter timers on verifying the location of local end nodes; for instance, if they are IP nodes, R1 might be able to ping them.

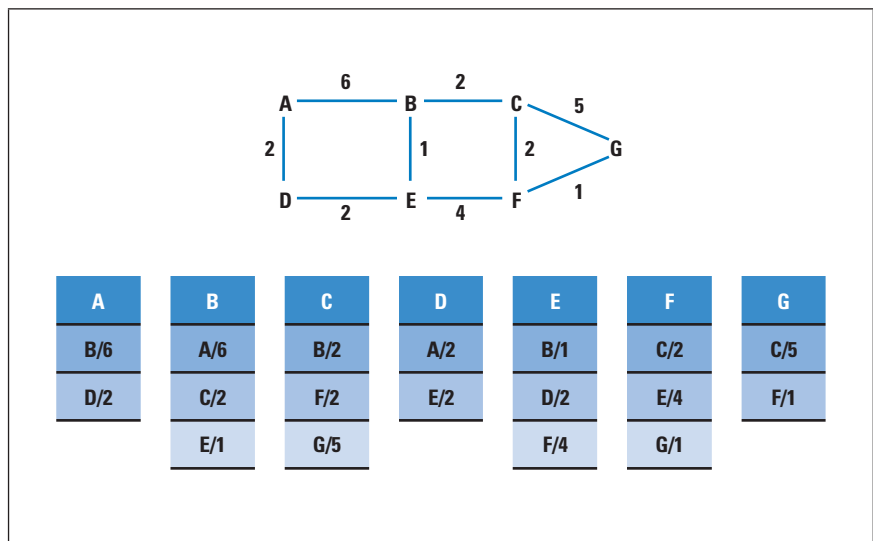
It is also possible to have a directory that lists not only (RBridge nickname, {set of attached end-node MAC addresses}) but also {(end-node IP address, end-node MAC address)} pairs. The first RBridge, or a *hypervisor*, or the end-node process itself, might query the directory about the destination, and encapsulate packets, rather than flooding, and thus also be able to bypass the *IPv4 Address Resolution Protocol* (ARP) and the *IPv6 Neighbor Discovery* (ND) protocols.

Link State Protocols

A *link state* protocol is a routing protocol in which each router R determines who its neighbors are, and broadcasts (to the other routers) a packet, known as a *Link State Packet* (LSP), that consists of information such as “I am R,” and “My neighbor routers are X (with a link cost of c1), Y (cost c2), and Z (cost c3).” The commonly deployed link state protocols are *Intermediate System-to-Intermediate System* (IS-IS)^{[2][9]} and *Open Shortest Path First* (OSPF)^[10]. IS-IS, designed in the 1980s to route DECnet, was adopted by the *International Organization for Standardization* (ISO). IS-IS can route IP traffic and is used by many *Internet Service Providers* (ISPs) to route IP. IS-IS was a natural choice for TRILL because its encoding easily allows additional fields, and IS-IS runs directly on Layer 2, so that it can autoconfigure, whereas OSPF runs on top of IP and requires all the routers to have IP addresses.

Figure 5 shows a small network (at the top), consisting of 7 routers. In the bottom half of the figure, the LSP database is shown; all the routers have the same LSP database because they all receive and store the most recently generated LSP from each other router. The LSP database gives all the information necessary to compute paths. It also gives enough information for all the routers to calculate the same tree, without needing a separate spanning-tree algorithm. As we will see, TRILL requires a tree (at least one tree) for distribution of multdestination packets.

Figure 5: Router Network and Link State



Acquiring Nicknames

Given that the most recently generated link state packet of each RBridge is broadcast to, and stored by, each other RBridge, it is possible to spread other information through the link state packets, such as a protocol for acquiring a unique nickname. Each RBridge chooses a nickname at random, avoiding nicknames already acquired by other R Bridges (as discovered by examining the LSP database).

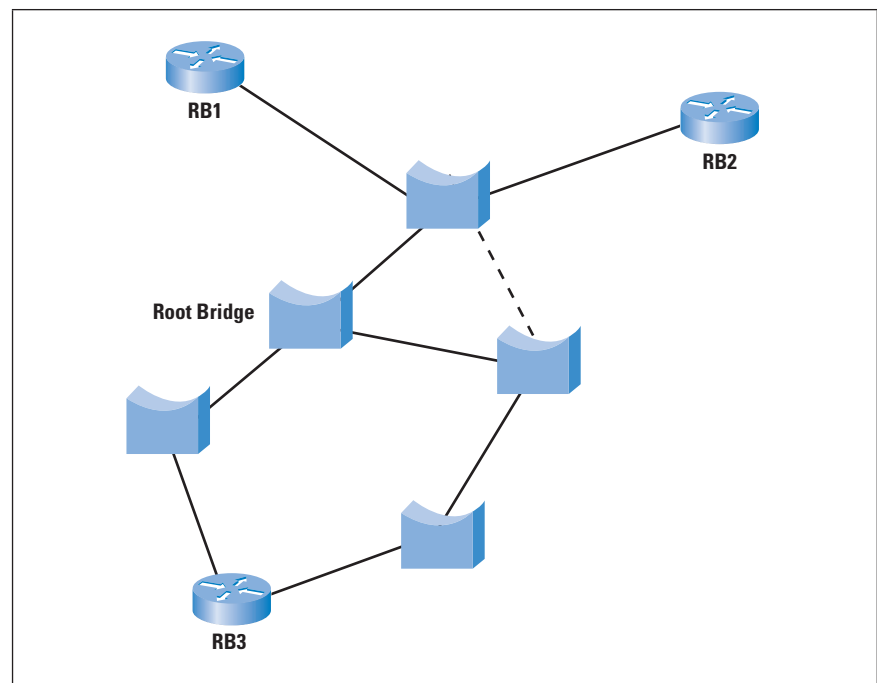
If two RBridges choose the same nickname, there is a tie-breaker, based on configured priority and 6-byte system ID. One of the RBridges gets to keep the nickname and the other RBridge has to choose another nickname that appears not to be in use.

It is possible to configure RBridges with nicknames, in which case a configured nickname takes priority over one that was randomly chosen. And in the case of misconfiguration, where two RBridges have been configured with the same nickname, again, ID and priority choose a winner, and the other one has to choose a different nickname.

Mixing RBridges with Bridges

TRILL is designed so that any subset of bridges in an Ethernet can be replaced by RBridges. A set of links connected by bridges will be perceived by RBridges as a single shared link connecting the RBridges on that link. The bridges inside that link will behave as ordinary bridges, forming a spanning tree and forwarding packets along that tree. Figure 6 illustrates an Ethernet connected by several bridges, with one port (indicated by the dashed line) selected by the spanning tree as being in backup.

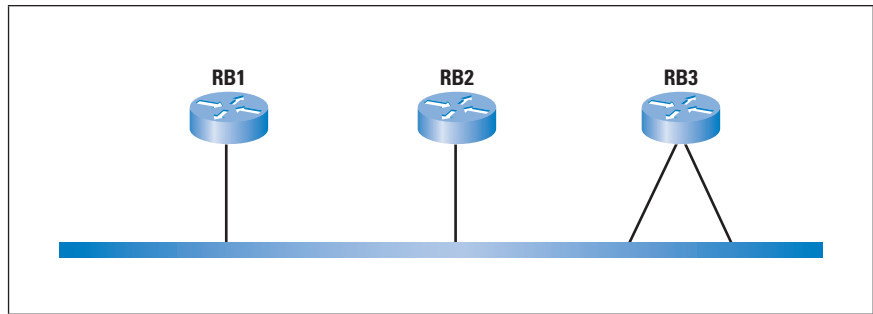
Figure 6: RBridges Connected by Bridged LAN



The RBridges RB1, RB2, and RB3 perceive the link as in Figure 7, a single shared link, on which RB3 has two ports.

Introducing RBridges into a bridged Ethernet partitions the spanning trees into smaller spanning trees. RBridges operate on a topology consisting of the RBridges themselves, connected with “links” that are either bridged Ethernets or point-to-point links.

Figure 7: Figure 6 as Perceived by RBridges: a Single Shared Link Where RB3 Has 2 Ports onto the Same Link



Link Types and the Hop-by-hop Header

In addition to the TRILL header, when RBridge R1 is forwarding a TRILL-encapsulated frame to neighbor RBridge R2, there is an additional header that is specific to the type of link connecting R1 and R2. Although TRILL carries Ethernet inside, a link between two or more RBridges could be an arbitrary type of link; for example, besides Ethernet, it could be a *Point-to-Point Protocol* (PPP) link^[13], an IP or *IP Security* (IPsec) tunnel, *Multiprotocol Label Switching* (MPLS) path, etc.

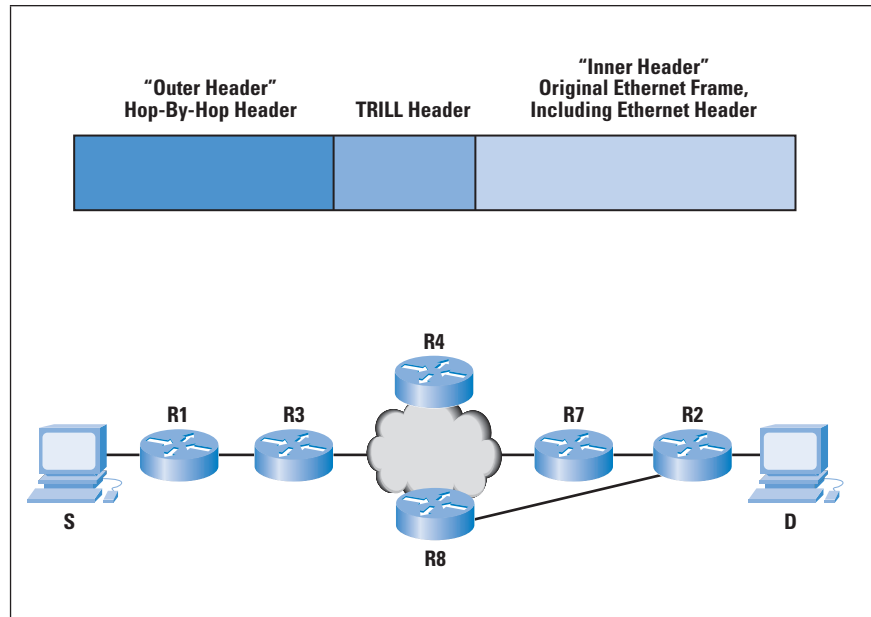
If the link is an Ethernet link, the “outer” header is an Ethernet header. If it is a PPP link, the outer header is a PPP header. The outer Ethernet header (on an Ethernet link) serves two purposes:

- If there are bridges on the link, they will perceive the packet as a normal Ethernet packet, and forward it through the spanning tree. The learning tables of the bridges on the link will see only the addresses of the RBridges on that link.
- It allows R1, when forwarding onto a link with multiple neighbors (say R2 and R3), to specify which of R2 or R3 is chosen by R1 to forward the packet by unicasting the packet to the chosen next-hop RBridge. For example, it could be that both R2 and R3 are equal costs to the destination, so R1 would need to specify which of them should forward the packet. Otherwise, both might forward the packet, and the packet would be duplicated.

So, as illustrated in Figure 8, a TRILL-encapsulated packet might have three headers:

- The outer header, or hop-by-hop header, which is stripped off at each hop, is specific to the type of link connecting neighbor RBridges, and, when forwarded between R1 and R2, it specifies R1 as source and R2 as destination
- The TRILL header, which similarly to a Layer 3 header remains in place as the packet travels from the first RBridge to the last RBridge, specifying the first RBridge (the one that encapsulated the packet with a TRILL header) as the ingress RBridge, and the last RBridge (the one that will decapsulate the packet) as the egress RBridge
- The inner Ethernet header, which specifies the communicating end-node pair as source and destination

Figure 8: TRILL Packet Headers



Again referring to Figure 8, assume S transmits an Ethernet packet to D. In the inner Ethernet header, Source = S, Destination = D.

R1 encapsulates it with a TRILL header, where ingress RBridge = R1 and egress RBridge = R2. R1 forwards it to R3, putting on a link header appropriate to the link. If the link is an Ethernet link, the outer Ethernet header will indicate S = R1, D = R3. When R3 forwards to R7, R3 leaves the TRILL header as is (other than decrementing the hop count), strips the outer header, and puts in a new outer header indicating S = R3, D = R7. Likewise, R7 forwards to R2. If it is a PPP link, there is no source or destination. When R2 forwards to D, R2 strips off the TRILL header and D sees the Ethernet packet exactly as transmitted by S.

VLANs

Ethernet has a concept known as a *Virtual LAN* (VLAN), which partitions communities of end nodes sharing the same infrastructure (links and bridges), such that end nodes in the same set can talk directly to each other (using Ethernet), whereas those in different VLANs have to communicate through a router. IP nodes, although generally unaware of Ethernet VLAN tags, perceive different VLANs to be different IP subnets.

Typically, a bridge is configured with a VLAN for each port, and the bridge adds a tag to the Ethernet header that indicates which VLAN the packet belongs to. A bridge with a port that is configured to be VLAN x will deliver only packets tagged as VLAN x to that port, and will usually strip the VLAN tag before forwarding.

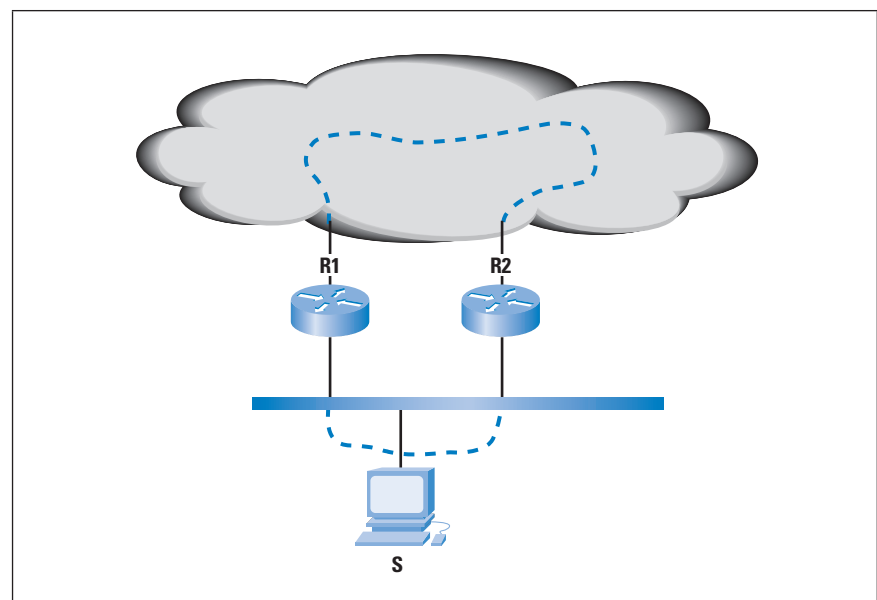
The original Ethernet did not have a VLAN concept. In today's Ethernet standard, each packet must be associated with a VLAN. A bridge might be configured with a default VLAN for a port, meaning that if no VLAN tag is in the packet, the bridge will treat it as if it is that default VLAN. A bridge B might be configured in various ways that make VLANs more complex:

- B might be configured to drop a set of VLANs rather than forward them onto a particular port, even though the port is a transit port.
- B might be configured to modify the VLAN tag to a different value when forwarding from one port to another.
- B might be configured to remove the VLAN tag when forwarding onto a particular port.

Appointed Forwarders

If there are multiple RBridges on the same link, together with end nodes, it is important that only one of them encapsulate a packet from an end node. As illustrated in Figure 9, if both R1 and R2 were to encapsulate a unicast packet from S, two copies would be delivered to the destination. However, if S were to transmit a multidestination packet (such as a multicast, or an unknown destination), then the copy that R1 encapsulates would be forwarded through the campus, received by R2 (which likely would not know that the packet originated on its port to R1), and R2 would decapsulate it. Then R1 would see a native packet from S, exactly as the first copy, and again encapsulate it and send it into the campus.

Figure 9: Link with Multiple RBridges.
Note: No Hop Count Protection on Native Frame.



The hop count in the TRILL header would not solve this loop, because the hop count does not exist while the packet is not encapsulated with a TRILL header.

IS-IS has an election protocol in which one of the RBridges is elected as the *Designated RBridge* (DRB). In order to allow load-splitting the task of encapsulating and decapsulating traffic, the DRB may delegate the job of encapsulation/decapsulation based on VLAN. In other words, if R1 is DRB, R1 can delegate to R2 the task of encapsulating/decapsulating traffic for a set of VLANs, say VLANs x, y, and z, and delegate to R3 a different set of VLANs, and R1 might handle the rest.

Implications of VLANs on TRILL

TRILL treats VLANs strictly as a way of partitioning the end nodes, in contrast with IEEE, which allows bridges to drop transit traffic based on VLAN. Consequently, an Ethernet link connecting TRILL RBridges R1 and R2 might be able to deliver packets tagged with VLAN x, but not deliver packets tagged with VLAN y.

It is important, as shown in Figure 9, that all the RBridges on a link know about each other; otherwise they might both encapsulate a packet.

The IS-IS election is done through Hello messages, whereby RBridges announce themselves. Unfortunately, possible configuration of bridges, whether intentional or by mistake, can partition a link for traffic marked as VLAN y, but have the link be connected for traffic marked as VLAN x. This situation complicates the IS-IS election. When transmitting a Hello message onto an Ethernet link, an RBridge R1 must assign it to a VLAN. If R1 chooses VLAN y, its neighbor R2 might not see the Hello message. And then, unaware that there were multiple RBridges on the link, both R1 and R2 might encapsulate a VLAN x packet.

TRILL handles this situation by having the DRB (by default) transmit Hello messages on all the VLANs for which it is enabled on the port. The DRB chooses a VLAN, say VLAN A, for inter-RBridge communication on the link, and informs the other RBridges on the link that they should use VLAN A. The other RBridges transmit IS-IS messages (including Hello messages and LSPs) and encapsulated TRILL packets, putting VLAN A in the outer header. The VLAN tag in the inner header is the one that represents the community that the end node belongs to. The VLAN tag in the outer header is only for the purpose of traversing an Ethernet hop between RBridges.

Additionally, (by default), an RBridge that is Appointed Forwarder for a VLAN, transmits Hello messages on that VLAN.

If it is known that there are no bridges, the RBridges (including the DRB) can be configured to send Hello messages only on the single VLAN specified by the DRB.

Modified Hello Protocol

IS-IS has an election protocol in which routers (or R Bridges in the case of TRILL) send Hello messages. Not only does the Hello message transmitted by R1 announce R1 to its neighbors, but the R1 Hello message contains a list of neighbors that R1 has heard Hello messages from. R2 will not consider R1 to be a neighbor unless R2 sees itself listed in the Hello messages of R1, indicating connectivity is two-way. When choosing a DRB, R2 ignores any routers for which connectivity to R2 is not two-way. Therefore, if there were a shared link with strange connectivity properties, the routers on the link might partition into cliques, each with its own DRB, each clique representing a separate link to the rest of the routers.

A surprising aspect of the use of IS-IS for TRILL was that the Hello protocol had to be modified slightly. In Layer 3 IS-IS, Hello messages are padded to the maximum size, because a possible hardware failure mode was that a link between R1 and R2 might be able to transmit small packets, but not large packets. In Layer 3, the IS-IS assumption was that R1 and R2 would rather not see that they were potential neighbors than use a flaky link. In IS-IS, LSP packets can be fragmented only by the source R1. All routers agree upon the maximum size of an LSP fragment that is guaranteed to be able to traverse all the links. Links that cannot forward packets of that size are not reported in the topology, and indeed, in Layer 3 IS-IS, would not even be discovered in the topology, because the Hello message (padded to that size) would not be seen by the neighbor router.

But with TRILL, it is important that only a single R Bridge be elected DRB, because the DRB determines which R Bridge will encapsulate/decapsulate packets for each VLAN. One of the first implementations of TRILL wound up forming a loop, where two R Bridges, R1 and R2, both performed encapsulation/decapsulation. This situation resulted because neighbors R1 and R2 did not see each other's Hello messages, because the R1 Hello, padded to classic Ethernet maximum size by R1, became too large to forward when a VLAN tag was added, so did not reach R2.

To ensure that only a single R Bridge on a link would be elected DRB, TRILL modified the Hello protocol as follows:

- Limit the size of Hello messages and do not pad them (in order to remove artificial impediments to receipt by neighbors).
- Elect a DRB based solely on priority (not two-way connectivity as in Layer 3 IS-IS). In other words, defer to a higher-priority R Bridge R1 even if R1 does not list you as a neighbor.
- Have a separate mechanism for probing, using packets of different sizes, to see what size packets can be forwarded on the link.

In addition to solving the multiple-DRB problem, this design enables TRILL to discover which links can handle jumbo-grams, so that paths can be engineered that can forward jumbo-grams.

If the link between R1 and R2 is not acceptable because it cannot handle the assumed LSP fragment size, or because connectivity is not two-way, the link is not reported in LSPs. The capability of a link to handle larger sizes can be reported in LSPs.

There was enough confusion about this minor change to the Hello protocol, and skepticism that the Hello mechanism, which has worked correctly for Layer 3 for decades, would need to be modified for TRILL, that an additional RFC was written [3] to specifically explain the TRILL Hello mechanism.

Multidestination Frames

Multiple Trees

The original design for TRILL had the RBridges compute a single, shared tree, based on the LSP database, and all multidestination traffic was forwarded along that tree. But, to be able to load-split the use of links for multidestination traffic, a facility for using multiple trees was added early in the development of the TRILL standard.

In TRILL, the RBridge with the highest priority to be a TREE root announces to the other RBridges (through its LSP) how many trees, and which trees, should be calculated. A tree is calculated as a tree of shortest paths from a given Root, with a deterministic tie-breaker so that all RBridges calculate the same tree. The Root can be an RBridge or a pseudonode. In some cases, a Root is particularly well-situated in the topology such that its tree forms good paths for all pairs of nodes, but it is desirable to have multiple different trees, choosing different tie-breaker links, calculated from the same Root. TRILL accomplishes this setup by having that Root acquire multiple nicknames, one for each tree, and using the tree number in the tie-breaker algorithm, so that although all the trees from that Root will still be shortest-path trees, different links will be chosen in the different trees.

When R1 encapsulates a multidestination frame, R1 sets the “multidestination” flag and specifies the tree Root nickname in the “egress RBridge” field in the TRILL header.

Filtering

A multidestination frame will be tagged with a VLAN (in the inner header). The frame need not be delivered to all RBridges—just those that are connected to a port with end nodes in that VLAN. So RBridges announce, in their LSPs, which VLANs they are attached to, where “attached to,” means that they are acting as Appointed Forwarder.

Additionally, TRILL provides for filtering based on Layer 2 MAC addresses derived from IP Multicast groups. RBridges announce the set of such MAC addresses they wish to receive. The first RBridge that accepts an IP Multicast control message, such as *Internet Group Management Protocol* (IGMP), snoops on it [5] and learns what multicast listeners or multicast router is attached. This snooping is used so R1 can report in its LSP the IP Multicast groups it wishes to receive (or all groups if a multicast router is attached).

One other refinement to multdestination is the *Reverse Path Forwarding* (RPF) check. To safeguard against loops, when R is calculating which subset of its ports belong to a particular tree, R also calculates, for each port, the set of ingress RBridges whose traffic on that tree should arrive on that port.

So, the processing of a multdestination frame received by R, with TRILL header indicating Ingress = R1 and Egress/tree Root = R2, is as follows:

- If the port on which R receives the packet is not included in the tree “R2,” discard the packet.
- If the port on which R receives the packet is in tree R2 but R1 is not listed in the RPF information for that port for tree R2, discard the packet.
- For each other port in R2, if the specified VLAN is reachable through that port and the IP Multicast address is requested by an RBridge along the path through that port, forward the packet on that port.

IS-IS Pseudonodes

If there is a link with N RBridges, rather than modeling the link as having on the order of N^2 links to be reported in LSPs, IS-IS has the DRB model the link as a pseudonode. The DRB gives the pseudonode a name, and the RBridges on the link report connectivity just to the pseudonode. The DRB generates an LSP on behalf of itself, reporting connectivity to the pseudonode, but additionally generates an LSP on behalf of the pseudonode, reporting connectivity to all the RBridges on the link. This portion of IS-IS is as designed from the beginning (from its origin as Phase V DECnet routing).

When IS-IS was originally designed, Ethernets tended to be very large shared links. But today, most Ethernets are simply point-to-point links (unless there are bridges making them appear to be shared links). So it would be wasteful for RBridges to always create a pseudonode for each Ethernet link. In Layer 3 it is not as unreasonable to always treat an Ethernet as a large shared link because an “Ethernet” link, as perceived by Layer 3, is likely to be a large collection of point-to-point links glued together with either bridges or RBridges.

But RBridges are likely to often see Ethernet links with just a single neighbor, especially in a topology with no bridges. So TRILL has the ability for the DRB to specify to its neighbor RBridges whether to report the link as a pseudonode or to report connectivity to all the RBridge neighbors as separate links. By default, the DRB R sets a flag known as the “bypass pseudonode” flag in its Hello message on the link, unless at some point since R rebooted R has seen two simultaneous neighbor RBridges on that link. With this mechanism, true point-to-point Ethernet links will be reported as a link between R1 and R2 rather than a pseudonode P, with links R1–P, R2–P, and P–R1 and P–R2 reported.

TRILL Implementations

TRILL is being widely implemented. TRILL fast-path hardware is included in chips available from all major merchant silicon manufacturers. A successful interoperability test was held at the University of New Hampshire *InterOperability Laboratory* in late 2010, and TRILL products are announced and shipping.

Future Potential TRILL Enhancements

Here are just three enhancements to TRILL being considered:

- Data centers require more VLANs than can be specified in 12 bits with a single VLAN tag. A TRILL extension to optionally include the ability to encode 24 bits of VLAN-like labeling in TRILL data frames is being considered.
- By optionally giving a pseudonode a nickname and having the appointed forwarder use that nickname in the ingress RBridge field, if the appointed forwarder changes, the end-node learning cache of distant RBridges will still be correct.
- A proposal is being made allowing IS-IS to be hierarchical in a TRILL campus. IS-IS hierarchy partitions the LSP database so that any single RBridge LSP database will be smaller, its path computation will be less computation-intensive, and it will lower the amount of LSP traffic. In particular, it shields the effects of a link that is cycling quickly from most of the campus, because only the RBridges in the region with the link will see reports of the state of that link.

Summary

The TRILL standard creates a cloud with a flat Ethernet address, so that nodes can move around within the cloud and not need to change their IP address. Although nodes attached to the cloud perceive the cloud as an Ethernet while the packet is traversing the cloud, it is encapsulated with a TRILL header, which like a Layer 3 technology, contains a source (ingress RBridge), destination (egress RBridge), and hop count. The addresses in the TRILL header are 16 bits, enabling a TRILL campus to support 64,000 RBridges. Transit RBridges do not learn about location of end nodes—only the existence of, and path to—other RBridges.

TRILL can use all the Layer 3 techniques, including shortest paths, *Equal Cost Multipath* (ECMP), and traffic engineering. It also supports VLANs and multicast. TRILL can calculate multiple trees, so that multidestination traffic can be split across links. Multidestination frames can be filtered based on VLAN and IP (v4 or v6) Multicast groups.

TRILL is compatible with existing Ethernet bridges (switches), so a bridged Ethernet can be gradually upgraded by replacing any subset of the bridges with RBridges. The more that are upgraded, the better the bandwidth usage, and the more stable the network becomes.

References

- [1] Perlman, R., Eastlake 3rd, D., Dutt, D., Gai, S., and A. Ghanwani, "Routing Bridges (RBridges): Base Protocol Specification," RFC 6325, July 2011.
- [2] "Information technology—Telecommunications and information exchange between systems—Intermediate System to Intermediate System intra-domain routing information exchange protocol for use in conjunction with the protocol for providing the connectionless-mode network service (ISO 8473)," ISO/IEC 10589:2002.
- [3] Eastlake 3rd, D., Perlman, R., Ghanwani, A., Dutt, D., and V. Manral, "Routing Bridges (RBridges): Adjacency," RFC 6327, July 2011.
- [4] ITU-T, "X.200: Information technology—Open Systems Interconnection—Basic Reference Model: The basic model," July 1994.
- [5] Christensen, M., Kimball, K., and F. Solensky, "Considerations for Internet Group Management Protocol (IGMP) and Multicast Listener Discovery (MLD) Snooping Switches," RFC 4541, May 2006.
- [6] Touch, J. and R. Perlman, "Transparent Interconnection of Lots of Links (TRILL): Problem and Applicability Statement," RFC 5556, May 2009.
- [7] W3C, "XML Base (Second Edition)," W3C Recommendation 28 January 2009,
<http://www.w3.org/TR/2009/REC-xmlbase-20090128/>
- [8] Perlman, R., "A Protocol for Distributed Computation of a Spanning Tree in an Extended LAN," *9th Data Communications Symposium*, Vancouver, 1985.
- [9] Callon, R., "Use of OSI IS-IS for routing in TCP/IP and dual environments," RFC 1195, December 1990.

- [10] Moy, J., “OSPF Version 2,” RFC 2328, April 1998.
- [11] Simpson, W., “The Point-to-Point Protocol (PPP),” RFC 1661, July 1994.
- [12] http://www.interfacebus.com/HDLC_Protocol_Description.html
- [13] Carlson, J. and Eastlake 3rd, D., “PPP Transparent Interconnection of Lots of Links (TRILL) Protocol Control Protocol,” RFC 6361, August 2011.

RADIA PERLMAN is a Fellow at Intel Labs, working on the design of various network routing and security protocols. She is the inventor of the Spanning Tree Algorithm, the designer of IS-IS, and the original concept for TRILL. She is the author of the textbook *Interconnections: Bridges, Routers, Switches, and Internetworking Protocols*. She is an IEEE Fellow and holds a Ph.D. from MIT.
E-mail: radiaperlman@gmail.com

DONALD EASTLAKE 3rd is Co-Chair of the IETF TRILL Working Group and a voting member of IEEE 802.1. He is the author of 56 IETF RFCs and a Principal Engineer with Huawei Technologies working on advanced network product research. Previously, he was a Principal Engineer at Cisco Systems and before that a Distinguished Member of Technical Staff at Motorola Laboratories, working on network protocols, security, and mesh networking.
E-mail: d3e3e3@gmail.com

The Case for IP Backhaul

by Jeff Loughridge, Brooks Consulting LLC

In any hierarchical network, designers must specify how the access layer delivers traffic to the core. In *Mobile Network Operator* (MNO) networks, the transport of voice and data from the cell sites to the wireless MNOs' core networks is called *backhaul*. *Time Division Multiplexing* (TDM) backhaul has dominated backhaul deployments since the inception of wireless communication. Leasing the backhaul access of multiple T1s/E1s for every cell site becomes prohibitively expensive in terms of operating expenses, particularly for providers that do not own the last mile. Today's 3G/4G cellular technologies have spurred a major change in the backhaul network: the transition from TDM to packet backhaul.

Ethernet is the most widespread packet-based backhaul technology. While this service is a vast cost and scale improvement over TDM backhaul, carrier Ethernet is a stepping stone in the evolution of backhaul networks. Expect MNOs to move to true IP backhaul networks to meet the scalability needs of their expanding networks. In this article, we will explain mobile backhaul evolution, shortcomings in carrier Ethernet backhaul, and how evolving service requirements will motivate cell site backhaul vendors to add IP-awareness to their networks.

Legacy Backhaul

Cellular systems were initially designed to carry only voice traffic. Since transporting digitized voice was a mature and well-understood technology, there was no need to take a divergent path for the backhaul of voice traffic in early cellular systems. Using TDM had obvious advantages among those being:

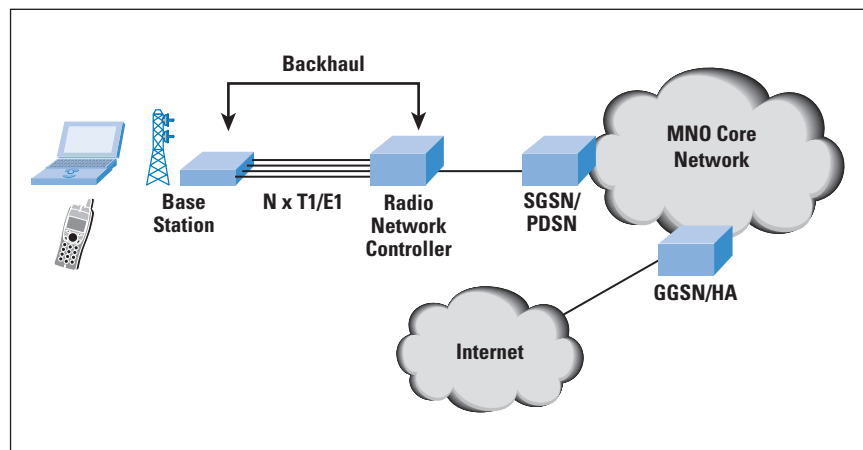
- Use of the same equipment used in wireline voice transmission
- Technical staffs' familiarity with TDM concepts and troubleshooting
- Ability to use existing *Operations, Administration, Maintenance, and Provisioning* (OAM&P) systems
- Ubiquity of the T1/E1 service

The initial work to offer data service on cellular systems naturally focused on adding data transmission to the existing voice infrastructure. Standards such as *Global System for Mobile Communications* (GSM) and *Interim Standard 95* (IS-95) took similar approaches in borrowing TDM time slots for data. The data services of the 1990s were very slow, even when compared to consumer modems of the time. Standards developed in the late 1990s and deployed in the early 2000s (*Enhanced Data rates for GSM Evolution* [EDGE] and CDMA2000) improved data transfer speeds.

TDM was clearly entrenched as a foundational technology for data communication in cellular networks going into the early 3G technology deployments (*Universal Mobile Telecommunications System* [UMTS] and *Evolution Data Optimized* [EV-DO]).

Figure 1 depicts the backhaul portion of the MNO network and how it fits into the broader architecture.

Figure 1: The Backhaul Network in the MNO Architecture



As data traffic usage for 3G networks grew, shortcomings of TDM backhaul began to materialize. The two prominent areas were bandwidth and cost. Cell sites with TDM access are typically equipped with multiple T1/E1s. With faster radio interfaces, the backhaul became the bottleneck in the network. Some smartphones became consumers of multi-megabyte data rates. User experiences were poor on some wireless networks as a result of a dearth of bandwidth in the backhaul segment. Continuing to increase the number of TDM lines or increase their capacity was not a viable option since the growth increments were too small and the operating expenses were too high.

The second limitation of TDM in 3G networks is cost. Although the cost of T1/E1s decreased considerably over the years, the costs piled up given the number of cell sites and number of T1/E1s per site. This figure became the highest contributor to the cost of the backhaul network. The MNOs that owned the last mile were at a distinct competitive advantage compared with the carriers who had to pay another party (often in a minimally competitive marketplace) for TDM access. For MNOs to continue their incredible traffic growth rates, a new access model was needed.

Carrier Ethernet Adoption

Ethernet quickly emerged as the most popular backhaul technology to replace TDM access infrastructure (other providers moved forward with microwave access with varying levels of success). The various iterations of Ethernet from 1970s to 2000s had trumped other LAN technologies in the market, and at the turn of the century gigabit Ethernet leveraged its success in the LAN to become popular in the WAN. The technology had several major advantages:

- *Large drop in cost per bit:* Ethernet would allow providers to drastically alter their access cost model by supplanting the aging and costly TDM infrastructure. With the price that consumers were willing to pay per month of data service staying relatively stagnant, this adjustment to the cost model was critical.
- *Ethernet can be carried over more underlying technologies:* *Synchronous Optical Networking/Synchronous Digital Hierarchy* (SONET/SDH), *Generic Framing Procedure* (GFP), *Dense Wavelength Division Multiplexing* (DWDM), and *Multiprotocol Label Switching* (MPLS) are a few examples. A key benefit Ethernet's ability to operate over these technologies was that many providers could consolidate their wireless access with their existing and speedier wireline access networks.
- *Ethernet interfaces ubiquitous and inexpensive:* Ethernet won the battle for LAN dominance. The technology was not restricted to traditional personal computers and servers—printers, phones, game consoles, *Digital Video Recorders* (DVRs), and home media center hubs are some examples of other equipment that often included Ethernet interfaces. This ubiquity in the business and consumer spaces results in a diverse supplier set and economies of scale for the vendors and suppliers.
- *Ease of bandwidth upgrade:* TDM circuits have an implementation time measured in months. This slow turn-around time for upgrades is a poor fit for an environment in which data usages is increasing at fast rates. Ethernet is much different. An increase in bandwidth to a network end-point will not require a change in equipment unless moving between the established tiers of 10, 100, 1000 Mb/s. Since the Ethernet service vendor likely uses a “policer” to keep customers within the purchased bandwidth level, a change in software configuration is usually all that is required to upgrade bandwidth. Another advantage is that bandwidth can be upgraded in granular increments. With the right back-end systems, an upgrade will take a matter of minutes. For companies looking to increase the velocity of service deployment, the ability to quickly move to high speeds is very favorable.

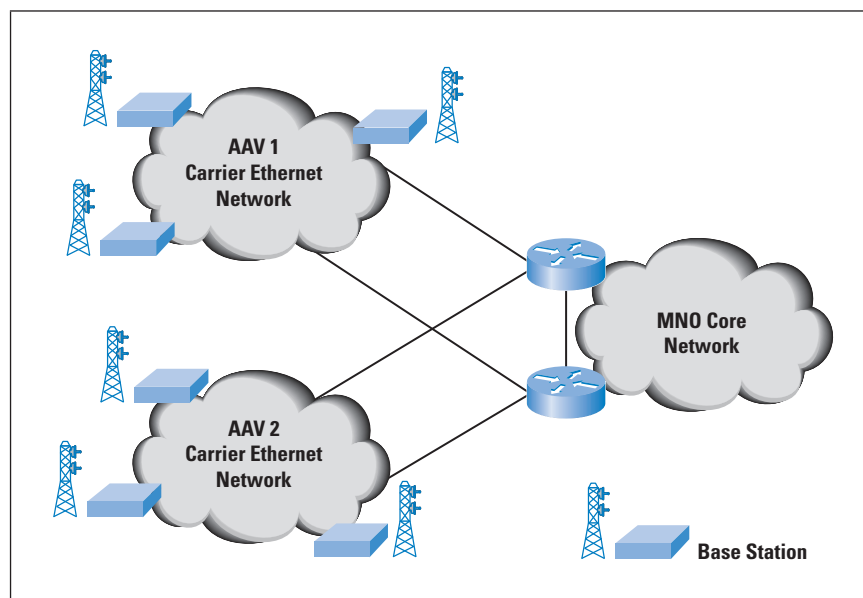
Established in 2001, the *Metro Ethernet Forum* (MEF) played a critical role in the acceptance of carrier Ethernet by wireless and wireline providers. The MEF is not a standards organization like the *Internet Engineering Task Force* (IETF). Instead, the MEF builds upon the work of standards bodies to establish common terminology, service requirements, and network interface requirements. The MEF created an architecture framework along with measurement and testing specifications. Although the MEF did not eliminate wireless providers' concerns about packet backhaul—particularly in the areas of jitter, delay, and packet delivery, the forum did increase the comfort level associated with metro Ethernet services. The MEF's E-LINE service definition established a connection-oriented path, a concept much more pleasing to traditional telcos than the perceived “anything goes” nature of packet switched networks. For more detail on the MEF's service definitions, see [0].

By the second half of the 2000s, many wireless providers were planning the deployment of Ethernet-based backhaul for new *High Speed Packet Access* (HSPA), *Worldwide Interoperability for Micro-wave Access* (WiMAX), and *Long-term Evolution* (LTE). In making this radical change, the providers often had to consider protecting existing revenue streams from voice and data (providers electing to move forward with greenfield deployments were at a luxury). Pseudowire technologies enabled the carriage of TDM traffic over IP/Ethernet networks, thus preserving investment in existing infrastructure.

Rather than build carrier Ethernet infrastructure, the MNOs that were not facilities-based (or had limited last mile footprints) purchased services from other parties, known as *Alternate Access Vendors* (AAV) in telco parlance. In the United States, the *Local Exchange Carriers* (LECs) and cable companies were well positioned for this business. MNOs often used multiple AAVs in a given market to cover the cell site footprint. Getting fiber to cell sites outside of major metropolitan areas was not always possible, which led some MNOs to use hybrid backhaul solutions that included microwave and TDM inverse muxing in addition to carrier Ethernet.

Figure 2 illustrates how MNOs rely on AAVs to cover their cell site footprint in a given market.

Figure 2: *Alternative Access Vendors*



The adoption of carrier Ethernet services by MNOs was not without challenges. Mobility gear such as *Radio Network Controllers* (RNC), base stations, and *Home Location Registers* (HLR) historically relied on T1/E1 interfaces for connection to the network. Telecom vendors had to implement Ethernet interfaces along with IP stacks. The providers had to completely revamp provisioning, service monitoring, performance monitoring, and service assurance systems and processes. Consider the following example.

For years, operations groups at telcos counted on near-immediate notification with an alarm indication signal in the *Time Division Multiple Access* (TDMA) frame. TDMA frames arrive every 125 μ sec (8,000 times a second). Packet-switched networks do not share the synchronous nature of TDM and do not have OAM fields in framing bits. The operators now had to rely on nascent specifications such as Y.1731 and 802.1ag for service monitoring.

Timing and synchronization—necessities in mobile networks—are gleaned from the physical layer in TDM networks. Asynchronous networks such as Ethernet/IP do not have an inherent mechanism for timing and synchronization. Keeping a single T1/E1 at the cell site is one method to ensure timing and synchronization in a carrier Ethernet scenario; however, the use of upper layer protocols is more appropriate, particularly for new builds that have no legacy TDM circuits. *Synchronous Ethernet* (SyncE), *Precision Time Protocol* (PTP, also known as IEEE 1588v2), and *Network Time Protocol version 4* (NTPv4) were deployed in backhaul networks to provide timing and synchronization. Note that SyncE transports timing information over the physical layer much like the TDM timing model, while PTP and NTP use IP for transport and are not dependent on an Ethernet physical layer.

The learning and flooding aspects of all Ethernet networks present inherent scaling challenges for very large networks. Spanning tree and its derivatives are commonly used to address these issues at low and medium scale. For larger networks that provide service to multiple customers, the service must scale in terms of its ability to offer service to multiple entities and in terms of the many switches required for an expansive footprint. Many protocols have arisen to solve one or both of these challenges. Examples are *Virtual Private LAN Service* (VPLS), *Multiprotocol Label Switching–Transport Profile* (MPLS-TP), and *Provider Backbone Bridging–Traffic Engineering* (PBB-TE). Being relatively new technologies, these can and do present challenges for operations groups. The breakages can occur in ways that are very difficult for the Carrier Ethernet provider and wireless provider to jointly troubleshoot.

The Next Step – IP Backhaul

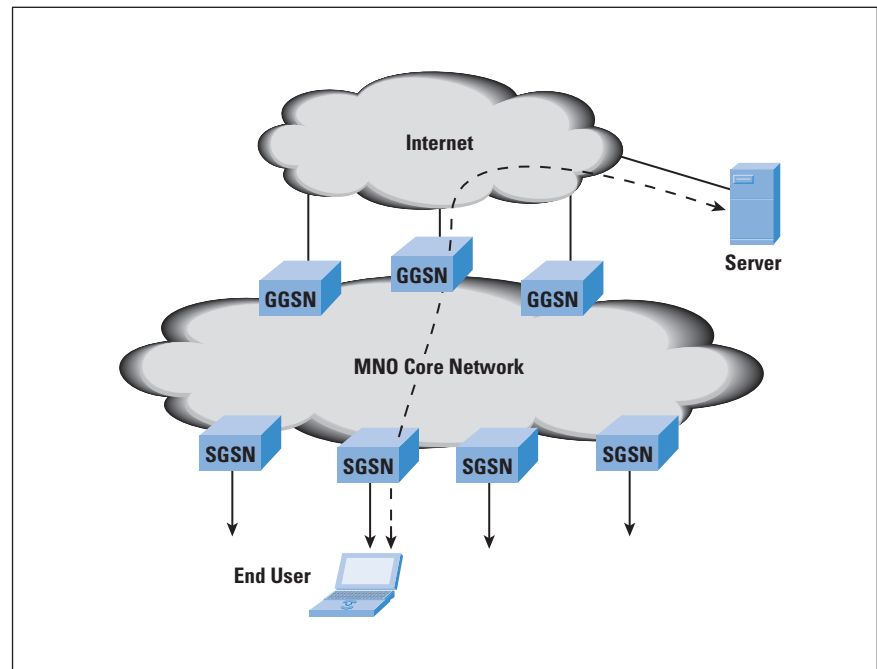
The phrase “all-IP” is frequently used to describe the most recent wireless technologies such as HSPA+, WiMAX, and LTE. This is applicable as the majority of network elements, including the handsets, are IP enabled. The existence of large-sized carrier Ethernet networks in the network architecture undermines the IP-centric argument. IP has superior scaling properties over Layer 2 networks. The footprint and number of nodes for carrier Ethernet networks continues to expand rapidly as the MNOs deploy 3G and 4G networks. The author sees evidence that protocols used to overcome Ethernet scalability issues will become increasingly complex and push MSOs and AAVs toward Layer 3-centric backhaul networks.

Before delving into the drivers of IP backhaul, let's examine a typical data traffic flow for today's wireless networks. We'll use the 3GPP's *GSM Packet Radio System* (GPRS) as this is the most common in world-wide deployments. Data flows are very centralized in this architecture. Macro-level mobility is controlled by two types of *GPRS Support Nodes* (GSN): *Gateway GPRS Support Nodes* (GGSN) and *Serving GPRS Support Nodes* (SGSN). GGSNs are typically deployed within the mobile core network at locations with Internet access. This is often at centralized mobile switching centers. SGSNs can be deployed closer to the network edge and multiple SGSNs can be served by a single GGSN.

The GGSN is the mobility anchor, much like the home agent in wireless networks that use Mobile IP. The SGSN is akin to the foreign agent in Mobile IP. GPRS network tunnel traffic between SGSN and GGSN using an IP-in-IP tunneling protocol called *Generic Tunneling Protocol* (GTP). Although GTP has several purposes in the GPRS core network, our focus will be on its tunneling of packets between SGSN and GGSN (called the *Gn* interface). The movement of the subscriber to a region served by another SGSN will trigger a macro-mobility event. A new GTP tunnel is formed using the original GGSN for session continuity [2].

Since all traffic from the *Mobile Subscriber* (MS) must traverse the GGSN as the mobility anchor, the traffic flow from the MS follows a very predictable path to a centralized location. Note that there is not a 1:1 relationship between SGSNs and GGSNs. As mentioned earlier, typical deployment of GGSNs is very centralized. Figure 3 depicts the flow.

Figure 3: Data flow in a GPRS Network



Although technologies like LTE are touted as flat IP networks, this only holds true from a *Radio Access Network* (RAN) perspective. What if a subscriber wants to communicate with another subscriber in the same building or local machine-to-machine traffic is highly sensitive to latency? The packets will be sent to the mobility anchor, perhaps hundreds of kilometers away. Routing decisions can be made in the RAN and core network; however, the decision is restricted since traffic must traverse the predefined tunnel endpoints.

Wireless networks will gradually decentralize and distribute mobility management. In 3G networks, some providers have been extending the core network closer to the subscriber as mobile gateways (GSNs and their equivalents in non-3GPP networks) become more cost-competitive. By deploying mobile gateways at what were previously aggregation *Points Of Presence* (POPs) and buying Internet connectivity at these locations, Internet-bound traffic exits the network quickly, consuming fewer resources for the provider. Other signs of this shift are evident in LTE and WiMAX. LTE's S1-flex interface allows the RAN to be connected to multiple core networks. The WiMAX reference model separates the *Network Access Provider* (NAP) and *Network Service Provider* (NSP). The NAP, which provides radio access functionality, can connect to multiple NSPs for Internet connectivity.

To fully realize the benefits of an IP-centric backhaul, steps must be taken to go beyond simply distributing mobility management. New solutions are needed to eliminate mobility anchoring via tunneling. Vendors, providers, and universities have already started to examine how to dispose of tunneling in the mobile environment [2].

The IP-centric backhaul network has many advantages over the carrier Ethernet networks that enable many of today's packet backhaul networks. Various advantages benefit the wireless providers, the IP backhaul provider, or both. These advantages are most prevalent when the MSOs have a highly distributed mobility management architecture.

- *Backhaul Offload:* Today's mobile elements at the cell tower have no ability to influence routing decisions; there is only one path to the core network. Adding egress points to the cell site or backhaul network reduces the distance and amount of traffic that must be backhauled. To accomplish the addition of egress points in a carrier Ethernet network, connection-oriented mechanisms such as Ethernet Virtual Circuits would require that the MSO and AAV modify multiple network elements' configurations. Offloading traffic with an IP network is substantially more simple and scalable. Offloading packets from the backhaul will represent a huge savings in access costs. The base station could be capable of hot potato routing traffic directly to an ISP instead of backhauling commodity Internet traffic to the MSO, where the costs of equipment, power, and software licenses quickly accumulate.

- *Multicast*: The reliance on tunneling as described earlier in this piece severely restricts the usefulness of multicast in current wireless networks. Distributing the mobility elements controlling the tunneling closer to the subscriber will mitigate these effects as would the elimination of mobility anchoring via tunneling techniques. The implementation of a true flat IP network would extend multicast capability into the RAN and position both MNOs and IP backhaul providers to realize the efficiency gains of multicast.
- *Localized Content and Peering*: With localized egress points, local content could be reached directly rather than traversing the core network. This would position wireless providers to peer with other providers at the local or regional level, a benefit that would be substantial for wireless providers operating in countries with non-meshy Internet infrastructure and expensive wide-area communications lines. In addition, caches could be implemented much closer to the subscriber to improve the user experience for video and other content types.
- *Machine-to-Machine (M2M) and Peer-to-Peer (PtP)*: When the communication is device to device in close geographic proximity, the traversal of the core network only adds latency, complexity, and cost. A distributed mobility management architecture and IP backhaul network engender an optimized path for M2M and PtP. The mobility anchor point could be placed at the cell tower or local aggregation point, providing a much improved communication path for subscribers and machines connected to the wireless network.
- *Uptime and Reliability*: Wireless providers have experienced challenges with carrier Ethernet service. Some of these problems can be chalked up to the relative newness of using carrier Ethernet for cell site backhaul. One has to wonder though, what experience exists in the industry for maintaining giant Layer 2 networks? The number of mobile devices will expand exponentially, triggering the deployment of thousands of new cell sites, microcells, and picocells. The author is less than confident that any underlying technology that enables carrier Ethernet will scale to the necessary degree while maintaining the uptime and reliability that users expect from their data service.

For large IP networks, the industry has over fifteen years' experience in designing, engineering, and operating IP networking carrying traffic at staggering capacities. The staff expertise, software maturity, and systems support exists today to maintain sizable IP networks. There are established best practices for Tier 1 ISPs that help ensure long uptime, speedy convergence upon failure, and sound network design.

Delivering an IP Backhaul Service

IP backhaul offerings could be delivered in a variety of ways. The simplest design for IP backhaul providers would be a shared IP transport network that commingles traffic between customers.

The wireless providers could then use protocols such as *Layer 2 Tunneling Protocol version 3* (L2TPv3) to build an MPLS/VPN-like overlay to provide logical separation and address overlap prevention. The preferred approach for MNOs would likely be a Layer 3 VPN service from the AAV, thereby offloading much of the routing complexity from the MNO.

An IP backhaul service must be capable of routing IPv6 packets, as the useful lifetime of an IPv4-only service is limited. MNOs cannot obtain new IPv4 addresses to number the base stations, and using RFC 1918 space is not a scalable approach. Using IPv6-only to address mobility equipment at cell sites (and equivalent radio interfaces) is the preferred method for overcoming the scarcity of IPv4 addresses.

The shift from carrier Ethernet to IP backhaul should not be a monumental one for many carrier Ethernet providers. The heavy lifting of installing fiber and deploying a packet switched infrastructure has already been accomplished. In addition, carriers that implement carrier Ethernet with protocols like VPLS already have an infrastructure that is ready for IP. The most challenging aspect of the transition will be the work needed to prepare OAM&P systems for an IP service. Of course, this may vary based on carrier Ethernet implementation and systems.

Conclusion

Carrier Ethernet service for cell site backhaul is a vast scale and cost improvement over TDM backhaul and has been extremely successful. OSI Layer 3 IP networks have superior scaling properties that will replace Layer 2 backhaul networks of today. Advances in wireless networking systems, the proliferation of new devices, and the development of new mobility services will be best served with a truly IP-centric backhaul network.

References

- [0] Santitoro, Ralph, “Metro Ethernet Services—A Technical Overview,” 2003, <http://metroethernetforum.org/metro-ethernet-services.pdf>
- [1] M. Grayson, K. Shatzkamer, and S. Wainner, *IP Design for Mobile Networks*, Cisco Press, 2009.
- [2] *Distributed Mobility Management in Future Wireless Networks* (DiMoWiNe), <http://conference.researchbib.com/print.php?category=event&id=10232&uid=6>

JEFF LOUGHRIDGE is the principal consultant and owner of Brooks Consulting LLC, a firm that specializes in Tier 1 ISP best practices and the design, engineering, and operations of large-scale wireline and wireless IP/MPLS networks. Prior to founding Brooks Consulting, Jeff spent over ten years supporting Sprint’s global IP network in both technical and managerial capacities. He earned a bachelor’s degree in computer science from Duke University and an MBA from the University of Phoenix—Northern Virginia campus.

E-mail: jeffl@brooksconsulting-llc.com

Global INET 2012

To help mark its 20-year-anniversary, the *Internet Society* (ISOC) is hosting a global forum that will bring together visionaries and thought leaders from around the world to focus on issues that will impact the future of the Internet.

The *Global INET 2012*, which is scheduled to take place in Geneva, Switzerland from April 22–24, will feature high-powered speakers, thought-provoking panel discussions, and interactive workshops to develop a vision for the explosive growth of the Internet over the next 20 years.

Thought leaders from across the Internet community will collaborate on topics critical to the global Internet's future, including privacy, net neutrality, IPv6, security, digital content and innovation, and human rights and freedom of expression.

Since its beginnings in 1992, ISOC has been dedicated to helping keep the Internet open, accessible, and defined by users—regardless of where they live, what they do, their abilities, or who they are.

Registration for Global INET 2012 is scheduled to begin in October 2011.

For more information:

- [1] Barry M. Leiner, Vinton G. Cerf, David D. Clark, Robert E. Kahn, Leonard Kleinrock, Daniel C. Lynch, Jon Postel, Larry G. Roberts, Stephen Wolff, "A Brief History of the Internet," December 2003, also published in ACM's *Computer Communication Review*, Volume 39, Number 5, October 2009.
<http://www.isoc.org/internet/history/brief.shtml>
<http://www.sigcomm.org/ccr/papers/2009/October/1629607.1629613>
- [2] "The Internet Society's Principles and Goals,"
<http://www.isoc.org/isoc/mission/principles/>
- [3] <http://www.isoc.org/isoc/conferences/inet/12/gva.shtml>

IPv6 Week

IPv6 Week will be a coordinated test of the new Internet Protocol, held February 6–12, 2012. Websites, content providers, Internet Services Providers, Network Service Providers, as well as end users are invited to participate. This is a Brazilian initiative, but anyone can participate.

For more information visit: <http://www.ipv6.br/IPV6/WeekIPv6>

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ contains standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2011 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

December 2011

Volume 14, Number 4

A Quarterly Technical Publication for
Internet and Intranet Professionals

In This Issue

From the Editor	1
Port Control Protocol	2
Challenges to DNS Scaling	9
Networking @ Home.....	15
IETF Tools	21
Fragments	25
Call for Papers.....	31

FROM THE EDITOR

Depletion of the IPv4 address space and the transition to IPv6 has been a “hot topic” for several years. In 2011, interest in this topic grew considerably when the *Asia Pacific Network Information Centre* (APNIC) became the first *Regional Internet Registry* (RIR) to start allocating addresses from its final /8 IPv4 address pool. Although depletion dates are difficult to predict accurately, there is no question that the day will come when it will no longer be possible to obtain IPv4 space from the RIRs. News stories about IP addresses being sold for considerable sums of money are becoming more common.

Numerous organizations have been working diligently to promote, test, and deploy IPv6 through efforts such as the *World IPv6 Day*, while the *Internet Engineering Task Force* (IETF) continues to develop solutions to aid in the transition. One such effort, the *Port Control Protocol* (PCP), is described in our first article by Dan Wing.

The *Internet Corporation for Assigned Names and Numbers* (ICANN) will soon begin accepting applications for new *Top-Level Domains* (TLDs). It is not yet known how many new TLDs will eventually be deployed, but the plans have prompted several studies focused on the resiliency and scalability of the *Domain Name System* (DNS). Bill Manning discusses some of the technical challenges associated with a vastly expanded TLD space.

The *IETF Homenet Working Group* “...focuses on the evolving networking technology within and among relatively small ‘residential home’ networks. For example, an obvious trend in home networking is the proliferation of networking technology in an increasingly broad range and number of devices. This evolution in scale and diversity sets some requirements on IETF protocols.” Geoff Huston gives an overview of some of the challenges facing this Working Group.

The product of the IETF is a set of documents, mainly protocol specifications and related material. These documents start life as *Internet Drafts* and proceed through a series of iterative refinements toward eventual publication as *Request For Comments* (RFCs). Over time, several *tools* have been developed to aid in the document development process, and they are now organized at the IETF Tools webpage. We asked Robert Sparks to give us an overview of some of the most important tools and the process involved in their development.

—Ole J. Jacobsen, Editor and Publisher

ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

Port Control Protocol

by Dan Wing, Cisco Systems

After the transition to *Internet Protocol Version 6* (IPv6), hosts will often be behind IPv6 firewalls. But before the transition, mobile wireless devices will want to reduce their keepalive messages, and hosts of all sorts will share IPv4 addresses using a variety of address-sharing technologies. To meet these needs, the IETF formed the *Port Control Protocol Working Group* in August 2010 to define a new protocol for hosts to communicate with such devices. The initial output of this Working Group is the *Port Control Protocol* (PCP)^[1]. Interoperability between two independently developed implementations of PCP was demonstrated at the IETF meeting in July 2011, highlighting the importance of this protocol to the industry. After it becomes a standard, PCP is expected to be deployed in various operating systems, IPv6 home gateways, IPv4 home gateways (*Network Address Translators* [NATs]), mobile third- and fourth-generation (3G and 4G, respectively) gateways (*Gateway GPRS Support Nodes* [GGSNs]), and *Carrier-Grade NATs* (CGNs).

Introduction to PCP

PCP performs two major functions: It allows packets to be received from the Internet to a host (such as to operate a server), and allows a host to reduce keepalive traffic of connections to a server. PCP can be extended in two ways: with new *OpCodes* or with new *Options*. The base PCP specification defines two OpCodes: MAP and PEER, and defines several Options that can be carried with those OpCodes.

To operate a server, packets are sent from a host on the Internet to a server. The IP model expects devices to be connected to a network and be able to exchange packets with each other. However, few deployed networks actually permit hosts to receive packets from the Internet because of business needs (for example, to protect wireless spectrum from malicious or accidental packets originated on the Internet) or because of technology restrictions (for example, IPv4 address-sharing devices such as *Network Address and Port Translators* [NAPT]). To operate a server, a host uses the MAP OpCode.

To reduce keepalives, a host needs to send traffic before a middlebox will destroy an idle connection. Many middleboxes, such as firewalls or NATs, maintain state and will destroy mappings if the connection has been idle. Today, in order to prevent destruction of mappings, hosts send keepalive traffic to keep those mappings alive. The keepalive traffic has several disadvantages, including reduction of battery lifetime, network chatter, and server scalability (servers have to discard the keepalive traffic). PCP allows a host to determine how aggressively a middlebox will destroy an idle connection, allowing the host to reduce its keepalive traffic with the PEER OpCode.

PCP is encoded in binary and carried over the *User Datagram Protocol* (UDP), which eases implementation on clients and servers. The client is responsible for retransmitting messages, and all messages are idempotent. The PCP client can be part of the operating system (much like a *Dynamic Host Configuration Protocol* [DHCP] client or a *Universal Plug and Play* [UPnP] *Internet Gateway Device Protocol* [IGD] client) or the PCP client can be coded entirely in an application (much like any other application-level protocol such as the *Network Time Protocol* [NTP]). A major feature of PCP is its flexibility and simple messaging, so it can be implemented easily in a variety of systems and at high scale.

Security

When installing an IPv4 NAT on a residential network, the NAT has a side effect: it prevents unsolicited incoming traffic from reaching hosts inside the home. Traffic that originates inside the home can traverse the NAT toward the Internet. This function is expected by many users to such a degree that when IPv6-capable routers were first installed on residential networks, users complained that their IPv6 hosts were seeing traffic from the Internet. This visibility meant that IPv6 printers, webcams, and other hosts had to be protected from malicious traffic from the Internet. Based on this experience, IPv6 *Customer Premises Equipment* (CPE) routers intended for installation in the residential market filter most unsolicited incoming traffic by default^[3]. Thus, IPv6 CPE routers provide filtering similar to what users experience today with IPv4 NAT devices.

With both IPv4 NAT and RFC 6092 IPv6 routers, outgoing traffic from a host creates a mapping that then allows bidirectional traffic to a specific (*Transmission Control Protocol* [TCP] or UDP) port on the internal host, meaning when a host sends a TCP SYN, a SYN ACK can be returned to the host. Neither IPv4 NAT devices nor RFC 6092 IPv6 routers have to do any additional filtering of that mapping, and after that mapping is created will allow traffic from any host on the Internet to reach the internal host—not just traffic from that particular host. This lack of filtering is necessary for certain applications to function.

PCP was built with a security model similar to that deployed on home networks. With PCP, a host can send a PCP packet requesting a mapping so that any host on the Internet can now initiate communications with the internal host. Similarly, without PCP, a host could send a TCP SYN from a specific port (for example, port 80), thereby creating a mapping nearly identical to a PCP mapping. As with sending a TCP SYN, PCP allows a host to open mappings only for itself, unless the network administrator has taken the extra step to enable the PCP THIRD_PARTY option.

You may wish to have additional restrictions for some networks. PCP is extensible to support authorization, and there is ongoing work to support authentication and authorization within PCP^[8].

PCP is extensible and there are already several proposed extensions to the protocol, including a way to control which IP address pool is assigned to a mapping^[5], bulk port allocation to optimize acquiring a large set of ports^[6], and rapid recovery after NAT failure or network renumbering^[7].

PCP Scenarios

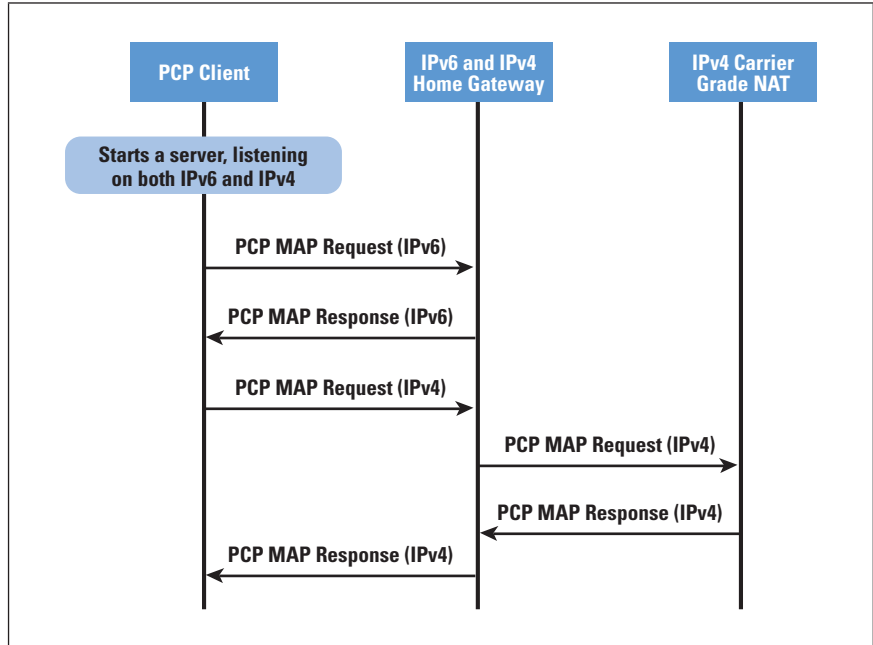
PCP works in all scenarios with IPv4 address sharing (using an IPv4 NAT or using other techniques), an IPv4 or IPv6 firewall, and NATs that translate from IPv6 to IPv4, IPv4 to IPv6, or IPv6 to IPv6. When working with nested NAT, such as a NAT in the home and a NAT operated by the *Internet Service Provider* (ISP), PCP can create the NAT mappings in both devices. When working with IPv6, PCP can create mappings in an IPv6 CPE router. In some networks we expect to see IPv6-only devices that IPv4 clients may need to access. For those devices to work, an IPv6/IPv4 translator (NAT64)^[10, 11] can translate between IPv6 and IPv4. PCP can work with an IPv6/IPv4 translator as well. In other scenarios IPv6/IPv6 translation may be necessary, and although translating IPv6 to IPv6 is far from desirable, PCP can also support IPv6/IPv6 (NPTv6)^[12].

A server, such as a one running on a sensor (for example, thermometer or electric meter), can use PCP to determine its publicly routable IPv4 or IPv6 address and port, and then populate a *Rendezvous* server with that IP address and port. For example, an IPv6-only thermostat might want to be accessible over IPv6 and IPv4, so it can be accessed by both the power company (to push new electricity rate information to the thermostat) and the homeowner (who might have IPv4 access only at work). The thermostat can use PCP to create a TCP mapping in the IPv6 CPE router (necessary because the IPv6 CPE router will, by default, filter unsolicited incoming IPv6 packets) and use PCP to create a TCP mapping in a NAT64 (necessary so the homeowner can access the thermostat). The IPv6 address and its TCP port, and the IPv4 address and its TCP port, can be published to the *Domain Name System* (DNS) (using DNS Server [SRV] records) or published to some other Rendezvous server. Then the power company or the homeowner can use the DNS (or the other Rendezvous server) to communicate directly with the thermostat.

Because PCP can inform the PCP client of address changes, network renumbering can be communicated immediately to hosts—something that cannot be done with most other NAT or firewall control mechanisms. Therefore, devices running on nomadic networks, such as in a connected vehicle, that use PCP will immediately learn when they have connected to a new network. This knowledge can allow them to update information in the DNS or in some other Rendezvous server so they remain accessible from the Internet.

PCP is expected to be implemented in home gateways and Carrier-Grade NATs, which provide value for both IPv6 (to operate a server and learn keepalive timeouts) and IPv4. Figure 1 shows how a dual-stack host would use PCP to operate an IPv6 or IPv4 server.

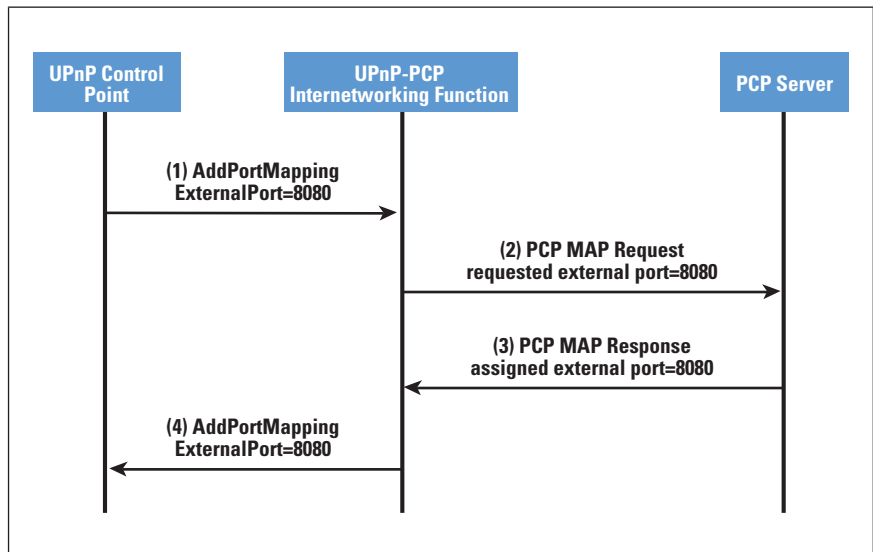
Figure 1: PCP Mapping IPv6 and IPv4



PCP Interworking with UPnP IGD

UPnP IGD Version 1 is widely available on residential-class NAT devices and host operating systems (Windows and OS X). However, because of security concerns it is often disabled by vendors, ISPs, or end users. UPnP IGD itself only works with a single layer of NAT, but it is possible to interwork between UPnP IGD and PCP^[4]. To do this interworking, a home gateway (NAT) processes UPnP IGD messages on its LAN interface and translates those messages to PCP messages on its WAN interface, as depicted in Figure 2.

Figure 2: UPnP-to-PCP Interworking, Showing AddPortMapping Success



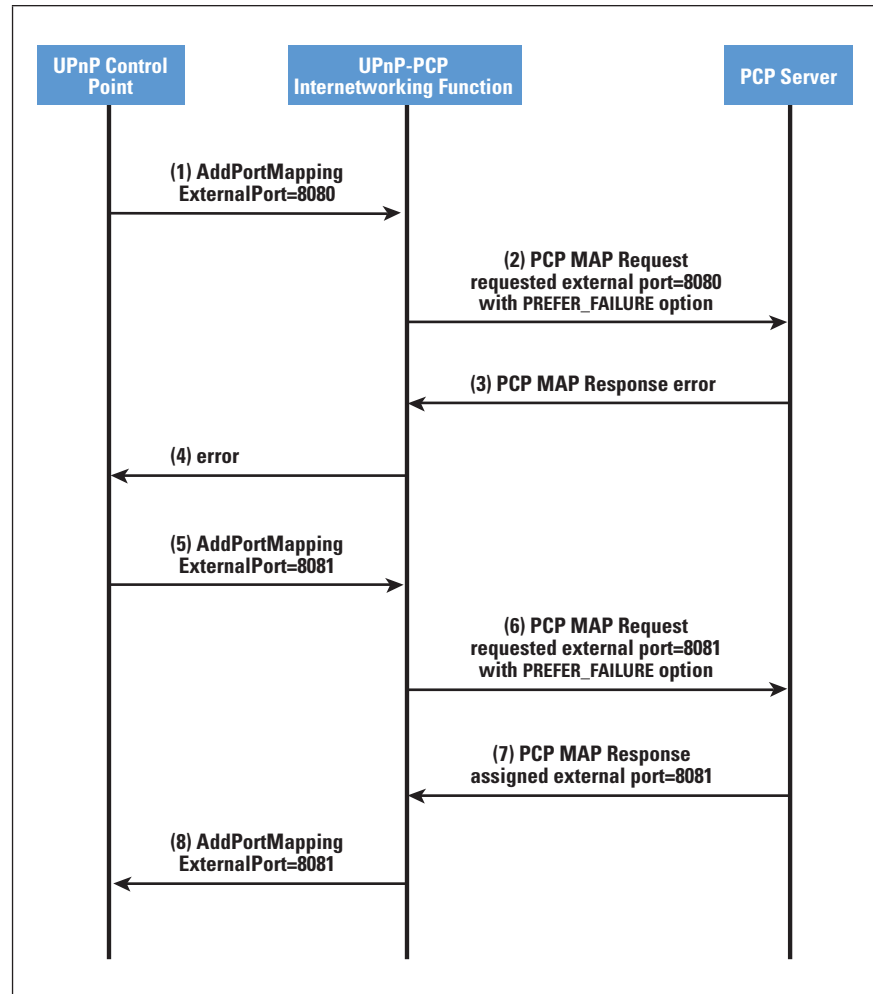
One difficulty with UPnP IGD is its *AddPortMapping* action, which maps a specific port on the home gateway. If that requested port is already mapped to another host, that port cannot be mapped to a new host (because it is already mapped to a different host). This problem exists today with UPnP IGD if two hosts in a home need the same port (for example, TCP port 80) because only one of them can map the port. In a CGN environment, where many subscribers share one IPv4 address, it is almost guaranteed that another subscriber has already mapped a “good” port (for example, 80 for HTTP, 8080 for HTTP, 5001 for Slingbox, 5060 for *Session Initiation Protocol* [SIP], etc.). Today, when a UPnP IGD port mapping is refused, the application may overwrite the first host’s mapping (causing significant problems), “hunt” for an available port, or simply give up and display an error to the user. The “hunting” is often sequential (trying the next-higher port number) but is sometimes random, and is done by the application itself, the operating system UPnP framework, or both.

UPnP IGD Version 2^[2] introduced the *AddAnyPortmapping* action, which avoids the need to “hunt” for an available port and allows the NAT to assign an available port. But UPnP IGD Version 2 is not yet widely available in home gateways, operating systems, or applications. Until IPv6 is ubiquitously available, applications (and users) will need to practice better port agility than has been practiced in the past, because “good” ports will simply not be available when IPv4 addresses are shared.

To ease the interworking with the UPnP IGD *AddPortMapping* action, the base PCP specification includes a *PREFER_FAILURE* option, which avoids creating a mapping if the requested port is unavailable. A message flow of this behavior is shown in Figure 3.

In a *Dual-Stack Lite*^[9] deployment, the home gateway is typically operated without a NAT function. In that configuration, the home gateway is expected to interwork between UPnP IGD (within the home) and PCP (toward the service provider’s CGN). The PCP packets sent by the home gateway will have the source IP address of the home gateway, rather than the IP address of the host that initiated the UPnP IGD action. To accommodate that situation, the home gateway populates the *THIRD_PARTY* option with the IP address of the internal host needing the mapping. The *THIRD_PARTY* option is useful in other scenarios as well, including interworking with other protocols (such as the *NAT Port-Mapping Protocol* [NAT-PMP]^[13]) to PCP, using PCP to create mappings for a device that does not support PCP (for example, an IP-enabled webcam), or using it as the protocol between a web portal operated by the ISP and its CGN.

Figure 3: UPnP-to-PCP Interworking,
Showing AddPortMapping Failure



Conclusion

PCP provides functions necessary for IPv6 hosts on home networks; it is a simple, scalable protocol that supports simple firewalling of IPv6 and IPv4 hosts, and to accommodate the transition to IPv6 also supports every conceived IPv4/IPv6 translation mechanism.

References

- [1] Dan Wing, ed., Stuart Cheshire, Mohamed Boucadair, Reinaldo Penno, and Paul Selkirk, "Port Control Protocol (PCP)," Internet Draft, work in progress, July 2011, **draft-ietf-pcp-base**
- [2] "UPnP Gateway committee: IGD:2 improvements over IGD:1," March 2009, <http://www.upnp.org/resources/documents/UPnPIGD2vsIGD1d10032009.pdf>
- [3] James Woodyatt, ed., "Recommended Simple Security Capabilities in Customer Premises Equipment for Providing Residential IPv6 Internet Service," RFC 6092, January 2011.

- [4] Mohamed Boucadair, Reinaldo Penno, Dan Wing, and Francis Dupont, “Universal Plug and Play (UPnP) Internet Gateway Device (IGD)-Port Control Protocol (PCP) Interworking Function,” Internet Draft, work in progress, February 2011, **draft-bpw-pcp-upnp-igd-interworking**
- [5] Reinaldo Penno, “PCP Support for Multi-Zone Environments,” Internet Draft, work in progress, June 2011, **draft-penno-pcp-zones**
- [6] Cathy Zhou, Tina Tsou, Xiaohong Deng, Mohamed Boucadair, and Qiong Sun, “Using PCP To Coordinate Between the CGN and Home Gateway Via Port Allocation,” Internet Draft, work in progress, July 2011, **draft-tsou-pcp-natcoord**
- [7] Stuart Cheshire, “PCP Rapid Recovery,” Internet Draft, work in progress, June 2011, **draft-cheshire-pcp-recovery**
- [8] Margaret Wasserman, Sam Hartman, and Dacheng Zhang, “Port Control Protocol (PCP) Authentication Mechanism,” Internet Draft, work in progress, October 2011, **draft-wasserman-pcp-authentication**
- [9] Alain Durand, Ralph Droms, James Woodyatt, and Yiu L. Lee, “Dual-Stack Lite Broadband Deployments Following IPv4 Exhaustion,” RFC 6333, August 2011.
- [10] Congxiao Bao, Christian Huitema, Marcelo Bagnulo, Mohamed Boucadair, and Xing Li, “IPv6 Addressing of IPv4/IPv6 Translators,” RFC 6052, October 2010.
- [11] Marcelo Bagnulo, Philip Matthews, and Iljitsch van Beijnum, “Stateful NAT64: Network Address and Protocol Translation from IPv6 Clients to IPv4 Servers,” RFC 6146, April 2011.
- [12] Margaret Wasserman and Fred Baker, “IPv6-to-IPv6 Network Prefix Translation,” RFC 6296, June 2011.
- [13] Stuart Cheshire, Marc Krochmal, and Kiren Sekar, “NAT Port Mapping Protocol (NAT-PMP),” Internet Draft, (expired), April 2008, **draft-cheshire-nat-pmp-03.txt**

DAN WING is the editor of the Port Control Protocol base specification and co-author of the PCP-UPnP interworking function specification. Dan has co-chaired the IETF’s BEHAVE Working Group since 2006. He is a Distinguished Engineer at Cisco Systems, where he works on IPv6 transition technologies.
E-mail: dwing@cisco.com

Infrastructure Challenges to DNS Scaling

by Bill Manning

This article looks a few steps beyond the *Root Scaling Study* report from 2009.^[1] In 2009, the *Internet Corporation for Assigned Names and Numbers* (ICANN) board commissioned a report to evaluate the effect of scaling the root zone from its current size to an undefined but larger root zone. Attributes considered were *Domain Name System Security Extensions* (DNSSEC), *Internet Protocol Version 6* (IPv6), *Internationalized Domain Names* (IDNs), and a larger number of entries in the zone. The report itself focused on the editorial processes and presentation of the finished root zone to the greater Internet. The report concluded that with prudence and with the addition of some “watch & warn” systems in place, the root zone could accommodate adding IPv6, DNSSEC, and IDNs along with other new *Top-Level Domain* (TLD) entries in a controlled manner. What the report did not consider was the effects of the deployed Internet infrastructure on the ability to get this new information into the rest of the *Domain Name System* (DNS) infrastructures of the Internet. Early experimental evidence^[7, 8] suggests that the current state of infrastructure deployment will create problems for the deployment of these attributes.

Until recently the root zone of the DNS has enjoyed two important stabilizing properties:

- It is relatively small—currently the root zone holds delegation information for 280 generic, country-code, and special-purpose TLDs, and the size of the root zone file is roughly 80,000 bytes.
- It changes slowly—on average, the root zone absorbs less than one change per TLD per year, and the changes tend to be minor.

The root system has therefore evolved in an environment in which information about a small number of familiar TLDs remains stable for long periods of time. However, the type, amount, and volatility of the information that is contained in the root zone are expected to change as a result of the following four recent or pending policy decisions:

- Support for DNSSEC, or “signing the root”
- The addition of IDN TLDs
- Support for the additional larger addresses associated with IPv6
- The addition of new TLDs

These changes are placed in a backdrop of an infrastructure that is fundamentally changing, removing a third attribute of a stable DNS that was the presumption of a common transport protocol with well-defined constraints.

Core Design Principles

The DNS was designed so that queries and responses would have the greatest chance of survival and broadest reachability by using an IPv4 default *User Datagram Protocol* (UDP) packet size of 512 bytes for the initial bootstrapping. Larger packet sizes are supported and the *Transmission Control Protocol* (TCP) was defined and used as an alternate transport protocol—but expected to be infrequently used.

With these core principles intact, the DNS was able to successfully evolve into a highly decentralized dynamic system. The geographic and organizational decentralization of the root system arises from a deliberate design decision in favor of diversity and minimal fate-sharing coordination, which confers substantial stability and robustness benefits on the global Internet.

Simple quantitative extrapolation from a baseline model of the current DNS does not predict realistic future states of the system beyond the very short term, because:

- Each part of the system adapts in different ways to changes in the quantity, type, and update frequency of information, while also responding to changes in the rest of the Internet.
- These adaptations are not—and cannot be—effectively coordinated.
- For some, if not all, of the actors, nonquantifiable considerations dominate their individual adaptation behavior (both strategically, in a planning context, and tactically, in an operations context).

The risks associated with adding DNSSEC and IPv6 addresses to the DNS simultaneously change the basic assumption for DNS Query/Response reachability. Signing DNS data would, by itself, immediately increase the size of any zone by roughly a factor of 4 and increase the size of the response message^[2]. The consequences of the second of these effects could be absorbed by replanning in order to recover lost headroom by adding bandwidth. Adding IPv6 addresses would in addition increase the size of any response. However, simply adding additional bandwidth may be insufficient when there are middleboxes, application layer gateways, or divergent transport options between the query path and the response path.

In these cases more information has to be carried in the packets that are returned in response to a query, meaning that the required amount of network bandwidth needed to support the operations of the server increases. As the DNS messages get bigger, they will no longer fit in single 512-byte packets forwarded by the UDP transport mechanism of the Internet. This situation will lead to clients being forced to resend their queries using UDP “jumbograms” or the TCP transport mechanism—a mechanism that has much more overhead and requires the end nodes to maintain much more state information. It also has much more overhead in terms of “extra packets” sent just to keep things on track. The benefit is, of course, that it can carry much larger pieces of information.

Moving the root system from its default UDP behavior to UDP “jumbograms” or TCP will not only have the undesirable effects mentioned previously, it will also affect the current trend of deploying servers using IP *anycast*^[10]. Anycast works well with single packet transactions (such as UDP), but is much less well suited to handle TCP packet streams. If TCP transactions become more prevalent, the anycast architecture may require changes.

The point of view from the client side is worth mentioning. In certain client configurations, where firewalls are incorrectly configured^[3], the following scenario can occur:

A resolver inside the misconfigured firewall receives a DNS request that it cannot satisfy locally. The query is sent to the root servers, usually over UDP, and a root server responds to this query with a referral, also over UDP. Today, this response fits nicely in 512 bytes. It is also true that for the past 6 years, the *Internet Systems Consortium* (ISC) has been anticipating DNSSEC and has shipped resolver code that, by default, requests DNSSEC data. After the root is signed, the response no longer fits into a 512-byte message. Estimates from the *National Institute of Standards and Technology* (NIST), using standard key lengths, indicate that DNSSEC will push the response to at least 2048 bytes or larger. This larger response will not be able to get past a misconfigured firewall that restricts DNS packets to 512 bytes, not recognizing the more modern extensions to the protocol that allow for bigger packets.

Upon not receiving the answer, the resolver on the inside will then retry the query, setting the buffer size to 512 bytes. The root will resend the response using smaller packets, but because it does not fit in a 512-byte packet, will fragment the response into a series of 512-byte replies, and the root server will set the “fragmented” and “truncated” flags in the packets, indicating to the resolver that the answer was fragmented and truncated, and encouraging the resolver to retry the query once more using TCP transport. The resolver will do so, and the root server will respond using TCP, but the misconfigured firewall also will reject DNS over TCP, because this transport has not been considered a normal or widely used transport for DNS queries.

In this worst case, a node will be unable to get DNS resolution after the root zone is signed, and the DNS traffic will triple, including one round in which TCP state must be maintained between the server and the resolver. There are of course ways around this problem, the most apparent ones being to configure the firewall correctly, or to configure the resolver to not ask for DNSSEC records.

Effect of IPv6 on Priming Queries

The basic DNS protocol specifies that clients, resolvers, and servers be capable of handling message sizes of at least 512 bytes. They may support larger message sizes, but are not required to do so.

The 512-byte “minimal maximum” was the original reason for having only nine root servers. In 1996 Bill Manning, Mark Kosters, and Paul Vixie presented a plan to Jon Postel to change the naming of the root name servers to take advantage of DNS label compression and allow the creation of four more authoritative name servers for the root zone. The outcome was the root name server convention as it stands today.

The use of 13 “letters” left a few unused bytes in the priming response, which were left there to allow for changes—which soon arrived. With the advent of IPv6 addressing for the root servers, it was no longer possible to include both an IPv4 “A” record and an IPv6 “AAAA” record for every root server in the priming response without truncation; AAAA records for only two servers could be included without exceeding the 512-byte limit. Fortunately the root system was able to rely on the practical circumstance that any node asking for IPv6 address information also supported *Extension Mechanisms for DNS* (EDNS0)^[4].

DNSSEC also increases the size of the priming response, particularly because there are now more records in the Resource Record set and those records are larger. In [5] the authors make the following observation: “The resolver MAY choose to use DNSSEC OK^[6], in which case it MUST announce and handle a message size of at least 1220 octets.”

EDNS and MTU Considerations

The changes described will also affect other parts of the Internet, including (for example) end-system applications such as web browsers; intermediary “middleboxes” that perform traffic shaping, firewall, and caching functions; and *Internet Service Providers* (ISPs) that “manage” the DNS services provided to customers.

Although modern DNS server software defaults to using EDNS0, current measurement^[7] collected from several of the RFC 1918^[11] servers suggests that EDNS0 usage has not yet reached generally accepted levels of usefulness. Over the 12-month study, the ratio of EDNS0 queries received at these nodes remained at roughly 65 percent of the total queries received, with about 33 percent being non-EDNS queries. In the “other” camp are queries that set EDNS0 but then restrict packet sizes to 512 bytes. These queries cannot use the larger, negotiable *Maximum Transmission Unit* (MTU) sizes for larger UDP responses and therefore must use TCP to support larger responses. Some evidence suggests that with signed data, there is a pattern of retransmission of queries when responses larger than 512 bytes are generated and blocked. Such retransmissions can take as long as 7 seconds before timing out.

Lack of EDNS0 support in DNS caches suggests that many parts of the Internet will be constrained to using the traditional UDP sizes or will fall back to using TCP. Even where EDNS0 is indicated as being available, there are increased difficulties in knowing or negotiating a consistent *Path Maximum Transmission Unit* (Path MTU)^[8].

The data supports an argument that the expectation of a useful UDP “jumbogram” or enough resources to manage hundreds of thousands or millions of TCP connections is unfounded because of historical expectations on “normal” DNS packet profiles. Clean, clear Internet paths that will allow larger packet sizes are rare, particularly when crossing the Internet. Locally, it is much more likely that larger packet sizes will be found and supported, raising the question for wide-scale deployment of IPv6 or DNSSEC because both attributes require larger packet sizes regardless of transport. If neither larger UDP packets nor TCP will be viable, what other choices are there?

Recent work inside the *Internet Engineering Task Force* (IETF) is exploring the use of the *Hypertext Transfer Protocol* (HTTP) as an alternative transport protocol for DNS messages.^[9] It might be possible to augment the deployed DNS base to understand the addition of a third transport protocol.

The augmentation of the DNS protocol to support multiple transport protocols will require additional logic on the part of the servers to keep track of which transport a query was received on and select that transport when sending back the response. It will also require more complex logic to determine failover selection from one transport to another.

With the efforts going into making the infrastructure of the Internet IPv6-capable, it is possible that the underlying MTU problems may be corrected faster than adoption of a new transport protocol for the DNS. Certainly MTU problems have been considered for many years and for slightly different reasons^[8] principally related to faster signaling rates and changes in the types of data being moved through the Internet. Regardless, this transition will take considerably more time than a simple DNS code refresh. Full support for larger packet sizes in the DNS will require changes in the equipment and code that comprise the baseline Internet infrastructure—and such changes may take decades.

References

- [1] Jaap Akkerhuis, Lyman Chapin, Patrik Fältström, Glenn Kowack, Lars-Johan Liman, and Bill Manning, “Report on the Impact on the DNS Root System of Increasing the Size and Volatility of the Root Zone, Prepared by the Root Scaling Study Team,” Version 1.0, September 2009.
- [2] “DNSSEC and Its Impact on DNS Performance,” 17 August 2009, <http://www.dnsops.gov/dnssec-perform.html>

- [3] Ray Bellis and Lisa Phifer, “Test Report: DNSSEC Impact on Broadband Routers and Firewalls,” SAC035, 16 September 2008,
<http://www.icann.org/en/committees/security/ssac-documents.htm>
- [4] Paul Vixie, “Extension Mechanisms for DNS (EDNS0),” RFC 2671, August 1999.
- [5] Peter Koch and Matt Larson, “Initializing a DNS Resolver with Priming Queries, Internet Draft, expired, July 2008,
<http://tools.ietf.org/id/draft-ietf-dnsop-resolver-priming-01.txt>
- [6] Roy Arends, Rob Austein, Matt Larson, Dan Massey, and Scott Rose, “DNS Security Introduction and Requirements,” RFC 4033, March 2005.
- [7] EDNS Support:
<http://www.ripe.net/data-tools/dns/as112/edns>
- [8] Matt Mathis, “The Case for Raising the Internet MTU,” July 2003, <http://staff.psc.edu/mathis/papers/Cisco200307/index.html>
- [9] Mohan Parthasarathy and Paul Vixie, “Representing DNS Messages Using XML,” Internet Draft, work in progress, September 2011, <http://www.ietf.org/id/draft-mohan-dns-query-xml-00.txt>
- [10] Ted Hardie, “Distributing Authoritative Name Servers via Shared Unicast Addresses,” RFC 3258, April 2002.
- [11] Yakov Rekhter, Robert G Moskowitz, Daniel Karrenberg, Geert Jan de Groot, and Eliot Lear, “Address Allocation for Private Internets,” RFC 1918, February 1996.

BILL MANNING has been in the network field since 1979, most recently with Booz Allen Hamilton. He has been an IETF Working Group chair, RFC author, and an ARIN Trustee, and he has been on numerous ICANN committees. He has worked as part of the teams that run Internet Root name servers, built the first Internet Exchange points, and worked on transitioning from NSFnet to commercial services. Current client work is focused on Internet Policy and Governance, Risk Analysis, and the future of naming systems. E-mail: bmanning@sfc.keio.ac.jp

Networking @ Home

by Geoff Huston, APNIC

One of the more interesting sessions at the *Internet Engineering Task Force* (IETF) meeting in Quebec City in July 2011 was the first meeting of the recently established *Homenet Working Group*^[1]. What is so interesting about networking the home? Well, if you regard challenges as “interesting,” then just about everything is interesting when you look at networking in the home!

It has been a very long time since the state of the art in home Internet involved plugging the serial port of the PC into the dialup modem. The *Asymmetric Digital Subscriber Line* (ADSL) modem, even when combined with some form of Wi-Fi base station, is looking distinctly passé these days. Today, the home network is seeing the intersection of a whole set of interests, including phone service, television service, home security services, energy management, utility service metering, other forms of home device monitoring, and, of course, connecting laptops and mobile devices to the net. The home network is not just a wired *Local-Area Network* (LAN), Wi-Fi home networks are commonplace, and there are also various Bluetooth devices. Maybe sometime soon it will be common for the home network to host some form of *Third-Generation* (3G) femtocell mobile cell phone repeater as well. But these days even that level of network complexity is not enough. Increasingly, the home office is part of the work office, and if numerous residents are at home, then the home network may be an endpoint for several corporate and institutional *Virtual Private Networks* (VPNs)^[2].

Within the home network we want sophisticated security. This security involves not just protecting the network from the neighbors; the security requirements include the ability for individuals to partition off their work-VPN part of the home network from other home users. For resiliency we might want a second network provider, so we might want to add site-based multihoming to the mix. And we need to make all this work for both IPv4 and IPv6.

That set of requirements represents a massive agenda. But to make this situation truly challenging, we cannot expect every home to come with an IT Operational Service Manager to ensure that all the various devices you bring into the home and connect to the network function as required for the particular requirements of the home. Indeed, we cannot expect any home to be so lavishly supported, nor can we afford to support home networking with a bevy of specialized call centers with on-demand support specialists, expert in the panoply of consumer devices that are being sold today.

With today's home networks, consumers are effectively on their own; and all this equipment better just work straight out of the box. No configuration, no buttons, it just has to work!

Routing @ Home

The evolution of networking at home has progressed from a single computer to a basic LAN, and from there to an Ethernet-bridged network with numerous Wi-Fi and wired LAN segments. All these environments have a single common architecture with a single “boundary” unit that acts as a point of demarcation between the *Internet Service Provider* (ISP) and the home network. This unit is generally called *Customer Premises Equipment* (CPE), and typically encompasses the functions of a modem; an IPv4 *Network Address Translator* (NAT); a *Dynamic Host Configuration Protocol* (DHCP) server for both IPv4 and IPv6; as well as security firewall, bridge, and rudimentary router functions.

But it is unrealistic to assume that home networks will continue to use a centralized model that places all of the management functions of the home network in a single unit. So how should we view home networks? Should home networks be a single bridged LAN, or are we seeing the evolution of home networks into multiple distinct domains with a routing fabric to glue them together? And if that is the case, what routing protocol should be used?

I have noticed in the low end of the CPE market it is not uncommon to see a rudimentary routing function supported by the *Routing Information Protocol* (RIP)^[3]. Thankfully, it is RIP Version 2, so the routing protocol can be configured with variable-length subnet masks, but even so, RIP is a very basic and simple routing protocol. But perhaps in this environment, that might be a positive factor rather than a liability in so far as RIP is simple enough to be auto-configurable. On the other hand, if there is an emergent need for more complex functions, then maybe we need to look a little harder at the available options.

One of these more complex functions is *subnet management*. In IPv6, the CPE will collect an IPv6 address prefix. This process differs from the conventional IPv4 environment where the CPE is typically assigned a single IPv4 address. So the ensuing question is: Is it possible to automate the distribution of IPv6 subnets across the entire home network? What form of management protocol is appropriate for this role?

Of course the situation gets much more complicated if the home network has two (or more) service providers. In the IPv6 environment, this task becomes a challenging one, not only with the distribution of multiple subnets across the home network, but also in the matter of exit path selection. If the home network is exercising due diligence to prevent source address spoofing, it is also necessary for the home routing infrastructure to deliver an outgoing packet to the “right” exit ISP, where the source address of the outgoing packet needs to match the address prefix provided by the corresponding ISP service. In other words, there is a requirement for source address routing in the home.

This challenge was not really addressed by the *Site Multi-Homing by IPv6 Intermediation Working Group* (SHIM6)^[4], despite the best of intentions, and it represents an even greater challenge if the intent is to provide mechanisms that can achieve such routing in an unmanaged home network environment.

I must admit to some concern here. We have managed to keep Internet routing working by using two principles. The first is to try to keep the routing task as simple as possible. Routing propagates a single “best” path to a destination. It does not necessarily do this propagation quickly, nor necessarily does it carry around with it a whole set of alternatives. It does just one job. The second principle is to admit that we have never really succeeded with the first principle of functional simplicity and we have always had expertise at hand to oversee the routing function and apply manual patches as required. The specialized requirements for the home network appear to be breaking both principles. The requirements are certainly not simple, and I see a mix of routing techniques—including various forms of policy-based routing requirements—entering the discussion. Secondly, there is no assurance that if things fail expertise is at hand to mend the failure. Indeed, the more complex the routing environment, the greater the potential for complex forms of failure. As we contemplate ever more complex requirements in the home network, we face a greater risk of encountering failure “by design,” where it is just not possible to design products for this environment that will “just work.”

Names @ Home

What should I call my printer? More to the point, how should I identify my Wi-Fi printer to all those devices at home that want to use it to print? I am sure that I would not like to use a proprietary naming scheme that requires me to add additional name resolution software to every device at home that wants to print something, nor do I want to transcribe IP addresses into everything. I would like my printer to get dynamically assigned IPv4 and IPv6 addresses when the device is plugged in and switched on, and have the name of the printer published via a generic name resolution mechanism, namely the *Domain Name System* (DNS).

But most of the time the rest of the world has no need to know the name of my printer at home, and I am not sure that it is a good move, securitywise, to gratuitously publish information in the public DNS. So what I would like for my printer is some form of “local” or “scoped” DNS, where I can name my printers, my disk servers, and other devices that I have at home in the context of my home and not have this information leak further afield. Is this scoped form of name resolution, split horizon DNS, or split views, possible in the context of the DNS without invoking further elements of configuration management?

Multicast DNS (mDNS) is perhaps one of the strongest candidates for this role. In essence, mDNS replaces the explicit client-server structure of the DNS with a scoped name subdomain of `.local` that is inherently scoped to the associated multicast domain.

This setup allows a client to perform DNS-like name resolution functions on a local network without the need to configure a conventional DNS server environment, and without the need to obtain global delegation of a site name in the global DNS.

An alternative approach is to use a conventional DNS delegation and conventional unicast DNS queries and responses. Clients are able to use DNS *Dynamic Updates*^[5] to provide the local DNS server with their details as they come online. This approach requires either open access from anyone to the nameserver or a security mechanism such as *Transaction SIGnature* (TSIG)^[6]. TSIG generally requires manual configuration, and alternatives are either little used—such as *Transaction KEY* (TKEY)^[7]—or involve further intricacies, such as Microsoft’s *Active Directory*, which uses other user authentication mechanisms to bootstrap the TSIG part using the *Generic Security Service Algorithm for Secret Key Transaction* (GSS-TSIG)^[8]. The DNS server itself can be advertised to all clients via the *Simple Service Discovery Protocol* (SSDP), as part of the larger *Universal Plug and Play* (UPnP) framework.

Sensing and Serving @ Home

Where to go from here? It is certainly the case that electronics has managed to pervade just about every device at home. Electricity meters are morphing into household energy-management systems, and many other household appliances are now controlled by internal processors. But individually configuring each of these devices is a forbidding task. Even adding an interface to allow manual configuration can often be a challenging objective.

The objective here is to define a standard mechanism to allow sensors to sense their local environment when powered up, obtain an IP address, advertise their existence and capabilities to the network, and, as appropriate, rendezvous with the sensor controller or controllers across the home network.

This example is another instance of a more generic class of automating the installation and use of services in “lightly” managed or even unmanaged networks, and it intersects significantly with the objectives encompassed with SSDP and UPnP. The potential volume of such devices places this example more squarely into a class of IPv6-only services, I suspect, which is a significant extension to the existing IPv4-centric UPnP frameworks.

What is needed is a bootstrap protocol that can provide a connecting device with:

- Address configuration
- Routing setup
- Name management and name server discovery
- Discovery of other services and controllers
- Security capabilities

Security @ Home

One of the most significant concerns with home networks lies in the area of security management. Host computers in a home network often want to place a very high level of implicit trust in their immediate network neighbors at the same home. It is not unusual for hosts in a home network to share printers, file servers, data, and even user profiles. Indeed, it is probably commonplace. But beyond this local security domain a host should become paranoid and treat all connection attempts with suspicion. But where does the local trust domain start and stop? What is the “local” security boundary?

This question is difficult to answer in an automated fashion. It is no longer the local LAN, particularly as home networks transition into routed networks. The security boundary is related to the local multicast scope, but this supposition assumes that it is possible to define a multicast scope that encompasses the local trust domain of the home network, and this assumption brings us back to the same question.

Even if you thought you might have a clean answer to the boundary question, you need to remind yourself about telecommuting. With telecommuting, there is a requirement to partition out an entire local network segment from the rest of the home environment and the home security domain and transplant it into the work security domain.

Everything @ Home

Home is certainly the new field of engagement for networked goods and services. However, it is one of the most challenging places to operate in from the perspective of attempting to deliver coherent services in a reliable and secure manner. The components are sourced from various vendors, and constructed incrementally over extended periods of time. It is an environment where older components need to coexist with new devices, and the overall engineering of the environment is at best piecemeal, and perhaps more often not engineered at all. In this environment out-of-the-box interoperability is of paramount importance, and therefore it is an environment where good standards really matter. Perhaps unsurprisingly, given these constraints, networking in the home is one of the environments that appear to raise the most challenges. It is an unforgiving environment where there is no real substitute for simplicity and reliability in a “plug-and-play” world.

The IETF Homenet Working Group has a lot of work to do. The Working Group will have to examine the diverse set of approaches in use today, add IPv6 functions, and produce a coherent set of outcomes in the form of standards that support robust, capable home networks that work in an unmanaged environment.

Ahhh home! There really is no place quite like it!

References

- [1] Homenet Working Group: <http://www.ietf.org/dyn/wg/charter/homenet-charter>
- [2] Paul Ferguson and Geoff Huston, “What Is a VPN?” (Part One and Part 2), *The Internet Protocol Journal*, Volume 1, No. 1 and No. 2, June and September 1998.
- [3] Gary Malkin, “RIP Version 2,” RFC 2453, November 1998.
- [4] Shim6 Working Group (concluded):
<http://wiki.tools.ietf.org/wg/shim6/charters>
- [5] Paul Vixie, ed., Yakov Rekhter, Susan Thomson, and Jim Bound, “Dynamic Updates in the Domain Name System (DNS UPDATE),” RFC 2136, December 1997.
- [6] Paul Vixie, Olafur Gudmundsson, Donald E. Eastlake 3rd, and Brian Wellington, “Secret Key Transaction Authentication for DNS (TSIG),” RFC 2845, May 2000.
- [7] Donald E. Eastlake 3rd, “Secret Key Establishment for DNS (TKEY RR),” RFC 2930, September 2000.
- [8] Stuart Kwan, Praerit Garg, James Gilroy, Levon Esibov, Randy Hall, and Jeff Westhead, “Generic Security Service Algorithm for Secret Key Transaction Authentication for DNS (GSS-TSIG),” RFC 3645, October 2003.

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005; he served on the Board of Trustees of the Internet Society from 1992 until 2001.
E-mail: gih@apnic.net

IETF Tools—Making It Easier to Make the Internet Work Better

by Robert Sparks

Many activities are associated with defining and refining an *Internet Engineering Task Force* (IETF) protocol, and all of them are detail-oriented. As IETF Working Groups are formed, mailing list discussions proceed, documents are written and reviewed, and interoperability is evaluated, participants encounter tasks that can be significantly simplified with the help of software tools. Fortunately, those participants frequently are also skilled software developers, and they create and share these tools as the need arises. A new paradigm has evolved recently: When a pressing need for a tool is identified—particularly one that has a large scope—the *IETF Administrative Oversight Committee* (IAOC) accelerates the creation of the tool by working with the community to gather requirements and financing the development of a solution. Comprehensive lists of available tools are maintained at [1] and at [2]. This article introduces a few important tools and discusses how you can help improve them or develop new ones.

Document Tools

The *Extensible Markup Language to Request For Comments* (XML2RFC)^[10] tool was developed to assist with Internet-Draft composition. Marshall Rose created and maintained the initial versions, capturing its input language and operation instructions in RFC 2629^[18]. This tool simplifies draft creation and maintenance by automatically producing documents that satisfy the RFC Editor's layout requirements, and assists in including the appropriate boilerplate as defined by the *IETF Trust*. It also simplifies the task of the *RFC Production Center*^[19, 20]. Starting with XML input rather than a draft in text form reduces the work required to create the RFC. The IAOC is currently funding a reimplementations of XML2RFC to reflect many years of user feedback, simplify maintenance—particularly of boilerplate handling—and make it easier for volunteers to contribute improvements. This reimplementations is currently available at [3]. Tony Hansen has been very active in gathering the requirements for and evaluating the reimplemented version. Julian Reschke also maintains *Extensible Stylesheet Language Transformations* (XSLT) code at [4] that translates RFC 2629-based input into several output formats.

After a new draft is prepared, Henrik Levkowetz' *Internet-Draft Nit Checker* (idnits) tool at [5] can scan it for any problems with the RFC Editor's checklist and guidelines and for other problems that drafts frequently encounter later in review. There are also tools for verifying sections of the document containing formal languages such as *Augmented Backus-Naur Form* (ABNF) or XML.

When an editor is satisfied that the document is ready to place in the repository, the automated *ID Submission tool*^[6] assists with an easy upload. At any point two versions of a draft can be compared with *rfcdiff*^[7], a flexible comparison program created by Henrik Levkowetz.

As a draft progresses, its history and current status can be tracked using the *Internet-Drafts Tracker* (ID Tracker) tool^[8]. This tool provides powerful search capabilities into the entire Internet-Draft repository, and a comprehensive view into the lifecycle of each Internet-Draft. With its roots in a tool to help the *Internet Engineering Steering Group* (IESG) keep track of drafts in IESG evaluation, the ID Tracker has evolved into a portal touching almost all aspects of IETF work. Each step of that evolution has improved efficiency and transparency, and has simplified access to the history of the development of each document.

Recent additions to the tracker allow for an easier capture of the details of Working Group processing. *Work in progress* will provide more visibility into the Working Group chartering and rechartering processes. The tracker is also used by other document streams. Many of the enhancements to the tracker are informed by the views into documents and Working Groups maintained by Henrik Levkowetz at [2]. The tracker continues to evolve through both IAOC-funded development efforts and volunteer contributions. An extension in progress will add visibility into the RFC Editor and *Internet Assigned Numbers Authority* (IANA) actions. When this extension is done the entire lifecycle of a Draft, from -00 submission to RFC publication, can be viewed in a single place.

Working Group and Meeting Tools

At each IETF meeting, a participant can build a custom view of the agenda using the tools at the datatracker and the tools sites. For example, [9] renders an interactive JavaScript-based calendar contributed by Adam Roach showing the *Real-Time Applications and Infrastructure* (RAI) meetings at IETF82. The pages at [10] provide a quick reference to the jabber rooms and audio streams of each Working Group meeting. The meeting materials tool facilitates uploading of agendas, slides, and minutes, which become available immediately through the agenda views.

Each Working Group has a Subversion Repository and an integrated instance of Trac^[21] at its disposal. The Subversion Repository can be used to maintain Working Group draft source, versioned instances of test documents, and even implementation code. IETF-specific customizations of the Trac system are described at [11]. Many Working Groups are already taking advantage of what the wiki Trac provides, and are using its ticketing feature to effectively track major Working Group document problems.

Notable examples are the problem tracking integrated into the *Hypertext Transfer Protocol Bis* (HTTPBIS) document status page at [12], and the summary of DISPATCH activity at [13]. The Trac wiki capability is also used by the Working Group Chairs at [14] and the IESG at [15].

IETF News

Keeping up with all of the activity across the IETF can be a challenge. One of the better tools for seeing what is happening is *The Daily Dose of the IETF*, created by Pasi Eronen, available at [16].

Again, this article is an introduction to just a few important tools. Comprehensive lists of available tools are maintained at [1] and [2].

Many of these tools were created because a person who needed them coded an initial version and contributed it to the community. Volunteers (and when needed, IAOC-funded efforts) then improve these tools over time. For several years, a group of volunteers have been meeting the Saturday before each IETF meeting for a day-long *Code Sprint*. If the existing tools need a minor tweak to make things work much better for you, or if you have an idea for a new tool you would like to start, please consider participating at the next Code Sprint. Between sprints, you can still help with the code. Refer to the sprint pages for an upcoming or recent sprint such as [17] and for information about getting started.

Whether or not you can contribute to the code, please discuss your ideas on the `tools-discuss@ietf.org` mailing list.

Several tool contributors have already been mentioned. Henrik Levkowetz deserves to be mentioned again. His herculean efforts maintaining `tools.ietf.org` and creating many of the tools there are of great benefit to the community.

References

- [0] Marshall T. Rose and Carl Malamud, "Writing Internet Drafts and RFCs Using XML," *The Internet Protocol Journal*, Volume 10, No. 1, March 2007.
- [1] <http://www.ietf.org/tools>
- [2] <http://tools.ietf.org/>
- [3] <http://xml.resource.org/>
- [4] <http://greenbytes.de/tech/webdav/rfc2629xslt/rfc2629xslt.html>
- [5] <http://tools.ietf.org/tools/idnits/>
- [6] <https://datatracker.ietf.org/submit/>
- [7] <http://www.ietf.org/tools/rfcdiff/>

- [8] <http://datatracker.ietf.org/>
- [9] <https://datatracker.ietf.org/meeting/82/agenda.html#RAI>
- [10] <http://tools.ietf.org/agenda/82/>
- [11] <http://trac.tools.ietf.org/misc/venue/wiki/IetfSpecificFeatures>
- [12] <http://tools.ietf.org/wg/httpbis/>
- [13] <http://trac.tools.ietf.org/wg/dispatch/trac/wiki>
- [14] <http://wiki.tools.ietf.org/group/wgchairs/>
- [15] <http://trac.tools.ietf.org/group/iesg/trac/wiki>
- [16] <http://tools.ietf.org/dailydose/>
- [17] <http://trac.tools.ietf.org/tools/ietfdb/wiki/IETF82Sprint>
- [18] Marshall T. Rose, “Writing I-Ds and RFCs using XML,” RFC 2629, June 1999.
- [19] Leslie Daigle, “RFC Editor in Transition: Past, Present, and Future,” *The Internet Protocol Journal*, Volume 13, No. 1, March 2010.
- [20] RFC Editor, “40 Years of RFCs,” RFC 5540, April 2009.
- [21] <http://trac.edgewall.org/about>

ROBERT SPARKS is an Area Director for the Real-Time Applications and Infrastructure Area (RAI) in the IETF. He previously chaired the IETF’s SIMPLE Working Group, which defines extensions to SIP for Presence and Instant Messaging, and the GEORPIV Working Group, which provides tools for applications to carry geographic location information and privacy rules to affect its use. Robert is a co-editor of the core SIP standard (RFC 3261), and several important SIP updates and extensions. He coordinates the premier real-time communications interoperability event, the SIPit. Robert is a Principal Software Engineer at Tekelec, and has held management and research positions at Estacado Systems, Xten (now Counterpath), dynamicssoft, Lucent, MCI Worldcom, and Texas A&M University. Robert holds a Master’s degree in Mathematics and a Bachelor’s degree in Computer Science from Texas A&M University. E-mail: rjsparks@nostrum.com

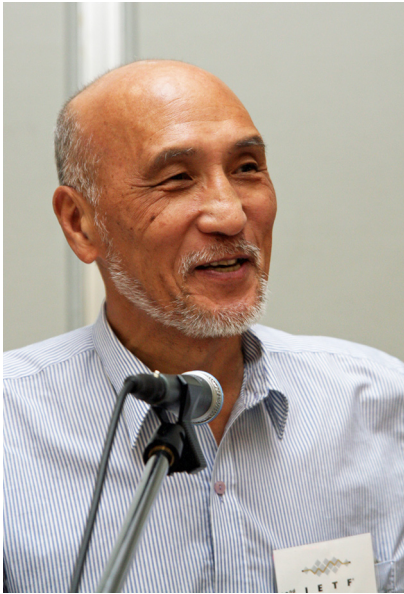


Photo: Peter Löthberg

Professor Kilnam Chon Receives 2011 Postel Service Award

The *Internet Society* (ISOC) recently announced that its prestigious *Jonathan B. Postel Service Award* was presented to leading technologist Professor Kilnam Chon for his significant contributions in the development and advancement of the Internet in Asia.

Professor Chon contributed to the Internet's growth in Asia through his extensive work in advancing Internet initiatives, research, and development. In addition, his pioneering work inspired many others to promote the Internet's further growth in the region. The international award committee, comprised of former Jonathan B. Postel award winners, noted that Professor Chon was active in connecting Asia, and that his efforts continue today in the advancement of the Internet in other regions.

The Postel Award was established by the Internet Society to honour individuals or organisations that, like Jon Postel, have made outstanding contributions in service to the data communications community.

Lynn St. Amour, President and CEO of ISOC, commented, "I met Professor Chon nearly fifteen years ago. He has long been a pioneer in the advancement of the Internet, striving to ensure its robust development. Beyond the amazing breadth of Professor Chon's work, perhaps his most remarkable achievement is his ability to inspire others. As a result of his work and the efforts of those he has motivated, Kilnam Chon has helped to ensure the global Internet is truly for everyone."

ISOC presented the award, including a US\$20,000 honorarium and a crystal engraved globe, during the 82nd meeting of the *Internet Engineering Task Force* (IETF) in Taipei, November 13–18, 2011.

The Internet Society is the world's trusted independent source of leadership for Internet policy, technology standards and future development. Based on its principled vision and substantial technological foundation, ISOC works with its members and Chapters around the world to promote the continued evolution and growth of the open Internet through dialog among companies, governments, and other organizations around the world. For more information about the Postel Service Award see: <http://www.isoc.org/postel/>

Alexandre Cassen and Rémi Després Receives 2011 Itojun Service Award

The third *Itojun Service Award* was presented to Alexandre Cassen and Rémi Després at the *Internet Engineering Task Force* (IETF) meeting held in Taipei, Taiwan in November 2011. The awardees were recognized for their design and implementation of "6rd," an IETF protocol that aims to speed the transition to global deployment of IPv6, which is critical to ensuring the continued growth and evolution of the Internet.

The 6rd protocol has been implemented by several *Internet Service Providers* (ISPs) around the world, including *Free Telecom*—the second largest ISP in France—as part of their efforts to deploy IPv6.

First awarded in 2009, the Itojun Service Award honors the memory of Dr. Jun-ichiro “Itojun” Hagino, who passed away in 2007 at the age of 37. The award, established by the friends of Itojun and administered by the *Internet Society* (ISOC), recognizes and commemorates the extraordinary dedication exercised by Itojun over the course of IPv6 development.

“Alexandre and Rémi’s efforts have helped to quickly bring a real IPv6 experience to hundreds of thousands of Internet users, demonstrating that IPv6 deployment can be effectively implemented on a large scale by commercial network providers,” said Jun Murai of the Itojun Service Award committee and founder of the WIDE Project. “On behalf of the Itojun Service Award committee, I am extremely pleased to present this award to Alexandre and Rémi for the significant work they have done to advance IPv6 development and deployment.”

The Itojun Service Award is focused on pragmatic contributions to developing and deploying IPv6 in the spirit of serving the Internet. The award, presented annually, includes a presentation crystal, a US\$3,000 honorarium and a travel grant.

Alexandre Cassen said, “It is truly an honor to have been selected to receive the Itojun Service Award. As a software developer myself, It is particularly touching to receive an award created in the memory of a coding legend such as Itojun. I would also like to thank the entire team at Free Telecom who, in 2007, implemented and deployed 6rd, allowing any subscriber who asked for IPv6 to have it with a single click. As I write this, Free Telecom has more than 1,500,000 subscribers using IPv6 every day, and all new subscribers have IPv6 enabled by default. IPv6 is happening Itojun!”

Rémi Després said, “The Itojun Award is the best possible recognition that long efforts to make IPv6 deployment practicable have been useful to the Internet community. Latecomer in IPv6 standardization, I was about to send my first email to Itojun on a technical issue when I heard of his death. I was even sadder since we undoubtedly would have otherwise enjoyed sharing our ideas and our enthusiasm. Sharing the honor of this award with Alexandre Cassen perfectly illustrates the great progress possible when a dynamic network operator with a pioneer spirit and talented engineers adopts an innovative and simple design. Making IPv6 operational on a large scale in only five weeks will be remembered as a milestone of both of our professional lives.”

More information on the Itojun Service Award is available at:
<http://www.isoc.org/itojun>

Internet Society Joins Opposition to Stop Online Piracy Act

The Internet Society Board of Trustees has expressed concern with a number of U.S. legislative proposals that would mandate *Domain Name System* (DNS) blocking and filtering by *Internet Service Providers* (ISPs) to protect the interests of copyright holders. While the Internet Society agrees that combating illicit online activity is an important public policy objective, these critical issues must be addressed in ways that do not undermine the viability of the Internet as a platform for innovation across all industries by compromising its global architecture. The Internet Society Board of Trustees does not believe that the *Protect-IP Act* (PIPA) and *Stop Online Piracy Act* (SOPA) are consistent with these basic principles.

Specifically, the Internet Society is concerned with provisions in both bills regarding DNS filtering. DNS filtering is often proposed as a way to block illegal content consumption by end users. Yet policies to mandate DNS filtering will be ineffective for that purpose and will interfere with cross-border data flows and services undermining innovation and social development across the globe.

Filtering DNS or blocking domain names does not remove the illegal content—it simply makes the content harder to find. Those who are determined to download filtered content can easily use a number of widely available, legitimately-purposed tools to circumvent DNS filtering regimes. As a result, DNS filtering encourages the creation of alternative, non-standard DNS systems.

From a security perspective, DNS filtering is incompatible with an important security technology called *Domain Name System Security Extensions* (DNSSEC). In fact, DNSSEC would be weakened by these proposals. This means that the DNS filtering proposals in SOPA and PIPA could ultimately reduce global Internet security, introduce new vulnerabilities, and put individual users at risk.

Most worrisome, DNS filtering and blocking raises human rights and freedom of expression concerns, and often curtails international principles of rule of law and due process. Some countries have used DNS filtering and blocking as a way to restrict access to the global Internet and to curb free expression.

The United States has been a strong proponent of online Internet freedoms and therefore has an important responsibility to balance local responsibilities and global impact, especially with respect to Internet policy. Given this commitment to global Internet freedom, it would be harmful to the global Internet if the United States were to implement such an approach.

“The Internet Society Board of Trustees is deeply concerned about the ramifications of the PIPA and SOPA bills on the overall stability and interoperability of the Internet,” said Raul Echeberria, Chairman of the Internet Society Board of Trustees.

“The Board recognizes that there can be misuses of the Internet; however, these are greatly outweighed by the positive uses and benefits of the Internet. We believe the negative impact of using solutions such as DNS blocking and filtering to address these misuses, far outweighs any short-term legal or business benefits.”

“The Internet Society believes that sustained, global collaboration amongst all parties is needed to find ways that protect the global architecture of the Internet while combating illicit online activities,” said Internet Society President and CEO Lynn St. Amour. “Mandating DNS blocking and filtering is simply not a viable option for the future of the Internet. We must all work together to support the principles of innovation and freedom of expression upon which the Internet was founded.”

For more details on DNS Filtering, visit:

<http://www.isoc.org/internet/issues/dns.shtml>

See also:

<https://www.eff.org/deeplinks/2011/12/internet-inventors-warn-against-sopa-and-pipa>

APNIC and JPRS Collaborate to Translate DNSSEC Technology Experiment Report

The *Asia Pacific Network Information Centre* (APNIC) has collaborated with *Japan Registry Services* (JPRS) to translate from Japanese into English the documents “DNSSEC Technology Experiment Report – Verification of Functionality and Performance” and “DNSSEC Technology Experiment Report – Operational Design.”

These documents contain the latest information on *Domain Name System Security Extensions* (DNSSEC) implementation, and provides information to those interested in implementing it. These reports are designed to introduce case studies to share knowledge and results gained through experiments conducted in 2010 that JPRS carried out in cooperation with Japanese ISPs, equipment vendors, and hosting providers.

APNIC would like to thank JPRS’s great initiative and all those involved in the process for making such an important contribution to DNSSEC awareness. APNIC also appreciates JPRS for making the documents available in English for wider distribution. The reports are available for download from:

<http://jprs.jp/dnssec/doc/DNSSEC-testbed-report-fpv1.0-E.pdf>

and

<http://jprs.jp/dnssec/doc/DNSSEC-testbed-report-odv1.0-E.pdf>

RFC Series Editor Appointment

The *Internet Architecture Board* (IAB) is pleased to announce the appointment of Heather Flanagan as the *Request For Comments Series Editor* (RSE). Ms. Flanagan will assume the responsibilities from the Acting RSE, Olaf Kolkman, and begin her tenure on January 1, 2012. The contract negotiated by the *IETF Administrative Oversight Committee* (IAOC) includes an initial term of two years and a presumptive renewal of two years.

Ms. Flanagan was selected by the *RFC Series Oversight Committee* (RSOC) based upon her experience, education, skills and energy she will bring to the position.

Ms. Flanagan is currently the Project Coordinator for the *COmanage* project, an effort funded by a grant from the *National Science Foundation* (NSF) and Internet2 to create a collaboration management platform, prior to that she was Director of Systems Administration, IT Services at Stanford University in Palo Alto, California. Her technical background is complemented by a Masters of Science of Library Science from the University of North Carolina, Chapel Hill that will prove invaluable in the accessing and indexing of RFCs.

Ms. Flanagan brings a high degree of energy and enthusiasm to the position. Her interpersonal skills as a facilitator and good listener will enable her to work well with the capable staff at the RFC Production Center and with the community in reaching consensus on a variety of issues facing the RFC Series.

The RSOC selection followed a lengthy process that included announcing the position inside and outside the community, several rounds of interviews, reference checks, and face-to-face interviews in Taipei at IETF 82. More than thirty-five applications were received, two-thirds of which were from outside the community.

We express our congratulations to Ms. Flanagan. We also want to extend our thanks to Ray Pelletier and the RSOC chaired by Fred Baker for their role in bringing the RSE selection process to a successful conclusion; to Olaf Kolkman for his service to the community as Acting RSE; to Joel Halpern for his ongoing work as editor of the “RFC Editor Model v2” document; and to the RFC Production Center for its customary diligence in the editing and publishing of RFCs this year, likely the second most productive in RFC publication history.

We look forward to working with the new RSE; we wish her well; and know that the community will work with Heather for the betterment of the RFC Series.

—For the IAB
Bernard Aboba, IAB Chair

2011 Global IPv6 Survey Results

On October 20, 2011 the *Number Resource Organization* (NRO) announced the publication of the “Global IPv6 Deployment Monitoring Survey 2011 Results,” initially previewed at the *Internet Governance Forum* (IGF) in Nairobi, Kenya, in September.

The findings from the survey drew on data supplied by around 1,600 international respondents, over 350 of which were from the *American Registry for Internet Numbers* (ARIN) region. On behalf of ARIN and GNKS Consulting, we would like to thank all who participated in the survey. Your feedback is crucial to expanding the understanding of where this community is moving, and what can be done to ensure readiness for the widespread adoption of IPv6. We hope you will take this opportunity to review the results at: http://www.nro.net/wp-content/uploads/ipv6_deployment_survey.pdf

The Public Switched Telephone Network in Transition

The United States *Federal Communications Commission* (FCC) recently held two workshops to examine the transition from the *Public Switched Telephone Network* (PSTN) to new technologies. Circuit-switched wireline voice technology has created a high standard for reliability, accessibility, and ubiquity. Consumers will continue to expect and demand these qualities, even as they shift from PSTN services to services provided over different networks. The transition away from the PSTN is already occurring, and is likely to accelerate. Through these workshops, the Commission will seek input on the technical, economic, and policy issues that must be addressed to minimize disruption during this transition, and to protect consumers, public safety, competition, and other important interests. For more information, visit: <http://www.fcc.gov/events/public-switched-telephone-network-transition-0>

Upcoming Events

The *North American Network Operators’ Group* (NANOG) will meet in San Diego, California, February 5–8, 2012. For more information see: <http://nanog.org>

The *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) will meet in New Delhi, India, February 21–March 2, 2012. For more information see: <http://www.apricot2012.net/>

The *Internet Engineering Task Force* (IETF) will meet in Paris, France, March 25–30, 2012. For more information see: <http://www.ietf.org/meeting/>

The *Internet Corporation for Assigned Names and Numbers* (ICANN) will meet in San Jose, Costa Rica, March 11–16, 2012 and in Prague, Czech Republic, June 24–29, 2012. For more information, see: <http://icann.org/>

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ contains standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2011 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

March 2012

Volume 15, Number 1

A Quarterly Technical Publication for
Internet and Intranet Professionals

In This Issue

From the Editor	1
Hacking Internet Security	2
DANE.....	12
Twenty-Five Years Ago	24
Letter to the Editor	36
Fragments	37

FROM THE EDITOR

Internet security continues to receive much attention both in the media and within the *Internet Engineering Task Force* (IETF) and similar organizations that develop technical solutions and standards. Last September, someone managed to break into a trusted *Certification Authority's* system and subsequently produced numerous fake *digital certificates*, files that comprise part of the architecture for what is generally referred to as “browser security.” In our first article, Geoff Huston describes what happened, the implications of this form of attack on the security of web-based services on the Internet, and what can be done to prevent similar attacks in the future.

In our second article, Richard Barnes describes the work of the *DNS-based Authentication of Named Entities* (DANE) working group in the IETF and explains how DANE, when deployed, can help prevent the sort of attack that is described in our first article.

This year I am celebrating 25 years in Internet technical publishing. Prior to launching *The Internet Protocol Journal* (IPJ), I was the editor of *ConneXions—The Interoperability Report*, published from 1987 until 1997 by Interop Company. With the generous support of *The Charles Babbage Institute* at the University of Minnesota, *ConneXions*, which was a paper-only publication, has been scanned and made available online. To mark the 25 combined years of *ConneXions* and IPJ, we asked Geoff Huston to examine the state of computer communications 25 years ago and give us his thoughts on where we have been and where we might be going in this rapidly developing technology landscape.

Please remember to check your subscription expiration date and take the necessary steps if you wish to continue receiving this journal. You will need your subscription ID and the e-mail address you used when you subscribed in order to access your record and renew online. Visit the IPJ website at www.cisco.com/ipj and click on the “Subscriber Services” link to get to the login page. The system will send you a URL that allows direct access to your record. If you no longer have access to the e-mail you used when you subscribed or you have forgotten your subscription ID, just send a message to ipj@cisco.com and we will make the necessary changes for you.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

Hacking Away at Internet Security

by Geoff Huston, APNIC

The front page story of the September 13, 2011, issue of the *International Herald Tribune* said it all: “Iranian activists feel the chill as hacker taps into e-mails.” The news story relates how a hacker has “... sneaked into the computer systems of a security firm on the outskirts of Amsterdam” and then “... created credentials that could allow someone to spy on Internet connections that appeared to be secure.” According to this news report, the incident punched a hole in an online security mechanism that is trusted by hundreds of millions of Internet users all over the network.

Other news stories took this hyperbole about digital crime and tapping into e-mail conversations on the Internet to new heights, such as *The Guardian’s* report on September 5, 2011, which claimed that the “... DigiNotar SSL certificate hack amounts to cyberwar, says expert.”^[1]

If application-level security is so vulnerable to attack, then this incident surely calls into question the basic mechanisms of trust and security upon which the entire global Internet has been constructed. By implication it also calls into question the trustworthiness of services operated by the major global Internet brands such as Google and Facebook, as much as it raises doubts about the levels of vulnerability for the use of online services such as banking and commercial transactions.

Just how serious is this problem? Are we now at the end of civilization as we know it?

Well, hardly!

Is digital cryptography now broken? Has someone finally managed to devise a computationally viable algorithm to perform prime factorization of massively large numbers, which lies at the heart of much of the cryptography used in the Internet today?

I really don’t think so. (At the very least, if someone has managed to achieve this goal, then that person is staying very quiet about it.).

Does this situation represent a systematic failure of security? Do we need to rethink the entire framework of cryptography and security in the Internet?

Not this time.

As far as I can tell, there has been no dramatic failure in the integrity of the digital technology used for security in the Internet today. Yes, some were surprised by this failure, including the Netherlands government, which uses certificates issued by the compromised certification authority, DigiNotar (<http://www.diginotar.com>) as part of its online service infrastructure. But the hacking incident was not based on a successful direct attack on the technology of cryptography by itself, and there is no reason to suppose that the strength of today's encryption algorithms is any weaker today than yesterday.

But in observing that the basic technology tools of the Internet security framework are still operating within acceptable bounds of integrity, and observing that this hacking attack did not create a gaping hole in our commitment to digital cryptography, what cannot be claimed is that the use of these cryptographic tools in today's Internet service environment is similarly trustworthy. The hacking attempt apparently was successful in so far as it provided the capability for third parties to impersonate trusted services and thereby capture users' private data, and evidently some people did indeed do precisely that, and that is not good at all.

Let's look a little more closely at this hacking episode and examine the way in which security is applied to the world of web browsing and the manner in which the vulnerabilities in this security framework were evidently exploited.

Securing a Connection

When I point my browser at my online banking service—or at any other secure website for that matter—a part of the browser navigation bar probably glows a reassuring green, and when I click it I get the message that I am connected to a website run by the Acme Banking corporation, and that my connection to this website has been encrypted to prevent eavesdropping. However, the website *certificate* was issued by some company that I have never even heard of. When I ask for more information, I am told the domain name, the company to whom the certificate for this domain name was issued, the identity of the certificate issuer, and the public key value. I am also reassuringly informed that the message I am viewing was encrypted before being transmitted over the Internet, and that this encryption makes it very difficult for unauthorized people to view information traveling between computers, and it is therefore very unlikely that anyone could read this page as it passes through the network. All very reassuring, and for the most part true, to the extent that we understand the strength of cryptographic algorithms in use today. The connection is using a *Transport-Layer Security* (TLS)^[2] connection and the traffic is encrypted using a private session key that should be impenetrable to all potential eavesdroppers.

But that is not the entire truth, unfortunately.

It may well be that your conversation is secure against eavesdropping, but it is only as secure as the ability of the other party to keep its private key a secret. If the other side of the conversation were to openly broadcast the value of its private key, then the entire encryption exercise is somewhat useless. So, obviously, my local bank will go to great lengths to keep its private key value a secret, and I rely on its efforts in order to protect my conversations with the bank.

But even then it is not quite the full story.

Am I really talking to my bank? Or in more general terms, am I really talking to the party with whom I wanted to talk?

The critical weakness in this entire framework of security is that the binding of certificates and keys to *Domain Name System* (DNS) names is not an intrinsic part of the DNS itself. It is not an extension of *Domain Name System Security Extensions* (DNSSEC)^[3, 4]. It has been implemented as an add-on module where third parties generate certificates that attest that someone has a particular domain name. Oddly enough, these *Certification Authorities* (CAs) may never have actually issued that particular domain name, because they are often disconnected from the DNS name registration business. Their business is a separate business activity where, after you have paid your money to a domain name registrar and secured your domain name, you then head to a domain name Certification Authority and pay them money (commonly they charge more money than the name registration itself) and receive a domain name certificate.

Certification Authorities

Who gets to be a Certification Authority? Who gets to say who has which domain name and what keys should be associated with that domain name?

Oddly enough the answer is, at a first level of approximation, just about anyone who wants to! I could issue a certificate to state that you have the domain name `www.example.com` and that your public key value is some number. The certificate I issue to that effect would not be much different from the certificates issued by everyone else. Yes, my name would be listed as the certificate issuer, but that is about all in terms of the difference between this certificate and the set of certificates you already trust through your browser.

So what is stopping everyone from being a Certification Authority? What is preventing this system from descending into a chaotic environment with thousands of certificate issuers?

For this situation the browser software folks (and other application developers of secure services) have developed a solution. In practice it requires a lot of effort, capability, diligence, and needless to say, some money, to convince a browser to add your Certification Authority public key to its list of trusted Certification Authorities.

You have to convince the browser developers that you are consistently diligent in ensuring that you issue certificates only to the “correct” holders of domain names and that you undertake certificate management practices to the specified level of integrity and trust. In other words, you have to demonstrate that you are trustworthy and perform your role with consistent integrity at all times. You then get listed with all the other trusted Certification Authorities in the browser, and users will implicitly trust the certificates you issue as part of the security framework of the Internet.

How many trusted Certification Authorities are there? How many entities have managed to convince browser manufacturers that they are eminently trustable people? If you are thinking that this role is a special one that only a very select and suitably *small* number of folks who merit such absolute levels of trust should undertake for the global Internet—maybe two or three such people—then, sadly, you are very much mistaken.

Look at your browser in the preferences area for your list of trusted Certification Authorities, and keep your finger near the scroll button, because you will have to scroll through numerous such entities. My browser contains around 80 such entities, including one government (“Japanese Government”), a PC manufacturer (“Dell Inc”), numerous telcos, and a few dedicated certificate issuers, including DigiNotar.

Do I know all these folks that I am meant to trust? Of course not! Can I tell if any of these organizations are issuing rogue certificates, deliberately—or far more likely—inadvertently? Of course not!

The structural weakness in this system is that a client does not know *which* Certification Authority—or even which duly delegated subordinate entity of a Certification Authority—was used to issue the “genuine” DNS certificate. When a client receives a certificate as part of the TLS initialization process, then as long as any one of the listed trusted Certification Authorities is able to validate the presented certificate, even if it is the “wrong” Certification Authority, then the client will proceed with the session with the assumption that the session is being set up with the genuine destination.

In other words the entire certification setup is only as strong—or as weak—as the weakest of the certification authorities. It really does not matter to the system as a whole if any single Certification Authority is “better” at its task than the others, because every certified domain name is protected only to the extent that the “weakest” or most vulnerable trusted Certification Authority is capable of resisting malicious attack and subversion of its function. Indeed, one could argue that there is scant motivation for any trusted Certification Authority to spend significantly more money to be “better” than the others, given that its clients are still as vulnerable as all the other clients of all the other Certification Authorities.

In other words, there is no overt motivation for market differentiation based on functional excellence, so all certificates are only as strong as the weakest of all the Certification Authorities. And therein lies the seed of this particular hacking episode.

The Hack

The hack itself now appears to have been just another instance of an online break-in to a web server. The web server in question was evidently running the service platform for DigiNotar, and the hacker was able to mint some 344 fraudulent certificates, where the subject of the certificate was valid, but the public key was created by the hacker. A full report of the hacking incident was published by Fox-IT^[5].

To use these fraudulent certificates in an attack requires a little more than just minting fraudulent certificates. It requires traffic to be redirected to a rogue website that impersonates the webpage that is under attack. This redirection requires collusion with a service provider to redirect client traffic to the rogue site, or a second attack, this time on the Internet routing system, in order to perform the traffic redirection.

So minting the fraudulent certificates is just one part of the attack. Were these fake certificates used to lure victims to fake websites and eavesdrop on conversations between web servers and their clients? Let's look at the client's validation process to see if we can answer this question.

When starting a TLS session, the server presents the client with a certificate that contains the server public key. The client is expected to validate this certificate against the client's locally held set of public keys that are associated with trusted certification authorities. Here is the first vulnerability. The client is looking for any locally cached trusted key to validate this certificate. The client is not looking as to whether a particular public key validates this certificate. Let's say that I have a valid certificate issued by the Trusted Certification Authority Inc. for my domain name, **www.example.com**. Let's also say that the server belonging to another Certification Authority, Acme Inc, is compromised, and a fake certificate is minted. If a user is misdirected to a fake instance of **www.example.com** and the bad server passes the client this fake certificate, the client will accept this fake certificate as valid because the client has no *presumptive* knowledge that the only key that should validate a certificate for **www.example.com** belongs to the Trusted Certification Authority Inc. When the key belonging to Acme Inc validates this certificate and ACME is a trusted entity according to my browser, then that is good enough to proceed.

Actually that is not the full story. What if I wanted to cancel a certificate? How do certificates get removed from the system and how do clients know to discard them as invalid?

A diligent client (and one who may need to check a box in the browser preference pane to include this function) uses a second test for validity of a presented certificate, namely the *Online Certificate Status Protocol* (OCSP)^[6]. Clients use this protocol to see if an issued certificate has been subsequently revoked. So after the certificate has been validated against the locally held public key, a diligent application will then establish a secure connection to the certification authority OCSP server and query the status of the certificate.

This secure connection allows for prompt removal of fraudulent certificates from circulation. It assumes of course that clients use OCSP diligently and that the Certification Authority OCSP server has not also been compromised in an attack, but in an imperfect world this step constitutes at least another measure of relative defence.

The OCSP server logs can also provide an indication of whether the fraudulent certificates have been used by impersonating servers, because if the certificate was presented to the client and the client passed it to an OCSP server for validation, then there is a record of use of the certificate. The Fox-IT report contains an interesting graphic that shows the geolocation of the source addresses of clients who passed a bad *.google.com certificate to OCSP for validation. The source addresses have a strong correlation to a national geolocation of Iran.

Obviously this attack requires some considerable sophistication and capability, hence the suspicion that the attack may have had some form of state or quasi-state sponsorship, and hence the headlines from *The Guardian*, quoted at the start of this article, that described this attack as an incident of cyberwarfare of one form or another. Whether this incident was a cyber attack launched by one nation state upon another, or whether this was an attack by a national agency on its own citizens is not completely clear, but the available evidence points strongly to the latter supposition.

Plugging the Hole?

This incident is not the first such incident that has created a hole in the security framework of the Internet, and it is my confident guess that it will not be the last. It is also a reasonable guess that the evolution of the sophistication and capability that lie behind these attacks points to a level of resourcing that leads some to the view that various state-sponsored entities may be getting involved in these activities in one way or another.

Can we fix this?

It seems to me that the critical weakness that was exploited here was the level of disconnection between domain name registration and certificate issuance. The holders of the domain names were unaware that fraudulent certificates had been minted and were being presented to users as if they were the real thing. And the users had no additional way of checking the validity of the certificate by referring back to information contained in the DNS that was placed there by the domain name holder.

The end user was unable to refine the search for a trusted Certification Authority that would validate the presented certificate from all locally cached trusted Certification Authorities to the one certification authority that was actually used by the domain name holder to certify the public key value. So is it possible to communicate this additional information to the user in a reliable and robust manner?

The last few years have seen the effort to secure the DNS gather momentum. The root of the DNS is now DNSSEC-signed, and attention is now being focused on extending the interlocking signature chains downward through the DNS hierarchy. The objective is a domain name framework where the end client can validate that the results returned from a DNS query contain authentic information that was entered into the DNS by the delegated authority for that particular DNS zone.

What if we were able to place certificates—or references to certificates—into the DNS and protect them with DNSSEC? The *DNS-based Authentication of Named Entities* (DANE) Working Group of the IETF^[0, 7] is considering this area of study. They are considering numerous scenarios at present, and the one of interest here does not replace the framework of Certification Authorities and domain name certificates, but it adds another phase of verification of the presented certificate.

The “Use Cases”^[8] document from the DANE working group illustrates the proposed approach. I will quote a few paragraphs from this document. The first paragraph describes the form of attack that was perpetrated in June and July this year on the DigiNotar CA. It is not clear to me if the text predates this attack or not, but they are closely aligned in time:

“Today, an attacker can successfully authenticate as a given application service domain if he can obtain a ‘mis-issued’ certificate from one of the widely-used CAs—a certificate containing the victim application service’s domain name and a public key whose corresponding private key is held by the attacker. If the attacker can additionally insert himself as a man in the middle between a client and server (for example, through DNS cache poisoning of an A or AAAA record), then the attacker can convince the client that a server of the attacker’s choice legitimately represents the victim’s application service.”^[8]

So how can DNSSEC help here?

“With the advent of DNSSEC [RFC 4033], it is now possible for DNS name resolution to provide its information securely, in the sense that clients can verify that DNS information was provided by the domain holder and not tampered with in transit.

The goal of technologies for *DNS-based Authentication of Named Entities* (DANE) is to use the DNS and DNSSEC to provide additional information about the cryptographic credentials associated with a domain, so that clients can use this information to increase the level of assurance they receive from the TLS handshake process.

This document describes a set of use cases that capture specific goals for using the DNS in this way, and a set of requirements that the ultimate DANE mechanism should satisfy. Finally, it should be noted that although this document will frequently use HTTPS as an example application service, DANE is intended to apply equally to all applications that make use of TLS to connect to application services named by domain names.”^[8]

Does DANE represent a comprehensive solution to this security vulnerability?

I would hesitate to be that definitive. As usual with many aspects of security, the objective of the defender is to expend a smaller amount of effort in order to force an attack to spend a far larger amount of effort. From this perspective, the DANE approach appears to offer significant promise because it interlocks numerous security measures and forces a potential attacker to compromise numerous independent systems simultaneously. Within the DANE framework the attacker cannot attack any certification authority, but must compromise a particular certification authority, and the attacker must also attack DNSSEC and compromise the information contained in signed DNS responses for that domain in order to reproduce the effects of the attack described here. This scenario seems to fit the requirement of a small amount of additional defensive effort by the server and the client, creating a significantly larger challenge to the attacker.

But many preconditions must be met here for this approach to be effective:

- DNSSEC needs to be ubiquitously deployed and maintained.
- Issued DNS certificates need to be published in the secure DNS zone using the DANE framework.
- Client DNS resolvers need not only to be DNSSEC-aware, but also to enforce DNSSEC outcomes.
- Applications, including browsers, need to validate the certificate that is being used to form the TLS connection against the information provided by a validated DNS response for the DANE credentials for that DNS zone.

It is probably not perfect, but it is a large step forward along a path of providing more effective security in the Internet.

Unfortunately, this solution does not constitute an instant solution ready for widespread use today—or even tomorrow. We could possibly see this solution in widespread use in a couple of years, but, sadly, it is more likely that securing the DNS for use in the Internet will not receive adequate levels of attention and associated financial resourcing in the coming years. It may take upward of 5 years before we see ubiquitous adoption of DNSSEC and any significant levels of its use by a DANE framework for certificates in the DNS. Until then there is the somewhat worrisome prospect of little change in the framework of Internet security from that used today, and the equally concerning prospect that this particular hacking event will not be the last.

Acknowledgement

I am indebted to Olaf Kolkman of NLnet Labs for a stimulating conversation about this attack and the implications for securing the Internet. NLnet Labs is one of a small number of innovative and highly productive research groups that has developed considerable levels of expertise in this area of security and the DNS.^[9]

Postscript

When you lose that essential element of trust, your continued existence as a trusted Certification Authority is evidently a very limited one. On Tuesday September 20, 2011, the Dutch company DigiNotar was officially declared bankrupt in a Haarlem court.

Disclaimer

The views of this article do not necessarily represent the views or positions of the Asia Pacific Network Information Centre.

References

- [0] Richard L. Barnes, “Let the Names Speak for Themselves: Improving Domain Name Authentication with DNSSEC and DANE,” *The Internet Protocol Journal*, Volume 15, No. 2, March 2012.
- [1] <http://www.guardian.co.uk/technology/2011/sep/05/digi-notar-certificate-hack-cyberwar>
- [2] William Stallings, “SSL: Foundation for Web Security,” *The Internet Protocol Journal*, Volume 1, No. 1, June 1998.
- [3] Miek Gieben, “DNSSEC: The Protocol, Deployment, and a Bit of Development,” *The Internet Protocol Journal*, Volume 7, No. 2, June 2004.
- [4] Roy Arends, Rob Austein, Matt Larson, Dan Massey, and Scott Rose, “DNS Security Introduction and Requirements,” RFC 4033, March 2005.

- [5] Fox IT, “DigiNotar Certificate Authority breach, ‘Operation Black Tulip,’”
<http://www.rijksoverheid.nl/bestanden/documenten-en-publicaties/rapporten/2011/09/05/diginotar-public-report-version-1/rapport-fox-it-operation-black-tulip-v1-0.pdf>
- [6] Michael Myers, Rich Ankney, Ambarish Malpani, Slava Galperin, and Carlisle Adams, “X.509 Internet Public Key Infrastructure Online Certificate Status Protocol – OCSP,” RFC 2560, June 1999.
- [7] <http://datatracker.ietf.org/wg/dane/>
- [8] Richard Barnes, “Use Cases and Requirements for DNS-Based Authentication of Named Entities (DANE),” RFC 6394, October 2011.
- [9] <http://nlnetlabs.nl/>
- [10] On March 26, 2012, at IETF 83 in Paris, France, a Technical Session with the title “Implementation Challenges with Browser Security” was held. The following presentations were given:
- Hannes Tschofenig: “Introduction”
 - Eric Rescorla: “How do we get to TLS Everywhere?”
 - Tom Lowenthal: “Cryptography Infrastructure”
 - Chris Weber: “When Good Standards Go Bad”
 - Ian Fette: “Lessons Learned from WebSockets (RFC 6455)”
 - Jeff Hodges: “It’s Not the End of the World”
- All of these presentations are available from:
<https://datatracker.ietf.org/meeting/83/materials.html>

GEOFF HUSTON B.Sc., M.Sc., is the Chief Scientist at *Asia Pacific Network Information Centre* (APNIC), the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. He is author of numerous Internet-related books, was a member of the Internet Architecture Board from 1999 until 2005, and served on the Board of Trustees of the Internet Society from 1992 until 2001. E-mail: gih@apnic.net

Let the Names Speak for Themselves: Improving Domain Name Authentication with DNSSEC and DANE

by Richard L. Barnes, BBN Technologies

Authentication of domain names is a fundamental function for Internet security. In order for applications to protect information from unauthorized disclosure, they need to make sure that the entity on the far end of a secure connection actually represents the domain that the user intended to connect to. For many years, authentication of domain names has been accomplished by having third-party *Certification Authorities* attest to which entities could represent a domain name. This system of external authorities, however, has recently come under heavy attack, and there have been several high-profile compromises^[0]. The *Domain Name System Security Extensions* (DNSSEC) offer an alternative channel for distributing secure information about domain names, through the *Domain Name System* (DNS) itself. The *DNS-based Authentication of Named Entities* (DANE) working group in the *Internet Engineering Task Force* (IETF) has developed a new type of DNS record that allows a domain itself to sign statements about which entities are authorized to represent it. End users' applications can use these records either to augment the existing system of Certification Authorities or to create a new chain of trust, rooted in the DNS.

Authentication

Without authentication, other security services are moot. There is little point in Alice's encrypting information en route to Bob if she has not first verified that she is talking to Bob and not an attacker Eve. In the context of Internet applications, authentication is about ensuring that users know whom they are talking to, and in most cases, that "whom," is represented by a domain name. For example, in the *Hypertext Transfer Protocol* (HTTP), the "authority" section of a *Uniform Resource Identifier* (URI) indicates the domain name of the server that will fulfill requests for that URI. So when an HTTP user agent starts a TCP connection to a remote server, it needs to verify that the server is actually authorized to represent that domain name^[1].

The most common security protocol used by Internet applications is the *Transport Layer Security* (TLS) protocol^[2]. TLS provides a layer above TCP that facilitates authentication of the remote side of the connection as well as encryption and integrity protection for data. TLS underlies *Secure HTTP* (HTTPS) and secure e-mail^[1, 3, 4], and provides hop-by-hop security in real-time multimedia and instant-messaging protocols^[5, 6]. In all of these applications, the server that the user ultimately wants to connect to is identified by a DNS domain name^[7, 8]. A user might enter `https://example.com` into a web browser or send an e-mail to `alice@example.com`.

One of the main purposes of using TLS in these cases is thus to assure the user that the entity on the other end of the connection actually represents `example.com`; in other words, to authenticate the server as a legitimate representative of the domain name. Note that these comments apply to *Datagram Transport Layer Security* (DTLS) as well, because it provides the same functions as TLS for *User Datagram Protocol* (UDP) packet flows^[9].

Today, a server asserts its right to represent a domain by presenting a *Public Key Infrastructure* (PKIX) digital certificate containing that domain^[8, 10]. A certificate is an attestation by a Certification Authority of a binding between a public key and a name—the entity holding the corresponding private key is authorized to represent that name. TLS ensures that only the holder of a given private key can read the encrypted data; the certificate ensures that the holder of the key represents the desired name.

Current TLS-based applications maintain a list of Certification Authorities whose certificates they will accept. Unfortunately, over time, these lists have grown very long, with major web browsers trusting nearly 200 Certification Authorities, representing a diverse range of organizations. Because any of these Certification Authorities can vouch for any domain name, a long list creates many points of vulnerability; a compromise at any point allows the attacker to issue certificates for any domain. Several recent attacks have taken advantage of this fact by targeting smaller Certification Authorities as a way to obtain certificates for major domains. For example, an attack through DigiNotar against Google is discussed in this issue^[0].

DNSSEC offers an alternative to Certification Authorities. In the DNSSEC system, each domain holder can act as an authority for subordinate domains. The IETF DANE working group has developed a DNS record format for “certificate associations,” so that domain holders can sign statements about which certificates can be used to authenticate as that domain. In effect, this scenario allows a domain to speak for itself, instead of through a third-party Certification Authority. DANE associations can be used either as a check on the current model (for example, to limit which Certification Authorities may vouch for a domain) or as an alternative trust path, rooting trust in a DNSSEC authority instead of a Certification Authority. Work on the protocol document is drawing to a close, and several prototype implementations are already in progress.

Background: PKIX and DNSSEC

At one level, the choice of which authentication technology to use is a choice of authorities and scoping. As mentioned previously, authentication is fundamental for security, but it is also very hard to accomplish scalably. For example, a web browser needs to be able to authenticate any website the user chooses to visit. It would clearly not work for each browser vendor to send a human representative to meet every website owner in order to find out what public key should be used for that website.

So instead of relying on having preestablished relationships with every entity we want to authenticate, we rely on centralized authorities to do identity checking. The authorities then create credentials that anyone else can check, so that if the credential is valid and you believe the authority is trustworthy, then the entity holding the credential has the indicated identity.

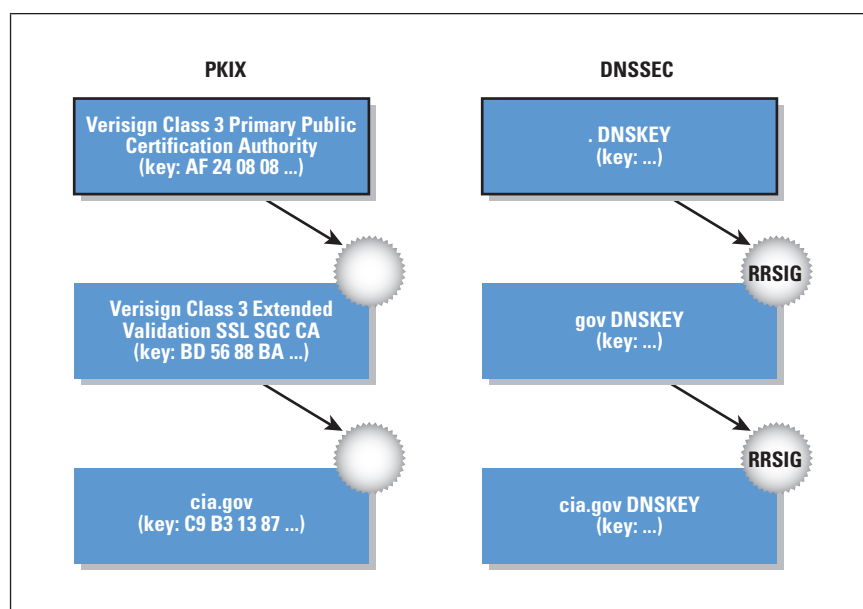
In a technical sense, an entity holds a credential if it holds the private key corresponding to the public key in the credential. The credential encodes a *binding* between the public key and the identity, asserted by the authority.

Authority is of course not a purely digital concept. If we want to know a person's name in real life we do not just ask them directly, because the person could lie. Instead we look to a credential issued by an authority, such as a driver's license or birth certificate. So the technology question here is how to manage authorities, and how to encode these credentials.

The IETF has defined two major cryptographic authority systems: PKIX, based on digital certificates^[10]; and DNSSEC, based on the DNS^[11]. Both of these systems allow authorities to associate public keys with identities, and both arrange these authorities hierarchically.

The hierarchy is important because it allows a *relying party* (someone who is verifying identities) to choose whom to trust. In these hierarchical systems, an authority's identity can itself be attested by a credential issued by another authority. When a relying party wants to verify a credential issued by an authority A, he then has to verify that A's credential is valid (under an authority B), and so on until he reaches an authority that he trusts. This sequence of credentials constitutes a logical path through the hierarchy, known as a "certification path" in PKIX terminology (Figure 1).

Figure 1: PKIX and DNSSEC Trust Hierarchies



In order to be useful as a given relying party to authenticate someone, a certification path has to end in a *trust anchor*, that is, an authority that the relying party trusts to make assertions. In the DNSSEC context, relying parties can in principle have only one trust anchor, namely the DNS root, although alternatives to the root have been proposed^[12]. The PKIX system, on the other hand, does not represent a single globally consistent hierarchy, so in order to be able to validate many certificates, relying parties often have to choose many trust anchors.

Crossing the Streams

Current TLS-based applications rely on PKIX for authentication of domain names, which has facilitated fairly broad deployment, but also created some vulnerabilities. PKIX is based on a very general digital certificate system called X.509, and because of this generality, it has no inherent binding to the DNS. This situation creates two problems when it comes to authenticating domain names.

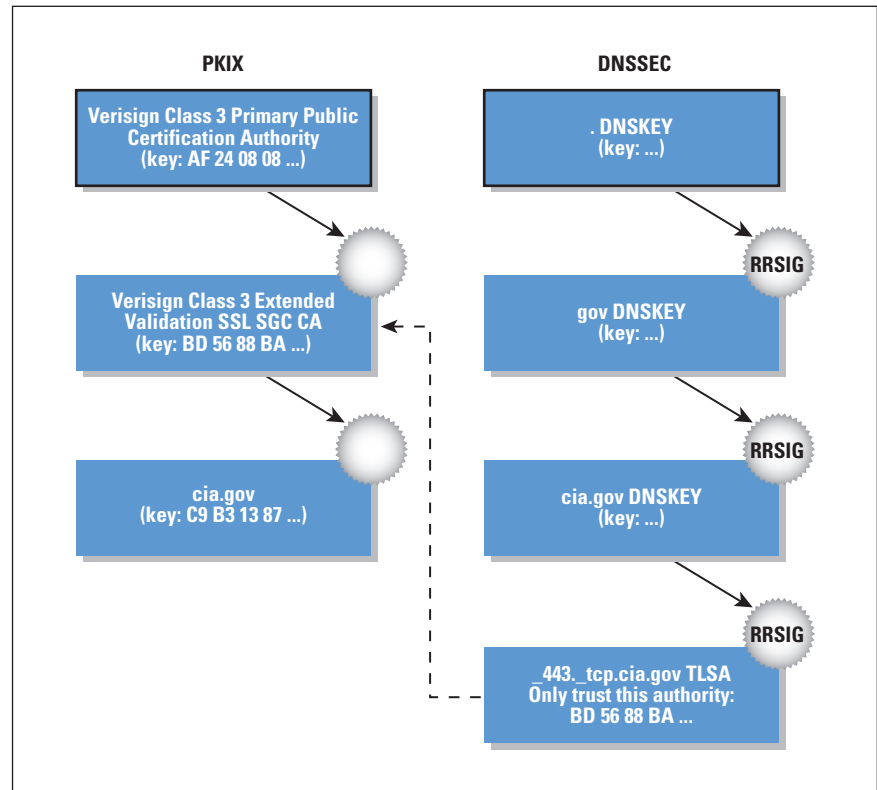
First, unlike the DNS, which has a single global root, there is no single authority under which all PKIX certificates can be verified. Indeed, there is an open marketplace of authorities, where each entity can choose which authority will sign its certificate, leaving relying parties with a choice: Either they must trust every authority that has signed a certificate for an entity it wants to authenticate, or they will be unable to validate the identities of some entities. In general, current software has preferred the former approach of trusting many authorities, to the extent that modern browsers and operating systems will trust up to 200 authorities by default. Users can add to this list, for example, using the “Accept this certificate?” dialogs in their browsers, but it can be very difficult to remove trust anchors from the default list^[13].

Second, PKIX authorities today are not constrained in the scope, so they can issue credentials for any name—even those for whom they have no real information (in contrast to the DNS—where each zone can vouch only for sub-domains; only the root can act with impunity). Conversely, there is no real way for a relying party to know what authority should be vouching for a site, so if a rogue authority were to issue a certificate to an unauthorized party, relying parties would have no way to detect it.

Given these vulnerabilities, any of the many authorities trusted within the PKIX system can attack any domain by issuing a false certificate from that domain. This false certificate can then be used to masquerade as the victim domain, for example, to perform a man-in-the-middle attack. Note that the authority itself is not necessarily the bad actor in this attack—it could be an external attacker that can obtain illicit access to the systems that issue certificates. The risks of having broadly trusted Certification Authorities have recently become clear, because attackers were able to break into two small Certification Authorities and create fraudulent certificates for Google and Facebook, among others^[14, 15].

The goal of DANE is to address some of the vulnerabilities of the current PKIX ecosystem by allowing DNSSEC—to “cross the streams” to allow domains to publish information secured with DNSSEC that can add additional security to PKIX certificates used for TLS. For example, a domain might use DANE to inform relying parties of which authorities can be trusted, as illustrated in Figure 2.

Figure 2: Using a DANE TLS Associations (TLSA) Record to Indicate Which PKIX Authority Should Be Trusted



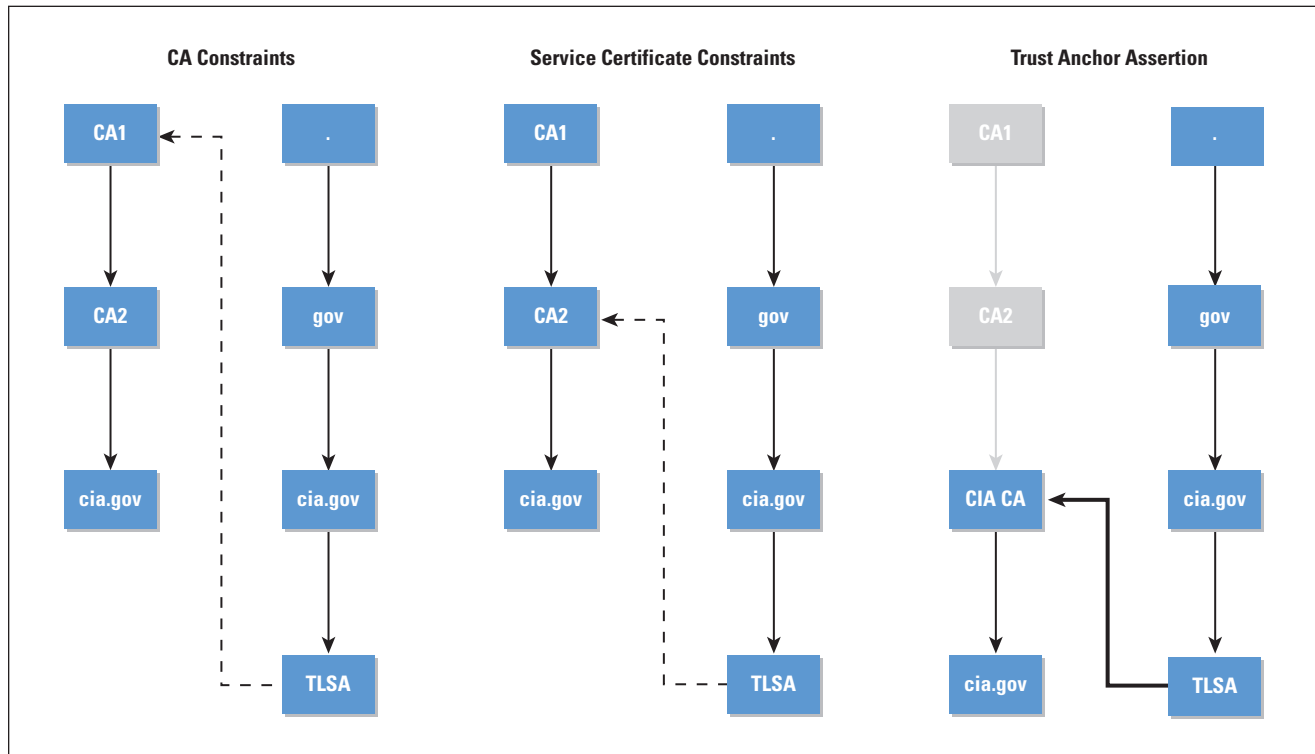
DANE Records

If the goal of DANE is to allow domain operators to make statements about how clients should judge TLS certificates for their domains, then what sorts of statements should DANE allow them to make? The DANE use cases document^[16] lays out three major types of statements (Figure 3):

1. *CA Constraints:* The client should accept only certificates issued under a specific Certification Authority.
2. *Service Certificate Constraints:* The client should accept only a specific certificate.
3. *Trust Anchor Assertion:* The client should use a domain-provided trust anchor to validate certificates for that domain.

All three of these statements can be viewed as constraining the scope of trust anchors. The first two types limit the scope of existing trust anchors, whereas the third provides the client with a new trust anchor (still within a limited scope). More on these anchors in a moment.

Figure 3: DANE Use Cases



The current draft DANE protocol defines a DNS Resource Record type TLSA for describing “*TLS Associations*”—statements about what certificates are “associated” to a domain^[17]. Each TLSA record has three basic fields:

- *Usage*: Which type of statement this record is making
- *Selector/Matching*: How a TLS certificate chain should be matched against this record (for example, by exact match, by public key, or by SHA-1 digest)
- *Certificate for Association*: The actual data against which the TLS certificate chain should be matched

These records are stored under the target domain with a prefix that indicates the transport and port number for the TLS server. So for example, if Alice runs a secure web service at **example.com** and wants to tell clients that they should accept only certificates from the Charlie’s CA, she could provision a TLSA record under **_443._tcp.example.com** with the following contents:

- *Usage*: CA constraint
- *Selector/Matching*: SHA-1 digest
- *Certificate for Association*: SHA-1 digest of Charlie’s certificate

When a client Bob wants to connect to **https://example.com**, he can find these TLSA records and apply Alice’s constraints when he validates the server certificate.

Adding Constraints to PKIX

The major objective of the CA constraints and service certificate constraints is to guard against “mis-issue” of certificates. A certificate is “mis-issued” when a CA issues a certificate to an entity that does not actually represent the domain name in the certificate. Mis-issue can come about in many ways, including through malicious Certification Authorities, compromised Certification Authorities (as in the Comodo and DigiNotar example discussed previously), or Certification Authorities that are simply misled as to the attacker’s identity through fraud or other means. Today, mis-issue can be difficult to detect, because there is no standard way for clients to figure out which Certification Authorities are supposed to be issuing certificates for a domain. When an attacker issued false certificates for the Google Gmail service under the DigiNotar Certification Authority, it was noticed only because a vigilant user posted to a Gmail help forum.^[18]

By contrast, domain operators know exactly which Certification Authorities they have requested certificates from, and, of course, which specific certificates they have received. With DANE, the domain operator can convey this information to the client. For example, to guard against the DigiNotar attack, Google could have provisioned a TLSA record expressing a Certification Authority constraint with its real Certification Authority (which is not DigiNotar) or a certificate constraint with its actual certificate. Then DANE-aware clients would have been able to immediately see that the DigiNotar certificates were improperly issued and possibly indicative of a man-in-the-middle attack.

Empowering Domain Operators

According to data from the EFF SSL Observatory, which scans the whole IPv4 address space for HTTPS servers and collects their certificates, around 48 percent of all HTTPS servers present self-signed certificates^[19]. An unknown number of other servers present certificates issued under Certification Authorities that are not in the major default trust anchor lists. For example, the United States Air Force web portal uses a certificate issued under a Department of Defense Certification Authority that is not trusted by Firefox^[20]. In the current environment, most clients cannot authenticate these servers at all; they have to rely on users manually checking certificates, hopefully with some out-of-band information. As a result, these servers and their users are highly vulnerable to man-in-the-middle attacks against their supposedly secure sessions.

DANE Trust Anchor Assertions enable the operators of a domain to advertise a new trust anchor, under which certificates for that domain will be issued. Using these records, clients can dynamically discover what trust anchors they should accept for a given domain, instead of relying on a static list provided by a browser or operating system.

It may seem odd to talk about a domain supplying a client with trust anchors, because trust anchor provisioning is typically a very sensitive activity. If an attacker is able to install a trust anchor into a victim's trust anchor store, then the attacker can masquerade under any name he wants by issuing certificates under that name. The PKIX working group even defined a whole protocol for managing trust anchors^[21].

DANE ensures that this trust anchor provisioning is secure by applying scoping and verifying that scoping using DNSSEC. DANE trust anchor assertions are scoped to a particular domain name, so even if an attacker can introduce a false trust anchor, he can use it to spoof only a single name. Furthermore, trust anchor assertions must be DNSSEC-signed, so clients can verify that the entity providing the trust anchor represents the domain in question. Ultimately, the client still has to have a list of trust anchors configured—but they are DNSSEC trust anchors instead of PKIX trust anchors.

Of course, in principle, a client needs only one trust anchor for DNSSEC, the root zone trust anchor. Because control of the DNS root does not change very often, it makes sense for this trust anchor to be statically configured!

The ability of a domain operator to explicitly indicate a trust anchor for a domain is obviously very powerful. It may be tempting to ask whether this case is really the only use case that DANE needs, that is, whether the constraint cases mentioned previously are needed at all. The answer is that the constraint cases are useful as a way to fold in PKIX validation with external Certification Authorities in addition to domain-asserted trust anchors. Most obviously, this feature is useful in transition, when not all clients will be DANE-aware. But even in the longer term, it is possible that Certification Authorities will be able to provide added value over DANE. For example, while DANE is made to bind certificates to domain names, Certification Authorities can vouch for bindings of certificates to other things, such as the legal identity and physical location attested in Extended Validation certificates^[22].

Transition Challenges

As described previously, DANE offers some valuable new security properties for TLS authentication. But as with most IETF technologies—especially security technologies—there are some challenges to be overcome and some new potential pitfalls.

The most significant constraint for DANE deployment is DNSSEC deployment. On the server side, this problem is not a significant one because DNSSEC support is spreading fairly rapidly. On the client side, it may be more difficult. Although there are DNS libraries with robust DNSSEC support, many of the major DNS *Application Programming Interfaces* (APIs) that applications use do not provide any information about the DNSSEC status of the results returned.

So in order to implement DANE, application developers may have to re-factor their DNS support in addition to querying for some new record types. If more sites come to rely on DANE, then this process could also draw increasing attention to the various types of intermediaries that cause DNSSEC breakage (for example, home gateways that set DNS flags improperly).

Adding DNSSEC to the TLS connection process can also add significant latency to the TLS connection process. In addition to completing the TLS handshake and certificate validation, the client has to wait for several DNS round trips and then validate the chain of DNSSEC signatures. These combined delays can add up to multiple seconds of latency in connection establishment. Especially for real-time protocols such as HTTPS, *Session Initiation Protocol* (SIP), or *Extensible Messaging and Presence Protocol* (XMPP), such delay is clearly undesirable.

One mechanism proposed to mitigate these delays is to have the server pre-fetch all of the relevant DNSSEC records, namely all of the DS, DNSKEY, and RRSIG records chaining back to the root^[27]. Then the server can provide a serialized version of the DNSSEC records in the TLS handshake, saving the client the latency of the required DNS queries. The details of this mechanism, however, are still being worked out among the DANE, TLS, and PKIX working groups^[23]. A prototype version is now available in the Google Chrome web browser^[24].

Security Considerations

From a security perspective, the major effect of DANE is the new role that DNS operators will play in securing Internet applications. Although DNSSEC has always meant that DNS operators would have more security functions, DANE deployment will give them an explicit effect on application security, acting as arbiters of who can authenticate under a given name in TLS. Especially if services use trust anchor assertions, DNS operators will play an analogous role to the one Certification Authorities play today—a compromise in a DNS operator will allow an attacker to masquerade as a victim domain (albeit for a more limited set of domains because of DANE constraints on names). So DNS operators are likely to inherit many of the security troubles that Certification Authorities experience today and will need to strengthen their security posture accordingly.

Another more subtle risk arises from the fact that the operator of a DNS zone is not always the same as the entity that is authorized to control the contents of the zone, which we will call the “domain holder.” We used the phrase “domain operator” previously because DNSSEC protects DNS information only between the operator’s name server and the client—it does not say that what is provisioned in the name server is authorized by the domain holder.

When a domain is operated by a third party, that third party is a point of vulnerability between the client and the holder of the domain. If the domain operator provides false DANE information through malice or compromise, then a client will not be able to distinguish it from genuine DANE information. To some extent, this risk is not really new; because many current Certification Authorities authenticate requests for domain certificates based on information that is under the control of the domain operator, domain operators can already influence the credentialing process. With DANE, however, the vulnerability is much easier to exploit, for example, because the DNS operator does not have to trick a third party. This vulnerability is also fundamental to protocols that rely on DNSSEC for security, and the implications for DANE are discussed in detail in the DANE use cases document^[16]. The main mitigation is simply increased care on the part of domain holders to ensure that domain operators are not behaving badly.

Conclusions

For many years now, Internet applications have relied on assertions by third-party PKIX Certification Authorities to ensure that a server holding a particular private key was authorized to represent a domain. The promise of DANE is a more direct interaction between clients and the domains they interact with, secured by DNSSEC. In the short run, DANE can be deployed as an adjunct to the current system of certificates and authorities, adding constraints to better protect domains. In the long run, DANE will also allow domain operators to vouch for their own names.

The transition and security problems that face DANE are largely the growing pains of DNSSEC. It is not that DANE is causing these problems itself; rather, the problems arise because DANE is the first real application of DNSSEC that is expected to be widely deployed. So although it may be difficult to mitigate some of the security problems that DANE raises, and to enable more robust DNSSEC support in applications and gateways, these changes will ultimately make it simpler for applications to use DNSSEC for other purposes.

The DANE working group is making consistent progress on its deliverables, and there are already some prototype deployment tools. Their use cases document has been published as RFC 6394^[16], and the corresponding document defining the TLSA record type is starting to mature^[17]. As of this writing, it is in Working Group Last Call. On the client side, a variant of DANE has already been implemented in Google Chrome; on the server side, prototype tools are available to generate DANE records and to generate “DNSSEC-stapled” certificates based on DANE records^[24, 25]. There is also an early-stage command-line tool for generating and verifying TLSA records^[26].

References

- [0] Geoff Huston, “Hacking Away at Internet Security,” *The Internet Protocol Journal*, Volume 15, No. 1, March 2012.
- [1] Eric Rescorla, “HTTP Over TLS,” RFC 2818, May 2000.
- [2] Tim Dierks and Eric Rescorla, Editors, “The Transport Layer Security (TLS) Protocol Version 1.2,” RFC 5246, August 2008.
- [3] Chris Newman, “Using TLS with IMAP, POP3 and ACAP,” RFC 2595, June 1999.
- [4] Paul Hoffman, “SMTP Service Extension for Secure SMTP over Transport Layer Security,” RFC 3207, February 2002.
- [5] Jonathan Rosenberg, Henning Schulzrinne, Gonzalo Camarillo, Alan Johnston, Jon Peterson, Robert Sparks, Mark Handley, and Eve Schooler, “SIP: Session Initiation Protocol,” RFC 3261, June 2002.
- [6] Peter Saint-Andre, “Extensible Messaging and Presence Protocol (XMPP): Core,” RFC 6120, March 2011.
- [7] Paul Mockapetris, “Domain Names – Concepts and Facilities,” RFC 1034, November 1987.
- [8] Peter Saint-Andre and Jeff Hodges, “Representation and Verification of Domain-Based Application Service Identity within Internet Public Key Infrastructure Using X.509 (PKIX) Certificates in the Context of Transport Layer Security (TLS),” RFC 6125, March 2011.
- [9] Eric Rescorla and Nagendra Modadugu, “Datagram Transport Layer Security Version 1.2,” RFC 6347, January 2012.
- [10] David Cooper, Stefan Santesson, Stephen Farrell, Sharon Boeyen, Russell Housley, and Tim Polk, “Internet X.509 Public Key Infrastructure Certificate and Certificate Revocation List (CRL) Profile, RFC 5280, May 2008.
- [11] Roy Arends, Rob Austein, Matt Larson, Dan Massey, and Scott Rose, “DNS Security Introduction and Requirements,” RFC 4033, March 2005.
- [12] <https://dlv.isc.org/>
- [13] <http://arstechnica.com/apple/news/2011/09/safari-users-still-susceptible-to-attacks-using-fake-diginotar-certs.ars>

- [14] <http://blogs.comodo.com/it-security/data-security/the-recent-ra-compromise/>
- [15] http://www.theregister.co.uk/2011/09/06/diginotar_audit_damning_fail/
- [16] Richard Barnes, “Use Cases and Requirements for DNS-Based Authentication of Named Entities (DANE),” RFC 6394, October 2011.
- [17] Paul Hoffman and Jakob Schlyter, “Using Secure DNS to Associate Certificates with Domain Names for TLS,” Internet Draft, work in progress, **draft-ietf-dane-protocol-16**, February 2012.
- [18] <http://www.google.co.uk/support/forum/p/gmail/thread?tid=2da6158b094b225a&hl=en>
- [19] <http://www.eff.org/observatory>
- [20] <https://www.my.af.mil/>
- [21] Russ Housley, Sam Ashmore, and Carl Wallace, “Trust Anchor Management Protocol (TAMP),” RFC 5934, August 2010.
- [22] http://cabforum.org/Guidelines_v1_2.pdf
- [23] Adam Langley, “Serializing DNS Records with DNSSEC Authentication,” Internet Draft, work in progress, July 2011, **draft-agl-dane-serializechain**.
- [24] <http://www.imperialviolet.org/2011/06/16/dnssec-chrome.html>
- [25] <https://dane.xelerance.com/>
- [26] <https://github.com/pieterlexis/swede>
- [27] Wikipedia, “List of DNS record types,” http://en.wikipedia.org/wiki/List_of_DNS_record_types

RICHARD BARNES has been with BBN Technologies since 2005. Richard is a member of BBN’s Internet standards security team. In that role, he currently leads BBN’s IETF standards efforts in the areas of geolocation, presence, and emergency calling. He is chair of the IETF GEOPRIV working group and a member of the IETF Security Area Directorate (SECDIR), and he is one of the program chairs of the Emergency Services Workshop. Prior to joining BBN, he was a student at the University of Virginia (United States), from which he received a B.A. and M.S. in Mathematics, with research focused on biologically based neural networks, quantum informatics, and network security. E-mail: rbarnes@bbn.com

A Retrospective: Twenty-Five Years Ago

by Geoff Huston, APNIC

The Information Technology business is one that rarely pauses for breath. Gordon Moore noted in 1965 that the number of components in integrated circuits had doubled every year from 1958 to 1965, and confidently predicted that this doubling would continue “for at least 10 years.” This feature has been a continuing feature of the silicon industry for the past 50 years now, and its constancy has transformed this prediction into *Moore’s Law*. The implications of this constant impetus for innovation in this industry have resulted in an industry that is incapable of remaining in stasis, and what we have instead is an industry that completely reinvents itself in cycles as short as a decade.

Looking back over the past 25 years, we have traversed an enormous distance in terms of technical capability. The leading silicon innovations of the late 1980s were in the Intel 80486 chip, which contained 1 million transistors on a single silicon chip with a clock speed of 50 MHz, and a similarly capable Motorola 68040 processor. Twenty-five years later the state of the art is a multicore processor chip that contains just under 3 billion individual transistors and clock speeds approaching 4 GHz. And where has all that processing power gone? In the same period we have managed to build extremely sophisticated programmed environments that have produced such products as Apple’s *Siri* iPhone application, which combines voice recognition with a powerful information manipulation system, and we have packaged all of this computing capability into a device that fits comfortably in your pocket with room to spare!

Given that the last 25 years in IT has been so active, to look back over this period and contemplate all that has happened is a daunting task, and I am pretty sure that any effort to identify the innovative highlights in that period would necessarily be highly idiosyncratic. So instead of trying to plot the entire story that took us from then to now, I would like instead just to look at “then.” In this article, to celebrate 25 combined years of *The Internet Protocol Journal* (IPJ)^[2, 3] and its predecessor *ConneXions—The Interoperability Report*^[0], I would like to look at the networking environment of the late 1980s and see what, if anything, was around then that was formative in shaping what we are doing today, and how it might influence our tomorrow.

The Computing Landscape of the Late 1980s

The computing environment of the late 1980s now seems to be quite an alien environment. Obviously there were no pocket-sized computers then. Indeed there were no pocket-sized mobile phones then. (I recall a visit from a salesman at the time who sported the very latest in mobile telephony—a radio setup that was the size of a briefcase!)

In 1987 the IT world was still fixated with the mainframe computer, which was basking in its last couple of years of viability in the market. IBM enjoyed the dominant position in this marketplace, and *Digital Equipment Corporation* (DEC) was competing with IBM with its VAX/VMS systems. These systems were intended to take the place of the earlier DEC-10 architectures, as well as offering an upgrade path for the hugely successful PDP-11 minicomputer line. The typical architecture of the computing environment was still highly centralized, with a large multiuser system at its core, and an attendant network or peripheral devices. These peripheral devices were traditionally video terminals, which were a simple ASCII keyboard and screen, and the interaction with the mainframe was through simple serial line character-based protocols.

Although it may not have been universally accepted at the time, this period at the end of the 1980s marked the end of the custom-designed mainframe environment, where large-scale computer systems were designed as a set of component subsystems, placed into a rack of some sort and interconnected through a bus or blackplane. Like many other human efforts, as far as the mainframe computer sector was concerned its final achievements were its greatest.

While the mainframe sector was inexorably winding down, at the other end of the market things were moving very quickly. The Zylogics Z80 processor of the mid-1970s had been displaced by the Intel 8080 chip, which evolved rapidly into 16-bit, then 32-bit processor versions. By 1987 the latest chip was the Intel 80386, which could operate with a clock speed up to 33 MHz. The bus was 32 bits wide, and the chip supported a 32-bit address field. This chip contained some 275,000 transistors, and was perhaps the transformative chip that shifted the personal computer from the periphery of the IT environment to the mainstream. This chip took on the mainframe computer and won. The evolving architecture of the late 1980s was shifting from a central processing center and a cluster of basic peripheral devices to one of a cluster of personal desktop computers.

The desktop personal computer environment enabled computing power to be treated as an abundant commodity, and with the desktop computer came numerous interface systems that allowed users to treat their computer screens in a manner that was analogous to a desktop. Information was organized in ways that had a visual counterpart, and applications interacted with the users in ways that were strongly visual. The approach pioneered by the Xerox Star workstation in the late 1970s and brought to the consumer market through the Apple Lisa and Macintosh systems were then carried across into the emerging “mainstream” of the desktop environment with Windows 2.0 in the late 1980s.

The state of the art of portability was still in the category of “luggable” rather than truly portable, and the best example of what was around at the time is the ill-fated Macintosh Portable, which like its counterpart in the portable phone space was the size of a briefcase and incredibly heavy.

Oddly enough, while the industry press was in raptures when it was released in 1989, it was a complete failure in the consumer market. The age of the laptop was yet to come.

One major by-product in this shift in the computing environment to a distributed architecture was a major shift in the attention to networking, and at the same time as there was a large-scale shift in the industry from mainframes to personal computers, there were also numerous major changes in the networked environment.

The Networking Environment of the Late 1980s

A networking engineer in the late 1980s was probably highly conversant in how to network serial terminals to mainframes. The pin-outs in the DB-25 plug used by the RS-232 protocol was probably one of the basic ABCs of computer networking. At that time much of the conventional networked environment was concerned with connecting these terminal devices to mainframes, statistical multiplexors, and terminal switches, and serial switch suppliers such as Gandalf and Micom were still important in many large-scale computing environments.

At the same time, another networking technology was emerging—initially fostered by the need to couple high-end workstations with mainframes—and that was *Ethernet*. Compared to the kilobits per second typically obtained by running serial line protocols over twisted pairs of copper wires, the 10-Mbps throughput of Ethernet was blisteringly fast. In addition, Ethernet could span environments with a diameter of around 1500 meters, and with a certain amount of tweaking or with the judicious use of Ethernet bridges and fibre-optic repeaters this distance could be stretched out to 10 km or more.

Ethernet heralded a major change in the networked environment. No longer were networks hub-and-spoke affairs with the mainframe system at the center. Ethernet supplied a common bus architecture that supported any-to-any communications. Ethernet was also an open standard, and many vendors were producing equipment with Ethernet interfaces. In theory, these interfaces all interoperated, at least at the level of passing Ethernet frames across the network (aside from a rather nasty incompatibility between this original Digital-Intel-Xerox specification and the IEEE 802.3 “standardized” specification!).

However, above the basic data framing protocol the networked environment was still somewhat chaotic. I recall the early versions of the multiprotocol routers produced by Proteon and Cisco supported more than 20 networking protocols! There was *DECnet*, a proprietary network protocol suite from the Digital Equipment Corporation, which at around 1987 had just released Phase IV, and was looking toward a Phase V release that was to interoperate with the International Organization for Standardization’s *Open Systems Interconnection* (OSI) protocol suite^[1] (more on this subject a bit later).

There was IBM's *Systems Network Architecture* (SNA), which was a hierarchical network that supported a generic architecture of remote job entry systems clustered around a central service mainframe. There was the *Xerox Network Services* (XNS) protocol used by Xerox workstations. Then there were Apollo's *Network Computing Architecture* (NCA) and Apple's *AppleTalk*. And also in this protocol mix was the *Transmission Control Protocol/Internet Protocol* (TCP/IP) protocol suite, used at that time predominately on UNIX systems, although implementations of TCP/IP for Digital's VAX/VMS system were very popular at the time. A campus Ethernet network of the late 1980s would probably see all of these protocols, and more, being used concurrently.

And there was the ISO-OSI protocol suite, which existed more as a future protocol suite than as a working reality at the time.

The ISO-OSI and TCP/IP protocol suites were somewhat different from the others that were around at the time because both were deliberate efforts to answer a growing need for a vendor-independent networking solution. At the time the IT environment was undergoing a transition from the monoculture of a single vendor's comprehensive IT environment—which bundled the hardware of the mainframe, network, peripherals, terminals, and the software of the operating system, applications, and network all into the one bundle—into a piecemeal environment that included a diverse collection of personal workstations, desktop computers, peripherals, and various larger minicomputers and mainframe computers in one environment. What was needed was a networking technology that was universally supported on all these various IT assets. What we had instead was a more piecemeal environment. Yes, it was possible to connect most of these systems into a common Ethernet substrate, but making A talk to B was still a challenge, and various forms of protocol translation units were also quite commonplace at the time. What the industry needed was a vendor-independent networking protocol, and there were two major contenders for this role.

ISO-OSI and TCP/IP

The ISO-OSI protocol suite was first aired in 1980. It was intended to be an all-embracing protocol suite that embraced both the IEEE 802.3 Ethernet protocols and the X.25 packet switching protocols that were favoured by many telephony operators as their preferred wide-area data services solution. The ISO-OSI network layer included many approaches, including the telephony sector's *Integrated Service Digital Network* (ISDN), a *Connection-Oriented Network Service* (CONS), a virtual circuit function based largely on X.75 that was essentially the “call-connection” function for X.25, and a *Connectionless Network Service* (CLNS), based loosely on the IP protocol with the use of the *End System-to-Intermediate System Routing Exchange Protocol* (ES-IS) routing protocol.

Above the network layer were numerous end-to-end transport protocols, notably *Transport Protocol Class 4* (TP4), a reliable connection-oriented transport service, and *Transport Protocol Class 0* (TP0), a connectionless packet datagram service. Above this layer was a Session Layer, X.215, used by the TP4 CONS services, and a Presentation Layer, defined using the *Abstract Syntax Notation One* (ASN.1) syntax.

ISO-OSI included numerous application-level services, including *Virtual Terminal Protocol* (VTP) for virtual terminal support, *File Transfer Access And Management* (FTAM) for file transfer, *Job Transfer And Management* (JTAM) for batch job submission, *Message Handling System* (MHS, also known as X.400) for electronic mail, and the X.500 Directory service. ISO-OSI also included a *Common Management Information Protocol* (CMIP). ISO-OSI attempted to be everything to everybody, as evidenced by the “kitchen sink” approach adopted by many of the OSI standardization committees at the time.

When confronted by many technology choices, the committees apparently avoided making a critical decision by incorporating both approaches into the standard. The most critical decision in this protocol suite was the inclusion of both connection-oriented and connectionless networking protocols. They also used session and presentation layer protocols, whose precise role was a mystery to many! ISO-OSI was a work-in-progress at the time, and the backing of the telephone sector, coupled with the support of numerous major IT vendors, gave this protocol an aura of inevitability within the industry. Whatever else was going to happen, there was the confident expectation that the 1990s would see all computer networks move inevitably to use the ISO-OSI protocol suite as a common, open, vendor-neutral network substrate.

If the ISO-OSI had a mantra of inevitably, the other open protocol suite of the day, the TCP/IP protocol suite, actively disclaimed any such future ambitions. TCP/IP was thought of at the time as an experiment in networking protocol design and architecture that ultimately would go the way of all other experiments, and be discarded in favor of a larger and more deliberately engineered approach. Compared to the ISO-OSI protocols, TCP/IP was extremely “minimalist” in its approach. Perhaps the most radical element in its design was to eschew the conventional approach at the time of building the network upon a reliable data link protocol. For example, in DECnet Phase IV, the data link protocol, *Digital Data Communications Message Protocol* (DDCMP), performed packet integrity checks and flow control at the data link level. TCP/IP gracefully avoided this problem by allowing packets to be silently dropped by intermediate data switches, or corrupted while in flight. It did not even stipulate that successive packets within the same end-to-end conversation follow identical paths through the network.

Thus the packet switching role was radically simplified because now the packet switch did not need to hold a copy of transmitted packets, nor did it need to operate a complex data link protocol to track packet transmission integrity and packet flow control. When a switch received a packet, it forwarded the packet based on a simple lookup of the destination address contained in the packet into a locally managed forwarding table. Or it discarded the packet.

The second radical simplification in TCP/IP was the use of real-time packet *fragmentation*. Previously, digital networks were constructed in a “vertically integrated” manner, where the properties of the lower layers were crafted to meet the intended application of the network. Little wonder that the telephone industry put its support behind X.25, which was a reliable unsynchronized digital stream protocol. If you wanted low levels of jitter, you used a network with smaller packet sizes, whereas higher packet sizes improved the carriage efficiency. Ethernet attempted to meet this wide variance in an agnostic fashion by allowing packets of between 64 and 1500 octets, but even so there were critics who said that for remote terminal access the smallest packets were too large, and for large-scale bulk data movement the largest packets were too small. *Fiber Distributed Data Interface* (FDDI), the 100-Mbps packet ring that was emerging at the time as the “next thing” as commodity high-speed networking used a maximum size of 4000 octets packets in an effort to improve carriage efficiency, whereas the *Asynchronous Transfer Mode* (ATM) committee tried to throw a single-packet-size dart at the design board and managed to get the rather odd value of 53 octets!

IP addressed this problem by trying to avoid it completely. Packets could be up to 64,000 octets long, and if a packet switch attempted to force a large packet through an interface that could not accept it, the switch was allowed to divide the packet into appropriately sized autonomous fragments. The fragments were not reassembled in real time: that was the role of the ultimate receiver of the packets.

As an exercise in protocol design, IP certainly showed the elegance of restraint. IP assumed so little in terms of the transmission properties of the underlying networks that every packet was indeed an adventure! But IP was not meant to be the protocol to support the prolific world of communicating silicon in the coming years. This protocol and the IP networks that were emerging in the late 1980s were intended to be experiments in networking. There was a common view that the lessons learned with experience of operating high-speed local networks and wide-area networks using the TCP/IP protocol suite would inform the larger industry efforts. The inclusion of IP-based technologies in the ISO-OSI protocol suite^[4] was a visible instantiation of this proposed evolutionary approach.

While these two protocol suites vied with each other for industry attention at the time, there was one critical difference: It was a popular story at the time that the ISO-OSI protocol suite was a stack of paper some 6 feet high, which cost many hundreds of dollars to obtain, with no fully functional implementations, whereas the TCP/IP protocol suite was an open-sourced and openly available free software suite without any documentation at all. Many a jibe at the time characterized the ponderous approach of the ISO-OSI approach as “vapourware about paperware,” while the IP effort, which was forming around the newly formed *Internet Engineering Task Force* (IETF), proclaimed itself to work on the principle of “rough consensus and running code.”

Local- and Wide-Area Networking

The rise of Ethernet networks on campuses and in the corporate world in the late 1980s also brought into stark visibility the distinction between local- and wide-area networking.

In the local-area network, Ethernet created a new environment of “seamless connectivity.” Any device on the network could provide services to any other device, and the common asset of a 10-Mbps network opened up a whole new set of computing possibilities. Data storage could be thought of as a networked resource, so desktop computers could access a common storage area and complement it with local storage, and do so in a way that the distinction between local resources and shared networkwide resources was generally invisible. The rich computing environment of visualizing the application, popularized by both the Macintosh and Windows 2.0, complemented a rich networked environment where rather than bringing a user into the location that had both the data and the computing resources, the model was invested, and the user was able to exclusively use the local environment and access the remote shared resources through networking capabilities integrated into the application environment. Local-area networking was now an abundant resource, and the industry wasted no time on exploiting this new-found capability.

But as soon as you wanted to venture further than your *Local-Area Network* (LAN), the picture changed dramatically. The wide-area networking world was provisioned on the margins of oversupply of the voice industry, and the services offered reflected the underlying substrate of a digital voice circuit. The basic unit of a voice circuit was a 64-kbps channel, which was “groomed” into a digital circuit of either 56 or 48 kbps, depending on the particular technology approach used by the voice carrier. Higher capacities (such as 256 or 512 kbps) were obtained by multiplexing individual circuits together. Even high-capacity circuits were obtained by using a voice trunk circuit, which was either 1.5 (T1) or 2.048 Mbps (E1), again depending on the digital technology used by the voice carrier. Whereas the LANs were now supporting an any-to-any mode of connection, these *Wide-Area Networks* (WANs) were constructed using point-to-point technologies that were either statically provisioned or implemented as a form of “on-demand” virtual circuit (X.25).

In the late 1980s users' patience was running thin over having to use an entirely different protocol suite for the wide area as distinct from the local area. Often the wide area required the use of different applications with different naming and addressing conventions. One approach used by many Ethernet switch vendors was to introduce the concept of an *Ethernet Serial Bridge*. This technology allowed a logical IEEE 802.3 Ethernet to encompass much larger geographic domains, but at the same time protocols that worked extremely efficiently in the local area encountered significant problems when passed through such supposedly "transparent" Ethernet serial bridges.

However, these bridge units allowed significantly larger and more complex networks to be built using Ethernet as the substrate. The Ethernet *Spanning Tree Algorithm* gained traction in order to allow arbitrary topologies of interconnected LANs to self-organize into coherent topologies that eliminated loops and allowed for failover resilience in the network.

What has changed, and what has stayed the same?

So what have we learned from this time?

In the intervening period ISO-OSI waned and eventually disappeared, without ever having enjoyed widespread deployment and use. Its legacy exists in numerous technologies, including the X.500 Directory Service, which is the substrate for today's *Lightweight Directory Access Protocol* (LDAP) Directory Services. Perhaps the most enduring legacy of the ISO-OSI work is the use of the "layered stack" conceptual model of network architectures. These days we refer to "Layer 2 Virtual LANs (VLANs)" and "Layer 3 Virtual Private Networks (VPNs)" perhaps without appreciating the innate reference to this layered stack model.

Of course the ISO-OSI protocol suite was not the only casualty of time. DECnet is now effectively an historic protocol, and Novell's *NetWare* has also shifted out of the mainstream of networking protocols. Perhaps it may be more instructive to look at those technologies that existed at the time that have persisted and flourished so that they now sit in the mainstream of today's networked world.

Ethernet has persisted, but today's Ethernet networks share little with the technology of the original IEEE 802.3 *Carrier Sense Multiple Access with Collision Detection* (CSMA/CD) 10-Mbps common bus network. The entire common bus architecture has been replaced by switched networks, and the notion of self-clocking packets was discarded when we moved into supporting Gbps Ethernets. What has persisted is the IEEE 802.3 packet frame format, and the persistence of the 1500-octet packet as the now universal lowest common factor for packet quantization on today's network. Why did Ethernet survive while other framing formats, such as *High-Level Data Link Control* (HDLC), did not?

I could suggest that it was a triumph of open standards, but HDLC was also an open standard. I would like to think that the use of a massive address space in the Ethernet frame, the 48-bit *Media Access Control* (MAC) address, and the use since its inception of a MAC address registry that attempted to ensure the uniqueness of each Ethernet device were the most critical elements of the longevity of Ethernet.

Indeed not only has UNIX persisted, it has proliferated to the extent that it is ubiquitous, because it now forms the foundation of the Apple and Android products. Of the plethora of operating systems that still existed in the late 1980s, it appears that all that have survived are UNIX and Windows, although it is unclear how much of Windows 2.0 still exists in today's Windows 7, if anything.

And perhaps surprisingly TCP/IP has persisted. For a protocol that was designed in the late 1970s, in a world where megabits per second was considered to be extremely high speed, and for a protocol that was ostensibly experimental, TCP/IP has proved to be extremely persistent. Why? One clue is in the restrained design of the protocol, where, as we have noted, TCP/IP did not attempt to solve every problem or attempt to be all things for all possible applications. I suspect that there are two other aspects of TCP/IP design that contributed to its longevity.

The first was a deliberate approach of modularity in design. TCP/IP deliberately pushed large modules of functions into distinct subsystems, which evolved along distinct paths. The routing protocols we use today have evolved along their own paths. Also the name space and the mapping system to support name resolution has evolved along its own path. Perhaps even more surprisingly, we have had the rate control algorithms used by TCP, the workhorse of the protocol suite, evolve along its own path.

The second aspect is use of what was at the time a massively sized 32-bit address space, and an associated address registry that allowed each network to use its own unique address space. Like the Ethernet 48-bit MAC address registry, the IP address registry was, in my view, a critical and unique aspect of the TCP/IP protocol suite.

Failures

What can we learn from the various failures and misadventures we have experienced along the way?

Asynchronous Transfer Mode (ATM) was a technology that despite considerable interest from the telephone operators proved to be too little too late, and was ultimately swept aside in the quest for ever larger and ever cheaper network transmission systems. ATM appeared to me to be perhaps the last significant effort to invest value into the network through allowing the network to adapt to the various differing characteristics of applications.

The underlying assumption behind this form of adaptive networking is that attached devices are simply incapable of understanding and adapting to the current state of the network, and it is up to the network to contain sufficient richness of capability to present consistent characteristics to each application. However, our experience has been quite the opposite, where the attached devices are increasingly capable of undertaking the entire role of service management, and complex adaptive networks are increasingly seen as at best meaningless duplication of functions, and at worst as an anomalous network behavior that the end device needs to work around. So ATM failed to resonate with the world of data networking, and as a technology it has waned. In the same way subsequent efforts to equip IP networks with *Quality of Service* (QoS) responses, or the much-hyped more recent *Next-Generation Networking* (NGN) networking efforts have been failures, for much the same basic reasons.

Fiber Distributed Data Interface (FDDI) also came and went. Rings are notoriously difficult to engineer, particularly in terms of managing a coherent clock across all attached devices that preserves the circumference of the ring, as measured in bits on the wire. From its earlier lower-speed antecedents in the 4-Mbps token, the 100-Mbps FDDI ring attracted considerable interest in the early 1990s. However, it was in effect a dead end in terms of longer-term evolution—the efforts to increase the clock speed required either the physical diameter of the ring to shrink to unusable small diameters or the clock signal to be locked at extraordinarily high levels of stability that made the cost of the network prohibitive. This industry appears to have a strong desire for absolute simplicity in its networks, and even rings have proved to be a case of making the networks too complex.

Interestingly, and despite all the evidence in their favor, the industry is still undecided about open technologies. TCP/IP, UNIX, and the Apache web platform are all in their own way significant and highly persuasive testaments to the power of open-source technologies in this industry, and a wide panoply of open technologies forms the entire foundation of today's networked environment. Yet, in spite of all this accumulated experience, we still see major efforts to promote closed, vendor-specific technologies into the marketplace. Skype is a case in point, and it is possible to see the iPhone and the Kindle in a similar light, where critical parts of the technology are deliberately obscured and aspects of the device behavior are deliberately sealed up or occluded from third-party interception.

The Next Twenty-Five Years

In wondering about the next 25 years, it may be interesting to look back ever further, to the early 1960s, and see what, if anything, has proved to be enduring from the perspective of the past 50 years. Interestingly, it appears that very little of that time, except for the annoying persistence of Fortran, and the ASCII keyboard as the ubiquitous input device, is still a part of today's networked environment. So over a 50-year time period much has changed in our environment.

But, interestingly, when we par down the period to the past 25 years, there is still much that has survived in the computing and networking environment. A Macintosh computer of the late 1980s looks eerily familiar, and although today's systems are faster, lighter, and a lot less clunky, there is actually very little that has changed in terms of the basic interface with the user. A Macintosh of that time could be connected to an Ethernet network, and it supported TCP/IP, and I suspect that if one were to resurrect a Mac system from 1988 loaded with *MacTCP* and connect it to the Internet today it would be frustratingly, achingly slow, but I would like to think that it would still work! And the applications that ran on that device have counterparts today that continue to use the same mechanisms of interaction with the user.

So if much of today's world was visible 25 years ago, then where are the aspects of change? Are we just touching up the fine-point details of a collection of very well established technologies? Or are there some basic and quite fundamental shifts underway in our environment?

It seems to me that the biggest change is typified in today's tablet and mobile phone computers, and in these devices it is evident that the metaphors of computing and interaction with applications are changing. The promise from 1968 in the film *2001: A Space Odyssey* of a computer that was able to converse with humans is now, finally, within reach of commodity computing and consumer products. But it is more than merely the novelty of a computer that can "talk." The constant search for computing devices that are smaller and more ubiquitous now means that the old paradigm of a computer as a "clever" but ultimately bulky typewriter is fading away. Today we are seeing modes of interaction that use gestures and voice, so that the form factor of a computer can become smaller while still supporting a functional and efficient form of interaction with the human user.

It is also evident that the pendulum of distribution and centralization of computing capability is swinging back, and the rise of the heavily hyped *Cloud*^[5, 6] with its attendant collection of data centers and content distribution networks, and the simultaneous shrinking of the end device back to a "terminal" that allows the user to interact with views into a larger centrally managed data store held in this cloud, appears to be back in vogue once more.

It is an open question whether these aspects of today's environment will be a powerful and persistent theme for the next 25 years, or whether we will see other aspects of our environment seize industry momentum, so they are very much just a couple of personal guesses. Moore's Law has proved to be truly prodigious over the past 50 years. It has allowed us to pack what was a truly unbelievable computing capability and storage into astonishingly small packages and then launch them into the consumer market with pricing each year that appears to be consistently lower than the previous year.

If this property of packaging ever greater numbers of transistors into silicon chips continues for the next 25 years at the same rate, then it is likely that whatever happens in the next 25 years, the only limitation may well be our imagination rather than any intrinsic limitations of the technology itself.

For Further Reading

- [0] The Charles Babbage Institute at the University of Minnesota has scanned the complete collection of *ConneXions—The Interoperability Report*, and it is available at this URL:
<http://www.cbi.umn.edu/hostedpublications/Connexions/index.html>
- [1] Starting in April 1989 (Volume 3, No. 4), *ConneXions* published a long-running series of articles under the general heading “Components of OSI,” which described almost every aspect of this protocol suite. The same journal also published articles on many of the other technologies mentioned in this article, including FDDI, AppleTalk, and ATM.
- [2] Vint Cerf, “A Decade of Internet Evolution,” *The Internet Protocol Journal*, Volume 11, No. 2, June 2008.
- [3] Geoff Huston, “A Decade in the Life of the Internet,” *The Internet Protocol Journal*, Volume 11, No. 2, June 2008.
- [4] International Organization for Standardization, “Final text of DIS 8473, Protocol for Providing the Connectionless-mode Network Service,” RFC 994, March 1986.
- [5] T. Sridhar, “Cloud Computing—A Primer Part 1: Models and Technologies,” *The Internet Protocol Journal*, Volume 12, No. 3, September 2009.
- [6] T. Sridhar, “Cloud Computing—A Primer Part 2: Infrastructure and Implementation Topics,” *The Internet Protocol Journal*, Volume 12, No. 4, December 2009.

GEOFF HUSTON B.Sc., M.Sc., is the Chief Scientist at *Asia Pacific Network Information Centre* (APNIC), the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. He is author of numerous Internet-related books, was a member of the Internet Architecture Board from 1999 until 2005, and served on the Board of Trustees of the Internet Society from 1992 until 2001. E-mail: gih@apnic.net

Letter to the Editor

Dear Editor,

Who knew? Twenty-five years ago I started a tiny company that grew into Interop to spread the technical word about this funny thing we called *The Internet* and this really obscure thing called “TCP/IP.” Back in the ’70s, when the basic protocols were being created and experimented with, you were a high school kid in Norway and I was running a tiny group at SRI International and I let you use my machine across the ocean by using the ARPANET, the precursor to the Internet, because you seemed both smart and polite. Fifteen years later I decided to hire you to start a newsletter, *ConneXions*—*The Interoperability Report*, about the burgeoning Internet because of those properties and the perceived need to communicate monthly about the ins and outs of these simple but far-reaching technical protocols. You had the technical knowledge and good sense to enlist the brains of the real engineers in the field with real experience to further the knowledge of “all things Internet.”

Who knew this would be still going on 25 years later? Your combination of passion and patience has produced an amazing record of ongoing expertise for the whole world to enjoy.

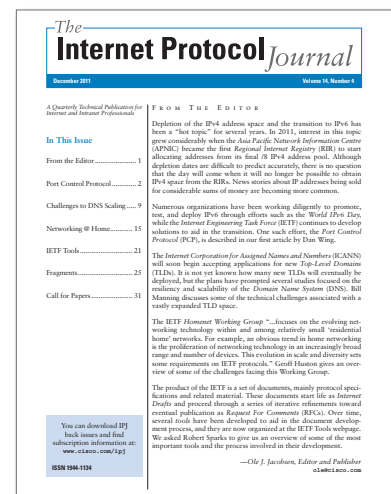
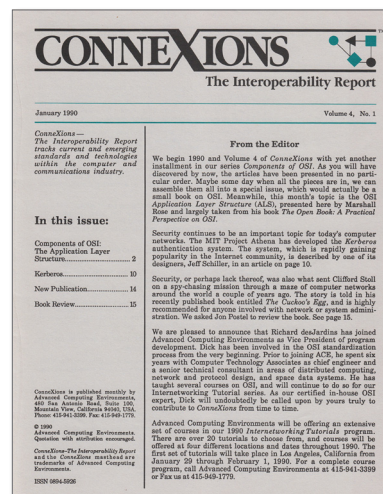
Thank you for being Ole!

—Dan Lynch
Founder of Interop
dan@lynch.com

Thank you, Dan!

I appreciate your very kind words. I also want to take this opportunity to thank all of the contributors to this journal. We could not do this without you!

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com



NIXI to Run New NIR in India

March saw the launch of a new *National Internet Registry* (NIR) for India, following the successful conclusion of talks between the *Asia Pacific Network Information Centre* (APNIC) and the Government of India.

The *Indian Registry For Internet Names And Numbers* (IRINN) will be run by the *National Internet Exchange of India* (NIXI) and serve ISPs within the country that wish to sign up. It is the result of a long collaboration between APNIC and NIXI, with APNIC staff sharing their expertise with NIXI, and NIXI officials putting together an impressive technical installation in preparation for the launch. The new registry was announced on the final day of APNIC 33, a technical conference conducted in conjunction with the *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT 2012).

APNIC Executive Council Chairman, Akinori Maemura said of the announcement, “We are extremely happy that this process is heading towards a positive conclusion; which, on the other hand, is also a commencement of a new relationship. I would like to thank the NIXI team for their support and the hard work they have demonstrated in making this a reality.”

Director General of APNIC, Paul Wilson commented, “We welcome the new National Internet Registry in India to the APNIC community. The Internet is a global community and IRINN, as the NIR is being called, should be part of that. I hope that many new Internet Services Providers will be formed in India, and they will always be able to choose between IRINN and APNIC for IP addresses. The market here is big enough and that kind of diversity will ensure better services and lower prices for all Indians.”

APNIC has over 300 members locally, mostly Internet Services Providers and Telecommunication Communications companies, and over 6 million *Internet Protocol version 4* (IPv4) addresses were allocated in 2011. There are already 6 National Internet Registries in Asia in South Korea: (KISA KRNIC), Japan (JPNIC), China (CNNIC), Indonesia (IDNIC), Vietnam (VNNIC) and Taiwan (TWINIC). This is out of 56 economies in the Asia Pacific region.

“It’s really about what is a better fit for the individual organizations. Typically we tend to see larger organizations prefer a regional service, especially those who operate in multiple economies to maintain an account with APNIC,” said Paul Wilson.

NIXI is a not-for-profit organization, set up for peering of ISPs among themselves for the purpose of routing domestic traffic within India, instead of routing it through international peering points, thereby resulting in reduced latency and reduced bandwidth charges for ISPs. NIXI is managed and operated on a neutral basis, in line with the best practices for such initiatives globally.

Internet Hall of Fame Advisory Board Named

The *Internet Society* (ISOC) recently announced that in conjunction with its 20th anniversary celebration, it is establishing an annual *Internet Hall of Fame* program to honor leaders and luminaries who have made significant contributions to the development and advancement of the global Internet.

Inaugural inductees will be announced at an Awards Gala during the ISOC's *Global INET 2012* conference in Geneva, Switzerland, April 22–24, 2012, www.internetsociety.org/globalinet

“There are extraordinary people around the world who have helped to make the Internet a global platform for innovation and communication, spurring economic development and social progress,” noted ISOC CEO Lynn St. Amour. “This program will honor individuals who have pushed the boundaries to bring the benefits of a global Internet to life and to make it an essential resource used by billions of people. We look forward to recognizing the achievements of these outstanding leaders.”

ISOC has convened an Advisory Board to vote on the inductees for the 2012 Internet Hall of Fame inauguration. The Advisory Board is a highly-qualified, diverse, international committee that spans multiple industry segments and backgrounds. This year's Advisory Board members include:

- Dr. Lishan Adam, ICT Development Researcher, Ethiopia
- Chris Anderson, Editor-in-Chief, *WIRED Magazine*
- Alex Corenthin, Directeur des Systemes d'Information, University Cheikh Anta Diop of Dakar/Chair, Internet Society Senegal Chapter
- William Dutton, Professor of Internet Studies, Oxford Internet Institute
- Joichi Ito, Director, MIT Media Lab
- Mike Jensen, Independent ICT Consultant, South Africa
- Aleks Krotoski, Technology Academic/Journalist/Author
- Loic Le Meur, Founder & CEO, LeWeb
- Mark Mahaney, Internet Analyst, Citigroup
- Dr. Alejandro Pisanty, Professor at National University of Mexico/Chair of Internet Society Mexico Chapter
- Lee Rainie, Director, Pew Research Center's Internet & American Life Project
- Jimmy Wales, Co-founder, Wikipedia

“We are extremely grateful to our distinguished Advisory Board members who have donated their time, energy, and expertise to this program,” St. Amour added. “The breadth of their experiences and the diversity of their perspectives are invaluable, and we truly appreciate their participation.”

The Internet Society is the trusted independent source for Internet information and thought leadership from around the world. With its principled vision and substantial technological foundation, the Internet Society promotes open dialogue on Internet policy, technology, and future development among users, companies, governments, and foundations. Working with its members and Chapters around the world, the Internet Society enables the continued evolution and growth of the Internet for everyone.

For more information, see: <http://www.internetsociety.org>

IETF Journal Now Available by Subscription

The *IETF Journal* provides anyone with an interest in Internet standards an overview of the topics being debated by the *Internet Engineering Task Force* (IETF), and also helps facilitate participation in IETF activities for newcomers.

The *IETF Journal* aims to provide an easily understandable overview of what is happening in the world of Internet standards, with a particular focus on the activities of the IETF Working Groups. Each issue highlights hot issues being discussed in IETF meetings and on the IETF mailing lists.

Visit *The IETF Journal* on the Web at www.internetsociety.org/ietfjournal to see the latest edition, or to subscribe to the e-mail edition or have it delivered as a hardcopy, visit:

<http://www.internetsociety.org/ietfjournal-subscribe>

The *IETF Journal* is an Internet Society publication produced in cooperation with the IETF.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2012 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

June 2012

Volume 15, Number 2

A Quarterly Technical Publication for
Internet and Intranet Professionals

In This Issue

From the Editor	1
Transition Space	2
December in Dubai.....	17
IP Fast Reroute	30
Letters to the Editor.....	35
Call for Papers.....	39

FROM THE EDITOR

Deployment of IPv6 took another step forward on June 6, 2012, when numerous website operators, *Internet Service Providers* (ISPs), and home router vendors participated in the *World IPv6 Launch*. Organized by the Internet Society, the event attracted significant media attention as the participants enabled IPv6 permanently and rendered it “on by default.” More information about the event is available from www.worldipv6launch.org

Migration to IPv6 is not a simple task, as outlined in many previous editions of this journal. Various tools and techniques have been developed, one being the use of so-called *Carrier-Grade NATs* whereby the end customers connect to the Internet using private (RFC 1918) addresses and the ISP provides translation for both public IPv4 and IPv6 addresses. In April of this year, the *Internet Engineering Task Force* (IETF) approved and the *Internet Assigned Numbers Authority* (IANA) allocated a new IPv4 address block (100.64.0.0/10), designated for use as “Shared Transition Space” in support of the IPv6 transition. We asked Wesley George to describe the rationale behind the use of this additional private address space and discuss the debate that resulted from this allocation.

The world of telecommunications has changed dramatically as a result of the rapid expansion of the Internet. Traditional telephone lines are being replaced by *Voice over IP* (VoIP) systems for both private and business use. These changes represent big challenges for traditional telephone carriers, and even for some countries whose income used to depend largely on telephone “settlement charges” for international phone calls. The *World Conference on International Telecommunications* (WCIT) will take place this coming December in Dubai. Geoff Huston discusses some of the proposed changes to the *International Telecommunication Regulations* that could affect the Internet in various ways and will be discussed at WCIT.

The IETF is concerned not only with IPv4-to-IPv6 migration, but also with recovery upon router or link failure. In our final article, Russ White describes *IP Fast Reroute*, a technique for providing fast traffic recovery when these failures occur.

As always, your feedback about anything you read in this journal is most appreciated. Please contact us at ipj@cisco.com and don't forget to renew your subscription and provide us with any postal or e-mail changes.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

Shared Transition Space: Is it necessary?

by Wesley George, Time Warner Cable

Recently, the *Internet Engineering Task Force* (IETF) approved^[1] and the *Internet Assigned Numbers Authority* (IANA) allocated^[2] a new IPv4 address block (100.64.0.0/10) designated for use as “Shared Transition Space” in support of the IPv6 transition. This decision was highly controversial within the different standards and policy bodies that discussed the idea. The author would like to note that people have been debating this topic for years, and nearly everyone within the broad stakeholder community seems to have a strong opinion on the matter, including me. Despite the best of intentions, some of my opinions and biases may appear within the article. I did not intend this article to be a definitive conclusion on the matter, but rather a summary of the recent discussion. Whether the standards bodies involved came to the “right” or “wrong” conclusion—as well as the veracity of the arguments on both sides—is an exercise for you, the reader.

Internet Service Providers (ISPs) and users have significant investments in equipment and applications that must be updated to support IPv6. Progress is accelerating with regard to IPv6 availability in hardware, software, and access, though broad availability remains a long-term problem. In the interim, IPv4 will continue to be an important capability for providing users with access to Internet resources. As a consequence, considerable effort has been expended in conserving the increasingly scarce IPv4 resources while maintaining “business as usual.” This conservation has taken the form of policies for address allocation and management^[3], as well as new protocols and technologies. It is likewise important to note that ISPs must manage IPv4 exhaustion in a way that is least disruptive to users while undertaking full IPv6 deployment—two completely different and parallel activities. Any business that relies entirely on efforts to extend the useful life of IPv4 without executing on an IPv6 deployment plan is merely delaying the inevitable effects on their customers and ultimately their profitability.

IPv4 “life extension” is an area that remains controversial. Some believe that any effort to extend the useful life of IPv4 and allow the IPv4 Internet to keep growing beyond its original design limitations will seriously affect the timeliness of reaching critical mass with IPv6. The idea that many opponents of the “life-extensions” methods are supporting is that IPv4 exhaustion and the resulting transition from IPv4 to IPv6 is going to be disruptive to customers and operations no matter when it actually occurs. From this perspective it is preferable to have a brief—but significant—disruption and transition completely to IPv6. This plan is akin to the idea that it is better to just rip the bandage off and have a moment of pain than removing it slowly in an attempt to reduce the pain.

The counterpoint to this argument is that we must look at the situation pragmatically with the goal of maintaining business continuity, growth, and customer satisfaction.

IPv4 Exhaustion

The impending IPv4 address exhaustion^[4] and the problems it will create has been the topic of much discussion in many different areas of the Internet community. The need to deploy IPv6 has figured prominently in the discussion, because it is the proper long-term solution. However, the unfortunate reality is that deploying IPv6 is a parallel activity to any work that provides continuity to the existing IPv4 network in order to keep it operational and able to grow to meet demand. As an Internet community, we are not where we need to be in terms of critical mass of our IPv6 deployments, in terms of either available, deployed equipment that supports IPv6 fully or applications that are able to use IPv6 when it is available.

IPv6 deployment is a requirement, but most ISPs do not have control over all variables affecting IPv6 deployment, and they have limited influence on progress outside of their network boundaries. This reality is especially true with residential services, where customers often purchase IP-enabled hardware directly from retailers to connect to their home networks. Consumers generally do not care about whether a device supports IPv4 or IPv6, so they do not make purchasing decisions based on such features. Customers should not be required to be technology experts in order to get their devices to work properly for their intended use. Customers generally are not interested in their ISP dictating the equipment that they may use in their home, and they do not like being told that they must replace “obsolete” gear, especially if they purchased it recently. The service provider sells “Internet” service, so customers expect their “Internet” devices to work—period. As a result, if an ISP wants to continue to grow, that ISP must continue to offer IPv4 services until the existing equipment without IPv6 support ages out of the network and is replaced.

The IETF recently released a *Best Current Practice* (BCP) document^[5] that provides some guidance for implementers that support for IPv6 on “IP-capable” devices is going to be a necessity, and the *Consumer Electronics Association* (CEA) now has a working group on IPv6 Transition^[6]. In conjunction with events like *World IPv6 Launch*^[7], there are near-constant improvements in the availability of IPv6-capable hardware, software, access, and services. The result of this situation should be that critical mass of IPv6 deployment will happen soon and reduce reliance on IPv4 and IPv4 life-extension technologies.

Because of the costs, operational complexities, performance concerns, and effects on customers that most IPv4 life-extension technologies create, service providers should focus on reaching IPv6 critical mass in essential areas.

When IPv6 has become sufficiently ubiquitous, the need for IPv4 life-extension technologies will be reduced along with the scale of deployments. Because a lot of the costs of deploying IPv4 life-extension technologies are initial costs, there is some truth to the argument that after they are deployed they are unlikely to disappear anytime soon. Why would a carrier invest significant time and money in deploying something only to pull it back out a short time later? Therefore the best method to reduce the cost of *Carrier-Grade NAT* (CGN) deployment is to work to deploy less of it.

ISPs are different when it comes to their expectations for growth, and their IPv4 addressing reserves or consumption rates differ accordingly. Some have areas of their internal network where they can make changes and reclaim globally unique IPv4 addresses for reuse to support customers, some have addresses that can be reclaimed via auditing and improved efficiency of allocation, and still others have already undertaken many of these projects and do not have much address space left to reclaim. Further, new IPv4 address availability as a combination of policies and demand may be different for each *Regional Internet Registry* (RIR). To summarize, the need for IPv4 address life-extension technologies is different on each network. The costs of deploying, the complexity of supporting, and the growth rate all figure into how widely service providers will have to deploy one or more technologies to extend their remaining IPv4 resources.

NAT444

Network Address Translation (NAT)^[28, 30] is already widely used for translating one IPv4 address to another, usually to provide separation or address sharing between a private network with multiple hosts and a public network or the Internet. In the context of IPv4 and IPv6 transition, these types of NAT are commonly referred to as *NAT44*, because they translate between IPv4 and IPv4 (vs. IPv4 to IPv6, IPv6 to IPv6, etc.). There is a proposed extension to NAT intended to preserve even more IPv4 resources. This proposal is called *Carrier-Grade NAT* (CGN)^[8]. The “Carrier Grade” in the name originates from the position of the NAT within the topology. Instead of NAT between a private and public network at the edge of a single network such as a home or business office, CGN is implemented inside of an ISP’s network and serves many customers simultaneously. These CGN implementations are typically scaled to handle thousands of simultaneous customer endpoints, often resulting in millions of simultaneous sessions. The RFC^[8] does not advocate the use of CGN; it describes how an ISP forced to deploy CGN can use it during IPv6 transition.

This sort of implementation addresses the need for an individual, globally unique IPv4 address for each of the ISP’s customers by allowing the ISPs to allocate each customer an IPv4 address that may not be globally unique and employ NAT to give them access to resources on the IPv4 Internet.

This sharing often allows ISPs to see oversubscription of public IPv4 addresses anywhere from 2:1 to more than 10,000:1 based on the type of applications behind the NAT and their simultaneous application layer port allocations and session counts. Most commonly, a CGN is used in conjunction with a local NAT on the customer's home network, creating two layers of NAT to traverse between the home network and the Internet. This model is commonly referred to as *NAT444*, because there is a translation layer between three sets of IPv4 addresses end to end.

A known problem with NAT is that it makes end-to-end communication and visibility between hosts more difficult, because it essentially hides hosts behind address translation. Because NAT is so common (nearly every home network and many commercial networks use NAT), networking applications have adapted so that they can discover the presence of a NAT and then change their behavior in order to maintain communications in the presence of NATs. However, the addition of this second layer of NAT often interferes with those workarounds, and undesirable or unpredictable results may occur^[9].

Over time it is likely that applications will again adapt to the impediments created by multiple layers of NAT, but it is not possible to anticipate and correct every potential problem that may be generated by adding this second layer of NAT. This reality should serve as a warning to those who provide services over an Internet connection: IPv6 support is extremely important. IPv6 is important because CGN means that ISP-controlled equipment will be actively involved in the path between content or application providers and their end users, making that relationship reliant on the service provider and the service provider's CGN vendor to an extent that was not necessary in the past. If the CGN implementation breaks something, it not only reflects on the CGN vendor and the service provider, it also reflects poorly on the relationship between the end customer and the service that that customer is using—and may cause that customer to form a negative opinion of the brand itself.

In other words, if a consumer uses an Internet-enabled application on a new Brand X smart TV and it does not work well, regardless of whether it is a problem with the CGN, the service provider, or something else entirely, the consumer may form the opinion and share via an online review that, "Brand X's TVs are ok, unless you try to use any of their fancy new features. I would not buy one if I were you, because Company X clearly does not know what it is doing." CGN represents a potentially significant increase in the amount of testing that must be done, especially in implementations that are uncommon, such as small, corner-case deployments, and closed architectures. Although using IPv6 is dependent on support at the client, the content or application provider, and the ISPs in between, if this support is present, it allows the content or application provider and client to bypass the service provider's CGN machinery—as well as any IPv4 NAT that may be present—and have a true end-to-end connection. This scenario restores control over the user experience back to the brand, and allows the ISP to resume supplying bit carriage.

IPv4 Addressing Requirements

Independent of the potential connectivity problems that NAT444 may create, it generates additional problems for the implementing ISP because of its need for IPv4 addresses. Because the CGN requires two sets of addresses—one for the inside (private) network and one for the outside (public) network—the ISP must identify address ranges to use for both. In order for its customers to be able to reach the Internet, the external pool must use globally unique IPv4 addresses. The number of addresses required will depend on the implementation of CGN, its scale profile, the topology of the network (how many hosts are behind each CGN instance), and the usage profile of the customer traffic. If the service provider has few or no available globally unique IPv4 addresses, it will have to either make changes in its network in order to reclaim addresses from elsewhere or make a request for a new allocation from its RIR^[29].

However, depending on the number of addresses that the RIR has available and its policies for justification, it may not be possible to obtain sufficient address space with this method. For example, in the Asia-Pacific region, the austerity policies in place mean that no matter how many IPv4 addresses they might have been able to justify using previous rules, most requesters are eligible for only a few hundred IPv4 addresses as their final allocation ever^[10]. This situation then requires the ISP to source IPv4 addresses via the IPv4 address transfer market^[11], adding additional cost to an already expensive deployment. In fact, if the service provider must source addresses via the transfer market, it may be more cost-effective to simply obtain more addresses and continue with business as usual without deploying CGN at all.

Internal Pool: Private Addressing Alternatives

When addresses are sourced for the public address pool, the service provider must also identify a pool of private addresses that is large enough for the provider to allocate one to each customer behind the CGN. Depending on the size and scale of the CGN, and how much the service provider is willing to segment and separate different sections of its network, this number could be a large block of addresses, perhaps even a /8 or more.

The most obvious choice might be to simply use address ranges reserved for private network use^[12], because there is a /8, a /12, and a /16 available for this purpose. However, this address space has some drawbacks. First, because of the prevalence of RFC 1918 addressing within most enterprise networks, there is a significant chance that the chosen address blocks may conflict with existing use of RFC 1918 space for management systems and other internal resources. Depending on the size of the CGN implementation, it may be necessary to instantiate multiple segments of the network where the entirety of RFC 1918 space is used, and in order for those segments to talk to one another or to talk to devices with conflicting numbering, significant additional complexity is required.

On the customer side, remote workers could experience problems where the address that they have been assigned is in a block that is already in use on their company's enterprise network, meaning that it may cause problems connecting to those hosts via a *Virtual Private Network* (VPN), or problems accessing some of the resources from the remote network. It may be possible to change the address assigned to the end user in an attempt to eliminate this conflict, but this approach is not necessarily scalable because it likely requires manual intervention in an automated address-assignment system, and there are limits to the number of times that a change of address can “fix” this problem without creating a problem for another user.

The other problem with the use of RFC 1918 space in the CGN is that it may conflict with the address space used by the customer's local network and NAT. For example, if a customer has a local network numbered out of **192.168.1.0/24** and the customer's router is allocated the external address of **192.168.1.85**, the router may fail to function properly because it has the same address range on both the internal and external interface. It may be possible through analysis to identify and carefully allocate addresses so that the portions of RFC 1918 commonly used by default in home gateway devices are not allocated. However, anecdotal evidence^[13] suggests that because of the wide variety of devices and implementations available—plus the fact that many users reconfigure their networks to use a different IP address range than the default configuration of the device—there simply may not be enough RFC 1918 addresses not in use to make this option viable.

“Squat” Space

Another alternative is to unofficially reuse one or more portions of the existing range of allocated globally unique IPv4 addresses as private addresses. In a network that does not talk directly to the Internet, such as a private network or VPN, the existing allocations of IPv4 space do not have any meaning, and so it is not strictly necessary to stick to RFC 1918 address space for numbering resources that are only internally accessible. Reuse of allocated IPv4 addresses has the benefit of not conflicting with in-use RFC 1918 addresses, but comes with its own set of problems. If the provider's own space is reused, the provider must carefully separate the private use from the public use to avoid conflicts, and managing this overlap may require additional complexity such as the use of VPNs as a method to separate the networks. The more common method is to reuse a block of addresses that is not currently allocated to the network using them; in other words, squatting on “someone else's” address space. Usually providers select space to use in this manner based on a low likelihood that either the owner will begin announcing the space on the global Internet or the users behind that network will need to connect to the users behind the ISP's NAT.

This method requires extreme care. The service provider must ensure that the routes for those prefixes are not inadvertently leaked to the global Internet, because such a leak could potentially cause a route-hijack denial-of-service attack, albeit an unintentional one. This method is even more risky if the ISP has one or more partners who have connections into the private portion of its network, because it may not have complete control of the announcement boundaries. Certainly there are safeguards such as tagging the announcements with *Border Gateway Protocol* (BGP) communities such as no-advertise or no-export^[14], but these solutions are not always practical, and they are not completely fail-safe. Depending on the chosen address space, the effects could be significant based on the true owner of that space—no service provider really wants to risk a public relations nightmare because it inadvertently caused an outage affecting the critical infrastructure of a large government agency or multinational corporation whose space it “borrowed” and then leaked to the Internet.

As a result of the IPv4 transfer market, it is quite likely that some of the address blocks that are not visible on the global Internet today and that some consider “safer” to squat on may end up being transferred to another party who plans to begin using them on the public Internet, and potentially requiring those squatting on the space to renumber to a different address block. ISPs can mitigate this risk somewhat by selecting multiple candidate blocks that are all preconfigured in the network such that it is relatively straightforward to make a rapid change from one block to another if the current block in use suddenly becomes unacceptable. Many ISPs use this method today, but because of the risks, it cannot be considered a real solution to the problem. Further, because it essentially encourages large service providers to violate the spirit—if not the letter—of the very policies that govern IP address allocation and use, standards bodies such as the IETF or policy organizations like RIRs cannot officially recommend such a solution.

Class E Addresses

A final alternative is to repurpose the reserved space in 240/4^[2] and make it available for this use. There have been several failed attempts to repurpose this reserved space within the IETF in the past few years^[15, 16]. The primary challenge with this alternative is that because the Class E space has been reserved for many years, many networking implementations are explicitly configured to reject this address space as invalid. Getting this problem fixed in software, and more importantly, getting those software upgrades deployed widely, may require a similar level of effort to that which is required to deploy IPv6, and deploying IPv6 would be a more effective use of the resources required to implement software and hardware changes.

Even in situations like a CGN where more of the implementation is under central control, this solution would be attractive only to a service provider that owns and operates the *Customer Premises Equipment* (CPE) routers for all of its customers such that it could work with a small number of vendors to get software patches to enable use of this space. Therefore this solution is also too limited in applicability to be seen as a general solution that a body like the IETF could recommend.

Shared Addresses

Although the solutions previously discussed may be acceptable in some applications, the risks and deficiencies make it necessary for other applications to find another source for the IP address blocks to be used on the private side of a CGN. It is possible to use “public” (globally unique) IPv4 addresses on the private side as well, but the challenges to obtaining additional public IPv4 addresses that were discussed previously are exacerbated by the even larger number of addresses required, so this solution is far from practical. Additionally, expecting each service provider that implements CGN to obtain its own address space for its inside pools would end up using a significant amount of the remaining IPv4 resources in a way that does not necessarily require globally unique addresses. However, because each service provider has different needs, growth rates, and applications, it is unclear that simply expecting each service provider to request space from the RIRs for its internal CGN pools would create a doomsday scenario where a few networks would use up all of the remaining available IPv4 space in a short time. Because CGN creates additional costs and complexity to implement and support, and could be viewed as “second-class” IPv4 service, most service providers are not likely to implement it across the entire network and all tiers of customers, instead preferring to implement it only as widely as absolutely necessary.

Service providers could choose to implement it only for net new customers (that is, growth above turnover); they could choose to implement it only in certain markets or for certain types of service where it is less likely to cause support problems and adversely affect the service. All of these things reduce the number of addresses that may be needed for the interior CGN address pool. Nevertheless, using globally unique addresses in an application that does not require unique addresses is not a good use of a very limited resource. That is why the idea of having a shared and reserved block of addresses specifically for use as an interior (private) pool on a CGN keeps resurfacing.

One alternative to formally reserving a shared transition space was to have a third party request a block of sufficient size from one or more of the RIRs and then make it available for use as a shared block by anyone who wishes to do so.

Given the “last /8” policies in effect at each of the RIRs, it would likely be quite difficult to justify sufficient space to be useful, and the cost involved in receiving and maintaining such a delegation would likely be prohibitive. There would also be challenges addressing potential abuse concerns.

Reserving a block via the standard IETF/IANA process meant that IETF would have a chance to document the problems and recommend best practices that must be considered when implementing something that uses this shared space. This policy would help to ensure that service providers and implementers are aware of these guidelines and recommendations. For example, many implementations make certain assumptions about address scope based on the address itself, such as assuming that RFC 1918 addresses are locally scoped, and then adapt their behavior accordingly. With things like squat space or an unofficially shared CGN space, implementers would not know that this space should be treated in a specific way, and the result may be more network breakage. The officially declared shared space must still wait for implementers to make changes to their products, and that may not always happen, but the chances are still better than if it had been done in an unofficial manner.

As you can probably see, this problem does not have a clear-cut and straightforward solution, and this situation has led to vigorous discussion within the standards and policy bodies that have discussed it. The next section gives a brief history of the activity in those bodies that ultimately led to the space being allocated.

Some History

Shared transition space proposals have been controversial each time a variant of the idea has come up for discussion. As IPv4 exhaustion became a reality and IPv6 deployment continued to lag, more people realized that IPv4 life-extension technologies such as CGN may be a necessary evil. When people saw CGN as a likely response to the gap between IPv4 exhaustion and wide IPv6 support, they began to understand the need for the shared transition space, and thus support for allocating that space has gradually grown.

Although variants of this discussion may be much older than the items discussed in the following paragraphs, this article focuses specifically on the history of the idea to allocate shared address space specifically for CGN. There was an unsuccessful proposal in 2005^[17] to update RFC 1918 with an additional three /8s, but this proposal was not specifically focused on CGNs, unlike some of the other proposals. The most recent set of proposals regarding shared CGN space first came up in the APNIC Policy *Special Interest Group* (SIG) in early 2008, where Policy Proposal 058 was discussed. APNIC members abandoned the proposal and recommended that the authors take the idea to the IETF, because that is the body that typically directs IANA to reserve IP address blocks for special uses such as this one^[18].

This recommendation resulted in a pair of Internet drafts^[19, 20], hereafter referred to as **shirasaki** in late 2008. The draft originally requested four /8s, with a minimum size of a /12, but subsequent revisions of the draft revised the request to only one /10. The draft never gained much traction within the IETF, but the authors continued to update it to keep the discussion going. In mid-2010, a second IETF draft^[21] was published, requesting that a full /8 be reserved for this purpose. It contained references to the **shirasaki** drafts, but provided additional justification and noted that a /10 may not be enough addresses for many of the large service providers.

The draft went through several revisions in the following months, eventually being replaced by a different draft^[22], hereafter referred to as **draft-weil**, which reduced the /8 requested down to a /10. Attendees of the IETF 79 meeting in Beijing, China, discussed the draft across two different working groups. People expressed strong opinions both in support of and in opposition to the idea, but the draft did not achieve clear consensus. With the future of the draft unclear, one of its authors submitted policy proposal 127 to the *American Registry for Internet Numbers* (ARIN)^[23]. The *ARIN Advisory Council* (AC) accepted this policy proposal as draft policy 2011-5^[24] in early 2011, and vigorously discussed it with participants at the ARIN XXVII public policy meeting and with members of the mailing list. At the conclusion of the discussion, the ARIN AC recommended the policy to the ARIN board for adoption.

This discussion took on additional urgency because during this time the IANA officially announced that it had exhausted the free pool of IPv4 addresses and delegated the last of the /8s to the RIRs in accordance with policy^[4]. The side effect of this exhaustion meant that it was no longer possible for IETF to direct IANA to reserve space unless IANA was directed to repurpose an existing reservation, because it had no unreserved address blocks of sufficient size to meet the request. Therefore, the IETF and one or more of the RIRs would have to work in concert to make a suitable IPv4 address block available, instead of it being solely under IETF's purview. ARIN staff reached out to the IETF's *Internet Architecture Board* (IAB) for guidance, because by strict interpretation^[25], ARIN was not authorized to make this allocation by itself. IAB reaffirmed this interpretation, and recommended that the matter be brought back to the IETF for (re)consideration^[26]. With this guidance, the authors revised **draft-weil-shared-transition-space-request** and reintroduced it for discussion. For a period of time, the document was split into two, with most of the long-form discussion of pros and cons being moved to a second draft^[27].

As of the publication date of this article, the secondary draft has expired without progressing, but most of the important information contained there was incorporated back into **draft-weil**. The document was not adopted by any IETF Working Group. Instead, an IETF Area Director sponsored it as an individual submission.

It went through its first IETF “Last Call” to gauge consensus and receive comments in August 2011. The subsequent discussion, revisions, and secondary last calls (October 2011 and January 2012) generated hundreds of messages on the IETF discussion list and a total of 12 versions of the document before it was approved for publication in February 2012.

The reason why the debate on this shared transition space was so spirited can be traced to a few critical concerns. First, although consensus-based RFCs documenting CGN^[8] were already approved, this draft allocating space specifically to facilitate its deployment became a referendum within the IETF on whether NAT444/CGN should even be used. If you believed that NAT444 and CGN were bad ideas, it was likely that you would also be against a shared transition space. From that perspective, shared transition address space provided a more complete solution to a problem that had been created by a “Bad Idea” that should not have been allowed to proceed in the first place. There was also resistance to what was deemed “waste” of the limited remaining blocks of IPv4 addresses to solve a problem that not everyone agreed was a real or important problem. Also, although IETF participants do not speak for their companies per se, this proposal had consistent support from numerous individuals employed by large residential broadband providers. As a result, some saw it as those service providers looking for a way to bail themselves out of a problem that they created by not deploying IPv6 rapidly enough to avoid having to use CGN. On the converse side of the argument, those in favor saw CGN as a largely foregone conclusion, and saw this proposal as simply a practical solution to a real problem.

The *Internet Engineering Steering Group* (IESG) ultimately sent a note to the IETF discussion list acknowledging the difficulty of coming to a decision on this matter and noting that some explanatory text would be added to RFC 6598:

“Colleagues,

The IESG has observed very rough consensus in favor of the allocation proposed in **draft-weil-shared-transition-space-request**. Therefore, the IESG will approve the draft. In order to acknowledge dissenting opinions and clarify the IETF position regarding IPv6, the IESG will attach the following note:

“A number of operators have expressed a need for the special purpose IPv4 address allocation described by this document. During deliberations, the IETF community demonstrated very rough consensus in favor of the allocation.

While operational expedients, including the special purpose address allocation described in this document, may help solve a short-term operational problem, the IESG and the IETF remain committed to the deployment of IPv6.”

In many ways, the final decision came down to the difference between theory and practice in the IETF's desire to make the Internet work better. Theoretically, making a CGN easier to implement has the potential to make the Internet work much more poorly, and could be seen as rewarding bad behavior (failing to deploy and support IPv6 in a timely fashion). However, in practice, making CGN harder to implement causes unnecessary pain and effort for operators and potentially for users, while having little or no effect on IPv6 deployment. Approving this shared transition space avoids the appearance that IETF is trying to punish operators or users for perceived past "sins" and helps to reinforce the idea that IETF is responsive to operational concerns and therefore still relevant to the operator community. It is unlikely that the result of this decision will have much bearing on an operator's plan for how widely, when, where, or even if it will deploy CGNs, and this article makes no such recommendations. However, I will reiterate that IPv6 is the long-term solution, and that the smallest CGN deployment possible will make for a less complex and less expensive network for the continued support of traditional IPv4 devices.

Acknowledgements

Special thanks to Ole Jacobsen for suggesting that I write this article, to Kirk Erichsen and Jason Weil for their review and comments, and to all involved in the discussion of the referenced IETF drafts and ARIN policy for giving me plenty to write about!

References

- [1] Victor Kuarsingh, Chris Donley, Jason Weil, Marla Azinger, and Christopher Liljenstolpe, "IANA-Reserved IPv4 Prefix for Shared Address Space," RFC 6598, April 2012.
- [2] IANA Address Assignments:
<http://www.iana.org/assignments/ipv4-address-space/ipv4-address-space.txt>
- [3] RIR Policies triggered by IPv4 Depletion:
ARIN:
https://www.arin.net/resources/request/ipv4_countdown.html
RIPE:
<http://www.ripe.net/internet-coordination/ipv4-exhaustion/reaching-the-last-8>
APNIC:
<http://www.apnic.net/community/ipv4-exhaustion/ipv4-exhaustion-details>
LACNIC:
<http://www.lacnic.net/en/politicas/manual111.html>
AFRINIC:
<http://www.afrinic.net/docs/policies/AFPUB-2010-v4-005-draft-05.htm>

- [4] Number Resource Organization (NRO), “Free Pool of IPv4 Address Space Depleted,” February 2011,
<http://www.nro.net/news/ipv4-free-pool-depleted>
- [5] Chris Donley, Christopher Liljenstolpe, Wesley George, and Lee Howard, “IPv6 Support Required for All IP-Capable Nodes,” RFC 6540, April 2012.
- [6] Consumer Electronics Association IPv6 Working Group,
http://www.ce.org/Press/CurrentNews/press_release_detail.asp?id=12139
- [7] World IPv6 Launch, <http://www.worldipv6launch.org/>
- [8] Sheng Jiang, Brian Carpenter, and Dayong Guo, “An Incremental Carrier-Grade NAT (CGN) for IPv6 Transition,” RFC 6264, June 2011.
- [9] Chris Donley, Lee Howard, and Victor Kuarsingh, “Assessing the Impact of Carrier-Grade NAT on Network Applications,” Internet Draft, work in progress, November 2011,
[draft-donley-nat444-impacts-03](#)
- [10] APNIC, “Policies for IPv4 address space management in the Asia Pacific region,”
<http://www.apnic.net/policy/add-manage-policy#9.10>
- [11] ARIN, “Understanding the IPv4 Transfer Market,”
https://www.arin.net/resources/transfers/transfer_market.html
- [12] Daniel Karrenberg, Yakov Rekhter, Eliot Lear, and Geert Jan de Groot, “Address Allocation for Private Internets,” RFC 1918, February 1996.
- [13] Akiro Kato, A sampling of RFC 1918 IP address usage in Japan,
<http://www.ietf.org/mail-archive/web/v6ops/current/msg06187.html>
- [14] Paul Traina, “BGP Communities Attribute,” RFC 1997, August 1996.
- [15] Vince Fuller, “Reclassifying 240/4 as usable unicast address space,” Internet Draft, work in progress, March 2008,
[draft-fuller-240space-02](#)
- [16] Paul Wilson, George Michaelson, and Geoff Huston, “Redesignation of 240/4 from ‘Future Use’ to ‘Private Use,’” Internet Draft, work in progress, September 2008,
[draft-wilson-class-e-02](#)

- [17] Tony Hain, “Expanded Address Allocation for Private Internets,” Internet Draft, work in progress, February 2005,
draft-hain-1918bis-01
- [18] Shirou Niinobe, Takeshi Tomochika, Jiro Yamaguchi, Dai Nishino, Hiroyuki Ashida, Akira Nakagawa, and Toshiyuki Hosaka, “Proposal to create IPv4 shared use address space among LIRs,” prop-058, January 2008,
<http://www.apnic.net/policy/proposals/prop-058>
- [19] Jiro Yamaguchi, Yasuhiro Shirasaki, Shin Miyakawa, Akira Nakagawa, and Hiroyuki Ashida, “NAT444 addressing models,” Internet Draft, work in progress, January 2012,
draft-shirasaki-nat444-isp-shared-addr-07
- [20] Ikuhei Yamagata, Shin Miyakawa, Akira Nakagawa, Jiro Yamaguchi, and Hiroyuki Ashida, “ISP Shared Address,” Internet Draft, work in progress, January 2012,
draft-shirasaki-isp-shared-addr-07
- [21] Jason Weil, Victor Kuarsingh, and Chris Donley, “IANA Reserved IPv4 Prefix for IPv6 Transition,” Internet Draft, work in progress, September 2010,
draft-weil-opsawg-provider-address-space-02
- [22] Victor Kuarsingh, Chris Donley, Jason Weil, Marla Azinger, and Christopher Liljenstolpe, “IANA-Reserved IPv4 Prefix for Shared Address Space,” Internet Draft, work in progress, February 2012. (Became RFC 6598^[1]),
draft-weil-shared-transition-space-request-15
- [23] ARIN Public Policy Mailing List (PPML), “Shared Transition Space for IPv4 Address Extension,” ARIN-prop-127,
<http://lists.arin.net/pipermail/arin-ppml/2011-January/019278.html>
- [24] “Shared Transition Space for IPv4 Address Extension,” ARIN policy 2011-5,
https://www.arin.net/policy/proposals/2011_5.html
- [25] Brian Carpenter, Fred Baker, and Michael Roberts, “Memorandum of Understanding Concerning the Technical Work of the Internet Assigned Numbers Authority,” RFC 2860, June 2000.
- [26] Internet Architecture Board, “Response to ARIN’s request for guidance regarding Draft Policy ARIN-2011-5,”
<http://www.iab.org/documents/correspondence-reports-documents/2011-2/response-to-arins-request-for-guidance-regarding-draft-policy-arin-2011-5/>

- [27] Stan Barber, Owen Delong, Chris Grundemann, Victor Kuarsingh, and Benson Schliesser, “ARIN Draft Policy 2011-5: Shared Transition Space,” Internet Draft, work in progress, September 2011,
draft-bdgks-arin-shared-transition-space-03
- [28] Geoff Huston, “Anatomy: Inside Network Address Translators,” *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [29] Daniel Karrenberg, Gerard Ross, Paul Wilson, and Leslie Nobile, “Development of the Regional Internet Registry System,” *The Internet Protocol Journal*, Volume 4, No. 4, December 2001.
- [30] Geoff Huston, “NAT++: Address Sharing in IPv4,” *The Internet Protocol Journal*, Volume 13, No. 2, June 2010.
- [31] *The Internet Protocol Journal*, Volume 14, No. 1, March 2011. This issue of IPJ is entirely devoted to the topic of IPv4 address depletion and IPv6 transition.
- [32] Michelle Cotton and Leo Vegoda, “Special Use IPv4 Addresses,” May 2012, **draft-vegoda-cotton-rfc5735bis-02**

WESLEY GEORGE has been working in IP networking for approximately 13 years, across operations, engineering and capacity planning, architecture, and design in large wired and wireless networks. He has been heavily involved in IPv6 evangelism and deployment for a surprisingly long time. He has been an active participant in IETF for 5 years, including serving as former co-chair of the IPv6 Renumbering (6renum) working group and current co-chair of the sunset4 working group. He was active in ARIN’s policy development process during the time that the policy discussed in this article was being addressed. He currently works for Time Warner Cable, but this article represents his views alone, and should not be mistaken for his current employer’s official stance on anything. Wes can be reached via twitter (@wesgeorge) or via wesley.george@twcable.com

December in Dubai: Number Misuse, WCIT, and ITRs

by Geoff Huston, APNIC

In November 1988, telephone companies from 178 nations sent their respective government representatives to the *World Administrative Telegraph and Telephone Conference* (WATTC) in Melbourne, Australia. At the time the generally cosy relationships between governments and their monopoly telephone companies often made it extremely difficult to see the difference between the government's representatives and those of the telephone company. The group resolved to agree to the rather grandly titled *International Telecommunication Regulations* (ITRs)^[1].

At this meeting the companies' national representatives agreed to a set of additional regulations that supplemented the binding regulations of the *International Telecommunication Convention*. The goals of these regulations were rather grand; they aspired to promote the "harmonious development and efficient operation of technical facilities, as well as the efficiency, usefulness and availability to the public of international telecommunication services." More practically, these ITRs defined the general principles for the provision and operation of international telephony services among signatories to the ITRs.

At that time the Internet was little more than a somewhat obscure experiment in advanced data communication protocols undertaken by a small number of researchers in North America and to a far smaller extent in Europe. However, since 1988 the Internet—and the world in which the Internet has flourished—has changed dramatically. If we view the rise of the Internet over the past 25 years as a product of an appropriately liberalized international regulatory regime as much as it was a product of the titanic shifts in computing and communications technologies that also occurred over this period, then we can make the case that the Internet of today is a product of these ITRs. And what a prodigious product it has been!

In Dubai, between the 3rd and 14th of December 2012, the nations of the world will convene at the 2012 *World Conference on International Telecommunications* (WCIT)^[2], and they intend to use this conference to review these 25-year-old ITRs and consider some proposed changes to this regulatory framework that underlie international telecommunications.

At the moment the international meeting cycle is ramping up to consider what aspects of the ITRs should be altered, what should stay the same, and what should be dropped. After all, much has happened in the past 25 years, and an argument could be made that the ITRs should be amended to better reflect today's world.

But the world is not exactly aligned at the moment about what should and what should not be folded into a new set of international regulatory obligations.

Some countries appear to be advocating for some quite specific measures to be added to the ITR to address what for them are characterized as otherwise unresolvable operational problems. Others are advocating a more general approach to have the ITRs explicitly embrace the Internet and fold references to the Internet in every place where specific carriage and service delivery technologies are referenced in the ITRs. It is when these two approaches intersect that the situation gets interesting.

In order to illustrate some of the underlying tensions that exist in this activity, I would like to take a specific example of a proposed amendment to the ITRs and consider in in terms of the broader context of telephony and the Internet.

The proposal I want to examine here concerns the topic that has been called “number misuse.” In telephony this term referred to an operating practice where a call to a dialled number is not routed to the destination subscriber who is located at that called number, but instead the call is re-routed to a different destination.

What we see in the “Number Misuse” proposal for a revision of the ITRs is an attempt to fold the concepts of “number misuse” and the Internet together, with a result that some countries want the ITRs to explicitly take on the concept of “IP Address and Routing Misuse” within the framework of national obligations through common regulatory action within the same scope as the telephony called number misuse. If successful, this effort would result in a regulatory obligation for governments to take necessary actions to investigate and prosecute such instances of so-called “number misuse.” The intended scope of such enforcement of such obligations would encompass not only the telephone network but also the Internet. Surely we all desire a global public communications network that operates with integrity, and surely we would want to see countries take the necessary actions to ensure that it happens. So why is this idea not exactly the best idea to appear in the ITR negotiation process so far?

Let’s look at the motivations behind number misuse in the world of telephone carriers and telephone services, and then look at how it could conceivably map in to the world of the Internet.

To understand the telephone world and where this problem of number misuse is coming from, it may be useful to understand a little of how money circulates in the phone world.

Telephony: Sender Pays

In many ways the telephone leaned heavily on the telegraph service for its service model, which, in turn, leaned on the postal service, establishing a provenance for the telephone service model that stretched back over some centuries to at least the 1680s and London’s Penny Post, if not earlier.

The postal service model that gained ascendancy over the preceding centuries was one in which the original sender of the letter paid for the entire service of letter delivery. If the postal service that received the letter in the first place needed to use the services of a different postal service to complete the delivery, neither the sender nor the intended recipient were aware of it. The postal services were meant to divide the money received from the sender to deliver the letter, and apportion it between themselves to compensate each service provider for undertaking its part in the delivery of the letter.

The telephone service, for the most part, operates in a very similar fashion. The caller pays for the entire cost of the call, and the called party pays nothing.

When both the caller and the called party are connected to the same carrier, the process is straightforward. The carrier charges the caller for the cost of the call and, presumably, some small (often not so small) margin for profit.

However, when we apply the same model to, say, international phone calls, the model is not so simple. The common desire on the part of the telephone operators was to preserve the same simple model: the caller pays. Now in this case the caller pays the presumably higher price of establishing a voice circuit from a carrier in one country in one part of the world to another carrier in another country in another part of the world. But now the caller's carrier should not keep all the revenue associated with the call. The other end, the *terminating carrier*, has also incurred costs in servicing this call. The arrangement that the telephone industry developed was the concept of "intercarrier call accounting financial settlements."

To explain this concept it may be useful to introduce the unit of a *call minute*, which is commonly used as a means of measuring a telephone call. What carriers establish between themselves on a bilateral basis is the intercarrier settlement cost per call minute of a telephone call that originates in one carrier and is terminated by the other carrier.

Now if both carriers can establish a value of a call-minute settlement rate where in both directions the call-minute termination costs roughly equate to the call-minute settlement rate, then in theory, at any rate, neither party is relatively advantaged over the other, irrespective of whether the callers are predominately located in one carrier or in the other carrier. In theory, such an arrangement should be financially neutral to both carriers.

However, although in theory practice and theory should align, in practice it rarely happens. What happened in the telephone case was that we saw some carriers set a call-minute call-termination settlement rate that was well above cost, while at the same time set its international call tariffs such that outbound calls were prohibitively expensive for local subscribers.

The result was that the local customers of these carriers found it cheaper to request that the other party call them—the desired outcome. The local carrier then generated income not by charging local subscribers but by revenue generated as an outcome of the call accounting settlement payments that were generated by the net imbalance of called versus calling call minutes.

Carriers all over the world played this game. For example, in France in the early 1990s it was some 5–10 times more expensive to call a U.S. number from France than it was to make a call between the same two numbers in the other direction. If you add in a further consideration, namely that in the 1980s many carriers were part of the public administration and were in effect government-operated national monopolies whose profits contributed to national revenue, then you get an outcome that is described in *Opinion No. 1* of the 1989 ITRs, under the heading “Special Telecommunication Arrangements,” namely: “...considering further that, for many Members, revenues from international telecommunications are vital for their administrations.”

Telephony Special Services and Number Misuse

It is often said that the only really major innovation in more than a century of the telephone service was the fax. Perhaps that is a little too unkind, but innovations in the delivered services industry were few and far between. However, there were many innovations that are important to this story, and the ones that are relevant here are *number redirect* and the so-called *premium* services.

The premium services attracted a higher call cost, and the carrier conventionally split the revenue from the service with the called service. These services traditionally included weather forecasts, sports results, new headlines (until the Internet became all but completely ubiquitous and decimated these services!), and so on. They also attracted the sex industry. However, in many countries such services were not permitted, so a conventional premium service was not an option for this industry.

As ever, we are naturally inventive, and some folks came up with a clever solution to use number redirect to redirect the call to this otherwise not-permitted premium service to another country. As part of this redirection, the premium service provider needed to reach an agreement with the new home carrier of the call-termination point to divide the international call accounting revenue provided by callers to this service between the carrier and the service provider. Not only did this arrangement effectively circumvent local regulations relating to locally provided premium services, it also leveraged off the international call accounting arrangements to the benefit of the premium service provider as well as the terminating carrier.

We may be inventive, but all too often we are greedy as well. The next step was to circumvent any arrangement with the destination carrier and redirect the call to an entirely different carrier.

One of the side effects of deregulation of the telephone industry in many countries was that in place of a single carrier that would receive all incoming international calls for a given country code there were numerous carriers that were ostensibly competing for these incoming calls. Instead of routing calls based solely on the dialled country code, carriers now could route calls based on number blocks within the country code, and use different transit routes based on number-block rules. What if a premium service provider took a number block from a country code and specified that all incoming calls were to be routed by a third-party carrier? That all sounds innocent enough, but what if this third party did not actually route the calls through to the country in question, but instead terminated the calls and still charged the calling carrier the international call accounting settlement rate? No doubt the service provider has gotten a better deal, so the service provider is happy, and the carrier that terminates the call is receiving a portion of the call settlement rate, so the terminating carrier is happy. But happiness is not universal here. The carrier in the called country code is getting nothing from this arrangement, even though its country call code is being used for these premium service calls. From the carrier's perspective it is being defrauded of what it might claim is legitimate international call accounting revenue through the "misuse" of the number block drawn from its country code.

If the country-code carrier could discover this unauthorized number-block diversion, then presumably it could withdraw the number block and stop the international call diversion. Unfortunately it does not always work. The carrier can withdraw the number block, but at times—and under perhaps somewhat shady circumstances—the premium service provider, and potentially the transit carriers, might still be able to convince local carriers that the number-block diversion is still legitimate. Although the country-code carrier might see the problem, the carrier's ability to enforce carriers in other countries to respect its authority regarding the use of number blocks drawn from its country code is not always clear. At times the carrier is effectively powerless to enforce a remedy.

And the scheme can be further refined. Why even enter into any form of discussion with the international carrier for a number block? Why not pick one or more of the more obscure national country codes, generate some number blocks from these codes, and then get a cooperative transit carrier to enter a number-block diversion request into the local carrier? The number block is perhaps drawn from a country code that already makes extensive use of third-party transit arrangements, the local carrier may not question the request, and the carriers in the countries from which the number blocks have been drawn may not have the resources to even detect that this event has occurred.

At this point we have arrived at the situation that is motivating some of the proposals to augment the ITRs in this round of negotiation. The position of the nations that have been highlighting this problem as being an important problem in the world of international telephony is that the unauthorized use of phone numbers drawn from their E.164^[3] telephone number block is, in their eyes, a case of “number misuse.”

The reason why they want to identify this situation and write it into the ITRs at this time is that they would like to involve governments in the role of enforcers of conformance with the conventions of management of telephone country codes. It appears that they would like to obligate governments to adopt a policy, as a common convention, that calls made to a country’s country code be directed such that the call request is sent to an authorized carrier located in the country, and to ensure that all authorized carriers essentially honor the integrity of the country codes of all other countries that use the E.164 country-code number plan.

It is also reasonable to ascribe the motivation for this measure as one that is intended to ameliorate the inexorable revenue leakage of the former rich money tap of international call accounting settlement payments. I am not sure that the various antics of the international premium service market are the true intended target of this measure. I suspect that the intended targets of this proposed regulatory measure are those carriers that have devised other methods to honor the intentions of their callers when they make an international phone call, and make the phone of the dialled number ring, yet at the same time bypass the traditional call accounting arrangements. Already *Voice over IP* (VoIP) trunking is commonplace, where the call is mapped into a VoIP call, and one way to bypass the conventional call accounting measures is to use a VoIP trunk to enter the dialled country, and then pass the call back into the *Public Switched Telephone Network* (PSTN) as a locally originated call, terminating it on the originally dialled number. The call is then subject to domestic intercarrier call-termination tariffs, which are generally far lower than their international counterparts.

The Internet and services such as Skype are exerting massive downward pressure on what carriers can charge for conventional phone services without encouraging all remaining customers to use Internet-based services. In an effort to retain some level of market share, it is now evidently more commonplace for carriers themselves to embrace IP-based approaches and bypass these imposed intercarrier international settlement charges. For many countries in the developing world, however, this shift represents a twofold financial blow. Not only are they seeing their foreign-sourced revenue stream disappear at the same rate as the call-termination minutes of conventional telephony vaporise, but they are also seeing this revenue stream being replaced by growing IP traffic volumes that represent a net cost to the national economy.

It should come as no surprise to see some countries attempt to advocate an international regulatory response that is intended to reverse this development, and restore the role of the international telephone network as a means of structural flow of monies from the business sector from the richer economies to the consolidated revenue stream of those poorer economies.

Internet Number Misuse

In and of itself, the previous discussion is by no means a novel discussion for the telephone world, and the tensions exposed by the continual erosion of the traditional telephone business through the onslaught of new technology is not at all surprising.

What is perhaps a bit surprising are the recent moves within the ITR preparatory activities that see numerous national delegations advocating pulling Internet addressing and routing into the same category of telephone-number regulation and also fold these factors into this matter of number misuse in a manner that would apply to both E.164 numbers and IP addresses.

Now some things do not readily translate from telephony to the Internet: there is no “National IP Address Plan” as a counterpart to the E.164 number plan, because the IP address plan is aligned to networks, as distinct from countries. However, you could take a broad view and find some form of mapping from the proposed recommendations regarding the use of E.164 networks to IP addresses. It would appear that the application of the proposals regarding number misuse would see a regulation to the effect that IP packets should be routed to the destination address specified in the packet, and not rerouted and terminated elsewhere. Surely this scenario describes part of the way the Internet works in any case. For the network to actually function, packets need to be passed to their addressed destination. Or so you would think.

And that is indeed what happens much of the time within the Internet. But by no means all of the time. As part of the normal course of operation of IP networks, many operators deploy equipment that intercepts packets and forms a synthetic response using the address of the intended destination. And many national administrations either operate—or mandate the operation of—equipment that inspects packets in transit and discards packets addressed to certain number blocks.

What is going on? Why do network operators regularly “misuse” IP addresses by deliberately intercepting packets and generating a synthetic response?

Packet Diversion

The most prevalent reason is the use of proxies, and, in particular, web proxies. These devices sit “on the wire” and intercept web fetches and cache the downloaded data.

When another user requests the same URL, the proxy uses the cached version of the content rather than forwarding the request on to the original site. This caching is by no means unusual: it is typical for web browsers to cache the most recently visited webpages and when the user returns to the page, the local cached copy is used rather than re-performing the download. For the browser and the network operator the rationale for this form of “address misuse” is the same: it is both a desire to improve performance for the end user and a desire to increase the efficiency of the network by reducing the data volumes being shifted across the transit links. So the outcomes are, on the whole, positive outcomes; users see improved performance and potentially lower costs for the service, using an interception technique that is generally transparent.

Is the deployment of a web proxy an instance of fraud?

Here is where another critical difference between the Internet and the telephone world comes into play. In the Internet the sender does not “pay all the way” to get a packet from its source to its intended destination. In general, every IP packet could be thought of as being partially funded by both the sender and the receiver.

The user who generated the packet pays for an *Internet Service Provider* (ISP) service, and the ISP may, in turn, purchase transit services from another ISP, and so on for sequenced transit services. However, at a peering exchange point, or within a provider network, the sender’s money runs out. The packet is not unfunded, however, for at this point the receiver’s services take over, and the packet transits a path that is funded by the receiver’s ISP’s transit services, and there to the receiver’s ISP and there to the receiver.

If a packet is diverted to a proxy, then who wins and who loses? Can we make the case that a party in this situation is being cheated?

As long as the proxy is a faithful proxy, then the user wins, insofar as the user experiences improved performance and the benefits of a more efficient network while still seeing precisely the same content. And the content provider wins, insofar as the content is delivered to the user without the incremental cost of packet handling at the content site. And the network service providers win, in so far as the amount of network traffic is reduced while the revenue levels remain constant. In this case there is no end-to-end service payment on the part of the user that would trigger an intercarrier settlement payment, so it is difficult to make the case that this action necessarily damages any party involved in the network transaction.

Given the widespread deployment of these proxy caching devices across the entire Internet, the beneficial outcomes of improved performance and network efficiency, and the option for content providers to use techniques that in effect mark content as not cacheable, it is extremely challenging to sustain a case that the use of proxies is a case of address misuse.

So the use of traffic diversion and intercepting proxies in the Internet is not generally regarded as an example of intentional fraud or even an accepted case of address misuse. It is just what we do today in the Internet.

Packet Interception

What about the deliberate interception and discarding of packets in flight? Surely this case is one of “misuse” of IP addresses?

That is a very hard case to make when you consider that such actions are exactly how firewalls work, and almost every network uses firewalls in some manner or other. The action of a firewall is to intercept all packets, and discard those that match some predetermined set of rules relating to acceptable and unacceptable packets.

Many users run firewalls that deliberately block all incoming connection requests unless they match quite specific rules.

Many ISPs run firewalls that deliberately block access to ISPs’ services from users who are not direct customers of the ISP.

Many countries have content regulations that block access to certain content, enforced either through government-operated facilities or through obligations imposed through the conditions associated with the carrier license within that country. The country I live in, Australia, imposes such constraints on its carriers for certain types of content, as does China through its much-reported national firewall facilities.

Users, service providers and carriers, and governments all use various forms of packet interception. Are we all guilty of number misuse? Should we support changes to the ITRs to obligate governments to stop this practice completely?

Aside from many other motivations for firewalls, security is a continuing concern in the Internet, and there is little doubt that although firewalls have not eradicated all forms of toxic traffic and associated abuse and attack, they are an important part of a larger story about securing the Internet. Irrespective of the various views that are expressed at a national level about censorship, intellectual property rights, and the position of common carriers and users, it seems counterintuitive to me that we would want to obligate governments to pull down our firewalls and filters as a necessary consequence of a revised set of ITRs.

Number “Misuse”

What this example illustrates is that the two networks—the traditional telephone network and the Internet—operate in very distinct and different ways. It not only encompasses differences between circuit and packet switching, but also reaches into the differences in the concepts of a network transaction, differences in the tariff structures, and, critically, differences in the way in which financial settlements are undertaken between service providers on the Internet.

Consider what could readily be acknowledged as an operating practice that defrauds operators in the world of telephony and negatively affects the services provided to telephone subscriber—that same practice in the Internet can result in positive outcomes used to enhance performance, reduce costs, and improve the operational efficiency of the service delivered to end users.

This case of attempting to regulate “number misuse” illustrates the fact that to take a stance of “one size fits all” when considering the topic of international regulation of telecommunications is a stance that has considerable risks of generating outcomes that are entirely inappropriate when translating a particular situation from telephony to the Internet.

WCIT and the ITRs—Where to Go from Here?

The international call accounting arrangements used by the telephone world, and the use of structurally embedded imbalances in call accounting settlement rates, are still major factors in the ITR discussions. This accounting imbalance is sanctioned in the resolutions of the 1988 World Administrative Telegraph and Telephone Conference, where *Resolution 3*, concerning the apportionment of revenue, provided for structural cross-subsidization of the developing world through asymmetric fixing of call accounting rates between the so-called developed and developing economies.

But in an increasing commercial world of telecommunications, where it is no longer a relatively exclusive collection of publicly funded monopolies that were an integral part of public utility service providers that in effect were an instrument of national governments, pushing the onus of an international developmental agenda onto an increasingly privatized commercial activity has been a less-than-comfortable fit. Private operators see this situation in a more dispassionate light as a business cost input, and seek to find ways to minimize this cost in order to improve the competitive positions of their businesses.

However, the changes in this industry over the past 25 years are so much larger than even this significant broad-scale shift in the onus of capital injection and operation from the public to the private sector. At the same time, we are seeing an even more fundamental shift in technology foundations, from circuits to packets with the introduction of the Internet into the picture. This shift has brought about profound shifts in the engineering of communications infrastructure and, as we have seen, it also has triggered profound shifts in the pricing of the consumer service, shifting from transactional pricing to a “connection rental” model where packet transit costs are bundled into the service. This bundling, in turn, has led to profound shifts in the manner in which money moves between the network operators themselves.

And perhaps of even greater and more lasting significance in this industry is the decoupling of carriage and content. We have now seen the rise of highly valuable content-centric enterprises that have business models that rely on a ubiquitous and abundant underlying communications infrastructure but are not financially beholden to the infrastructure operators. They have been able to forge direct relationships with consumers without having to deal with any form of mediation or brokerage imposed by carriage providers. The current values of these content enterprises dwarf the residual value of the carriage service sector, and the outlook for this sector is one of continuing shift in value away from carriage service providers and into the areas of content-based services.

Given the sheer scale of these changes in this industry over the past quarter century, it seems to me that the view that you can simply fold the Internet transparently into the current framework of the ITRs by the prolific insertion of “and the Internet” into the text of the regulations is simply not viable.

Packets are not circuits, and the mechanisms used to engineer packet networks are entirely different from those used with the circuit switches that supported traditional telephony services. This difference encompasses far more than engineering. The ways in which users pay for services differ, and this shift in the retail tariff structure of the Internet service implies a forced change in the way in which carriers interact to support a cohesive framework of network interconnection. The concept of a “call” really has no direct counterpart in the Internet. To extend this thought further into the area of “call accounting” and “caller pays” is again an extension that does not clearly map into the Internet. So when the existing ITRs refer to intercarrier call accounting financial settlements, there is no clear translation of such a concept into the Internet. When we extend this intercarrier interconnection framework into structural imbalances in call accounting settlement rates, and extend this framework further into the concepts of number misuse, all forms of connection between traditional telephony and the Internet are completely lost.

However, this conclusion should not imply that the ITRs are now an historic relic, completely overtaken by comprehensive shifts in both the technology and service models of today’s global communications network. Irrespective of the fine level of detail in these 25-year-old documents, the ideals behind the ITRs are indeed worthy ideals, and they should not be discarded lightly.

Ultimately, what we are dealing with here is the role of individual nation states with respect to a public communications service for the entire world. In setting forth a framework for supporting an efficient, effective, and capable global communications system, the obligations stated in the current ITRs relating to the promotion of international telecommunications services, and the endeavours to make such services generally available to the public, all remain thoroughly worthwhile objectives.

The concept that widely respected technology standards are critical to worldwide technical interoperability of any telecommunications service is also an important aspect, and again the recognition of this factor in the ITRs is a worthwhile consideration.

But, as we both review the changes of the past quarter century and try to peer into what may emerge over the next quarter century, perhaps less is best in this area of regulatory measures.

Rather than seeking to explicitly add various regulations that attempt to address specific incidents of number misuse, and instead of making rather clumsy efforts to include the Internet into the already detailed provisions relating to intercarrier settlement models of the increasingly historic traditional telephone network, perhaps the best set of ITRs we could have for tomorrow's world are national obligations that support a lightweight common regulatory framework.

This framework should be both more minimal with respect to describing or relying on particular technologies and service frameworks and more encompassing in scope in stating the overall objectives and common aspirations all nations share in supporting this unique, incredibly valuable common resource of a common communications service that truly embraces the entire world.

Postscript: "It's All Just Telecoms"

I received a comment soon after I wrote an early draft article that I thought would provide some further insight to the WCIT process, so here is the comment and some further thoughts on the topic:

The comment was in the form of a report from a preparatory meeting for WCIT earlier in 2012. Evidently there is a mood within certain parts of the ITR drafting process to simply say: "The ITRs should apply to the Internet in full, because the Internet is nothing more than a telecom service and should be treated that way."

In one sense it is true that the Internet is nothing more than a telecommunications service, but in the same way that the post, radio, television, and of course the telephone are also all just telecommunications services. But the nature of the particular service has many consequences, and the attempt to lump telephony and the Internet into the same form of regulatory handling is at best a somewhat misguided effort.

I truly wonder if, more than a century ago, the counterparts of today's government delegates, in a meeting of that august body, the *Universal Postal Union* (UPU), would have argued that a telephone conversation was just an exchange of letters without the artifice of paper, and that the telephone was indeed just a part of the postal service, because it is just "a communications service."

Indeed I am pretty sure their counterparts did precisely that, and for the next 80 years or more in many countries the Postmaster General operated the telephone service, and operated the wireless spectrum administration and regulated radio and television broadcasts, as well as operating the national postal service, the telegraph service, and telex services, all because “it’s all just communications.”

But, ultimately we changed this paradigm. We created distinct entities to administer different communications media and services because it is actually not “all just communications”—nor is it “all just telecoms.” Effective regulatory handling of these different communications mechanisms, using distinct forms of investment and finances, and at times entirely distinct regulatory frameworks and often distinct organizations and associated participatory arrangements, allows us to realize the true potential of these various services and do so efficiently and effectively. This recognition of a need for distinction in the regulatory frameworks for various services avoids the unfortunate situation of the stultifying dead hand of history misapplying one form of regulation on an entirely distinct and very different medium.

I suspect the best thing the postal folks, in the form of the UPU, ever did was to tell the telephone folks “hail and farewell” and let them get on with their role using an organization specifically designed to meet their collective needs in supporting telephony.

It may be well and truly time for the telephone folks, in the form of the *International Telecommunications Union* (ITU), to come to a similar arrangement in its dealings with the Internet!

Disclaimer

These views do not necessarily represent the views or positions of the Asia Pacific Network Information Centre.

Further Reading:

- [1] The current International Telecommunication Regulations (1988):

http://www.itu.int/dms_pub/itu-t/oth/3F/01/T3F010000010001PDFE.pdf

- [2] World Conference on International Telecommunications (WCIT-12),

<http://www.itu.int/en/wcit-12/Pages/default.aspx>

- [3] Geoff Huston, “ENUM—Mapping the E.164 Number Space into the DNS,” *The Internet Protocol Journal*, Volume 5, No. 2, June 2002.

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001.

E-mail: gih@apnic.net

Behind the Curtain: IP Fast Reroute

by Russ White, Verisign

The field of network and protocol engineering has three watchwords: *faster*, *bigger*, and *cheaper*. Although we all know the joke about choosing two out of the three, the reality of networking is that we have been doing all three for years—and it doesn't look like there is any time on the horizon when we will not be doing all three.

In that spirit, *IP Fast Reroute* addresses all three of these watchwords. Fast—you are probably thinking—is obvious, but what about bigger and cheaper? Fast Reroute provides the network designer with some trade-offs in the space of redundancy through additional backup links against deploying protocol changes, and network stretch against the size of a failure domain, so you can—in theory—build larger, less-redundant failure domains with Fast Reroute than without.

But to understand these effects, we need to go behind the curtain, understanding Fast Reroute as more than a few configuration options. This article first looks at the motivation behind IP Fast Reroute, and then discusses four different techniques, or stages, in the Fast Reroute story.

What Is Your Motivation?

To really discuss network speed, we need to be able to define how fast “fast” really is. In the 1980s, a network was fast if it could converge in 90 seconds or less (the longest time the *Routing Information Protocol* [RIP] could take to converge). As we moved into more advanced Distance-Vector and Link State protocols (*Enhanced Interior Gateway Routing Protocol* [EIGRP], *Open Shortest Path First* [OSPF], and *Intermediate System-to-Intermediate System* [IS-IS]), 5-second convergence became the norm. We learned to tweak timers to get to convergence times faster than 1 second.

But what if we need convergence that is faster than less than 1 second? What if we need to converge so fast that the only packets lost are either in flight or in a buffer waiting to be serialized onto the link? And what if we need to be able to handle a large number of prefixes with minimal network disruption due to link or device failures?

IP Fast Reroute techniques come into play in this situation.

Preinstalled Backup Paths

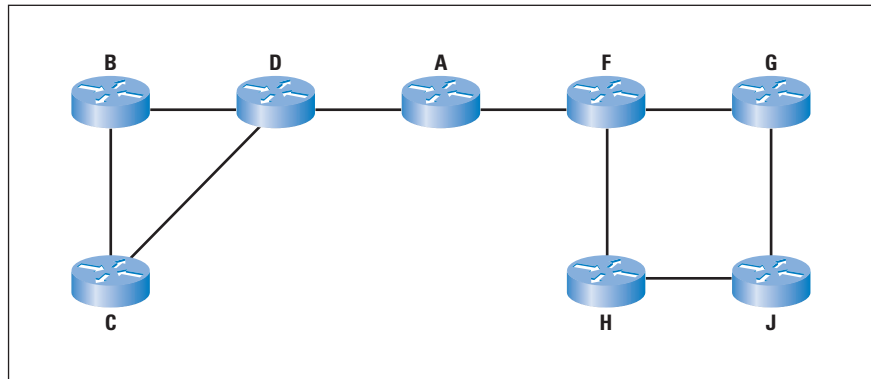
Although it is often sold as a Fast Reroute technique, *preinstalled backup paths* really are not; rather they support other Fast Reroute techniques at the protocol level. If the protocol has calculated a loop-free path that is an alternate to the current best path, this alternate path can be installed in the forwarding table so it is readily available for use in case the best path fails.

This solution does provide immediate failover at the hardware level, but the alternate path must be calculated to be installed. How is this alternate path computed?

Loop-Free Alternates

The first mechanism available for calculating an alternate path is with *Loop-Free Alternates*. To understand this mechanism, we must make a short detour into graph theory (or geometry, if you prefer). Use the following network as an example:

Figure 1: Network for Loop-Free Alternates



Assume:

- A is the destination.
- B's best path is through D to A.
- G's best path is through F to A.

What is the key to allowing B to forward traffic through C toward A if the B → D link fails? B must know the traffic it forwards to C (for A) will not be forwarded back to B itself. How can B know C will forward the traffic to D, rather than to B itself? By examining the metric at C toward A.

In EIGRP, B knows C's metric toward A because the routing protocol includes this information in the update. In a link state protocol (OSPF or IS-IS), B can calculate C's cost to A directly by running *Shortest Path First* from C's perspective (given B and C share the same link state database).

Loop-free alternates are simply calculating whether any given neighbor will forward traffic to any particular destination back to you, or on toward the destination. If a neighbor would forward the traffic on toward the destination, then it is a loop-free alternate.

Under what conditions would C forward traffic sent from B back to B? *If C is using B as its best path (or one of its best paths) toward A.*

What about G? If it forwards traffic to J toward A, will J return the traffic to G itself? In this four-hop ring, there are two possible configurations:

- J is using H as its best path. In this case, traffic forwarded by G to A through J will be correctly forwarded. Note, however, that in this case H cannot use J as an alternate path toward A, because any traffic H sends to A through A will loop back to H itself.
- J is using G as its best path. In this case, J can use G as a loop-free alternate, but G cannot use J as a loop-free alternate.

No matter how you work the metrics in the four-hop ring case, there will always be at least one device that does not have a loop-free alternate path to A.

Split Horizon and Loop-Free Alternates

If the concept of loop-free alternates is difficult to understand by considering the problem in this way, another useful way to look at the problem is through the distance-vector idea of *split horizon*. To review, the split horizon rule states:

Do not advertise a route to a destination toward a neighbor you are using to forward traffic to that same destination.

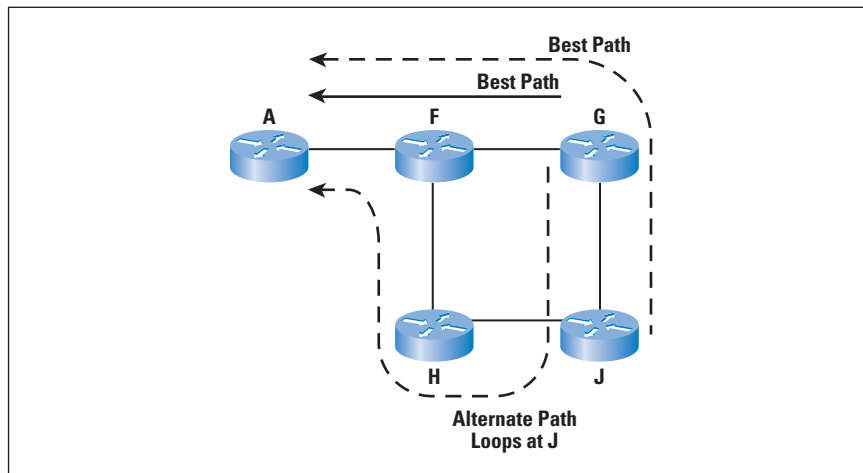
If C is forwarding traffic toward A to B, then C will not advertise A to B, meaning B will not even know about this alternate path, preventing a loop even if B's best path to A fails. If you always consider where a distance-vector protocol will split horizon, you will always be able to see where loop-free alternates will fail to provide an alternate path to any given destination.

Getting Around the Loops

If we want to design a system that will find every possible alternate path toward a given destination, rather than just finding those that are not normally taken out by split horizon anyway, what must we do? We need to find a way to route through a neighbor to some distant next hop without that neighbor actually forwarding the traffic back to the originating router.

To put this concept in more concrete terms, examine the following network as an example:

Figure 2: Alternate Path Loops



If G wants to use the path through J as an alternate path, then it must somehow figure out how to forward traffic to J without J returning the traffic to G itself. How can this process be done? G can tunnel the traffic through J to some device somewhere beyond J; therefore, every mechanism beyond loop-free alternates must use some form of tunneling to resolve the Fast Reroute problem. Calculating the point to which G needs to tunnel is the topic of the remaining mechanisms.

Not-Via

Even though we might be working with a link state protocol, it is easiest to understand *Not-via* in terms of a distance-vector protocol and split horizon. Not-via essentially begins with the observation that G does not have an alternate path to A through J in this case because J will not advertise such a route. *J is, in fact, using G as its best path toward A, so the path from G through J to A cannot be viable.*

The solution is not just simply having J advertise the route to A because traffic forwarded by G toward A through J will simply be looped back to G itself. So what is the solution?

In the case of Not-via, F advertises a route to itself through H only (not through G). This route will be advertised through H, then J, and finally to G. When G receives this route, it can determine that this path is an alternate path to A because its best path to A is normally through F. Any path that can reach F not through (or not via) its best path to F must, necessarily, be a loop-free alternate path to F. To reach A through F, however, G must tunnel to F directly, thereby avoiding the problem of J returning traffic destined to A back to G.

The address F advertises through H only is called “F Not-via G,” and that is why this system is called “Not-via.” This mechanism works in every topology (so long as an alternate path exists). The one downside to Not-via is that for each protected link or node, a new advertisement must be built and advertised through the network.

Disjoint Topologies

The problem of finding a next hop that passes over the split-horizon point can also be solved using the ability to form multiple disjoint topologies—multiple topologies that do not share the same links (or nodes, in some cases) to reach the same set of destinations. If this information sounds complex, that is because it is complex; a lot of hours and thought have gone into various systems to build and use multiple disjoint topologies within a single physical network. But there is a moderately simple way, referring back to Figure 2. In this network, G can take the following steps:

1. Remove the $G \rightarrow F$ link from its local database temporarily (just for this calculation).
2. Calculate the best path to F.
3. If an alternate path to F exists, mark this alternate path as a second topology.

4. If its path to F fails, place all traffic that would normally pass across $G \rightarrow F$ on this alternate topology.

It might not be obvious from this set of actions, but these actions will actually cause G to discover that it is, in fact, on a ring, and that it can place traffic on the opposite direction on this ring to get traffic to the same destination. Placing the traffic it would normally send to F via $G \rightarrow F$ on a separate topology overcomes the forwarding table at J, a process that would loop the traffic back to G itself. You could use a tunnel to F instead of a separate topology; tunnels are, in effect, a disjoint topology seen in a different way.

Conclusion

What advantage does IP Fast Reroute provide the network designer? The ability to reduce the amount of physical redundancy while maintaining the same actual level of redundancy in the network. Moving to Not-via or disjoint topology solutions removes the need to manually manage link costs as well, while adding only moderate complexity at the protocol level.

IP Fast Reroute is an interesting technology just on the edge of adoption that will be useful in campus, data center (through Layer 2 routing), and standard Layer 3 network designs.

For Further Reading

Work is currently active on the disjoint topology mechanism within the research community and the IETF; in particular, the following drafts will be of interest to anyone who wants to learn more:

- [1] Alia Atlas, Robert Kebler, Maciek Konstantynowicz, Andras Csaszar, Russ White, and Mike Shand, “An Architecture for IP/LDP Fast-Reroute Using Maximally Redundant Trees,” Internet Draft, work in progress, October 2011,
draft-atlas-rtgwg-mrt-frr-architecture-01
- [2] Alia Atlas, Gabor Envedi, and Andras Csaszar, “Algorithms for Computing Maximally Redundant Trees for IP/LDP Fast-Reroute,” Internet Draft, work in progress, March 2012,
draft-envedi-rtgwg-mrt-frr-algorithm-01
- [3] Stefano Previdi, Mike Shand, and Stewart Bryant, “IP Fast Reroute Using Not-via Addresses,” Internet Draft, work in progress, December 2011,
draft-ietf-rtgwg-ipfrr-notvia-addresses-08
- [4] Clarence Filsfils and Pierre Francois, “LFA applicability in SP networks,” Internet Draft, work in progress, January 2012,
draft-ietf-rtgwg-lfa-applicability-06

RUSS WHITE is a Principle Research Engineer at Verisign. He has co-authored numerous technical books, RFCs, and software patents. He focuses primarily on network complexity, network design, the space where routing and naming intersect, control-plane security, protocol design, protocol operation, and software-defined networks. E-mail: riwhite@verisign.com

Letters to the Editor

Ed.: We received several letters in response to the article “A Retrospective: Twenty-Five Years Ago,” by Geoff Huston, published in the previous issue of this journal. Here is some of the feedback:

Hi Geoff,

Just wanted to show my appreciation for your nice article. As an ex-DEC who moved to WorldCom after my MSc in Computer Engineering & Telecoms with a Master’s project on IP signaling over ATM, I can certainly relate to a large part (not all ;-)) of what you wrote.

I normally don’t read such long articles, but had to make an exception as I kept interested until the end!

Thank you!

—Pedro Paiva, Etoy, Switzerland
`pedro.paiva@a3.epfl.ch`

Greetings Geoff,

I just wanted to let you know that I really enjoyed your recent article, “A Retrospective: Twenty-Five Years Ago,” published in *The Internet Protocol Journal*. I lived through most of the history that you talked about as I came up through the telecom industry and then finished off my career at Cisco.

It certainly is interesting to reflect back on all the past controversy around network infrastructure design and how competing ideas and philosophies played out. (Talk about losers, remember *Switched Multi-megabit Data Service* (SMDS) driven by the *Regional Bell Operating Companies* (RBOCs)? While at Nortel, I remember once in a design review meeting that one of our BNR geeks put up a slide (overhead foil back then) that showed various network evolution scenarios. The last one was an “oh-by-the-way, there’s this theory that the Internet could take over the world” (of network infrastructure). All the room snickered. Who’s laughing now?

There was as much energy, maybe more, put into defending architectures based on market control as there was on technological elegance. Still, it is a fascinating and dynamic industry full of extremely smart people with clever ideas, and I enjoyed every minute of it.

I started at “the phone company” in the late 1960s and it has been quite a journey from relay-driven switches controlling tip and ring loops to the current *Multiprotocol Label Switching* (MPLS) backbone networks, terabit switching, and hitching rides on photons.

Thanks for your insight and for your well-written article.

Best regards,

—Marc Williams
willimarc@gmail.com

The author responds:

Hi Marc,

Thanks for your note and your recollections from some 25 years ago.

I recall SMDS as well. If I recall correctly, this was an invention coming out of a university in Western Australia. Elsewhere in the world it was marketed as a 34-Mbps product. In Australia it was marketed in 2-Mbps and 10-Mbps forms (evidently the telco thought that we primitive Aussies were not “ready” for any higher speed!). I was a customer of their 10-Mbps product, and experienced some disappointment when it became evident that 10 Mbps was a theoretical peak that was simply unachievable because the inline PCs that were used for packet accounting slowed the throughput of any SMDS link down to just 3 Mbps! So in Australia SMDS was largely killed by the telco and it was never really used for high-speed digital trunk services.

I experienced a similar reaction to the Internet in the late 1980s as you have observed, when, in response to suggesting that the universities were about to build a national IP network, many of the telco managers did the polite snicker performance and then suggested that we should “get with the times,” sign up as customers of their national ATM network, and leave the engineering to them. I’m glad the universities saw through it and supported me in persisting along the path to a national IP network. It was a strange moment some 6 years later when the same telco came knocking on our door to make an offer to buy the network from the universities because their own efforts to construct an IP product were simply getting nowhere at the time.

It has indeed been quite a journey, and I too have enjoyed every bit of it!

Kind regards,

—Geoff, Chief Scientist, APNIC
gih@apnic.net

Hello Geoff!

I haven't chuckled that much in years; what great memories. A few of my strong memories:

- Lack of documentation for new functions in software required an off-net test network and a Sniffer. The amount of hours spent figuring exactly what the function was doing or wasn't doing could fill an ocean. Absolutely my favorite activity and still is.
- I inherited a stat-mux system that was transporting ASCII terminals back to a centralized DEC terminal server arrangement. Hated it with a passion. One day, after a couple of beers, a light bulb came on that Ethernet is a stat-mux, so I bought a couple of Cisco AGS units, remotely installed a terminal server and an AGS, hauled it back to the other AGS in the central location and danced a jig, and then I started ripping out the old WAN stat-mux the following week.
- Anything relying on a token for timing is pure evil. You never know when you've engineered a TTL exhaust until it happens, and that can be based on Distance + Nodes or pure application coincidence. Ring resets are the devil's work. Token-based systems are not stat-muxs, but Ethernets are; that's why Ethernet survived and is the "last man standing."
- I totally agree with your comments surrounding the "cloud." I can remember that the distributed-versus-centralized fad has occurred at least four times over the past 25 years ...
- Z80: I built my first PC with a Z80; thank goodness for the peek-and-poke function!
- OEM would claim anything was portable as long as it had a carrying handle attached, even if it took two people to carry it.
- I fell in love with TCP/IP very early for the simple reason that it has the best of both worlds: a tightly coupled connection and connectionless protocol. It is much faster to troubleshoot or modify because IP requires a different expertise than TCP, and when you run across individuals who can work across the layers, hire them!

So, a lot of fond memories. I started out as a telemetry engineer on the Apollo project and I thought that was challenging and fulfilling. But, it doesn't hold a candle to the 1984–1995 period.

Oh, one other thing; I take umbrage to "...the annoying persistence of FORTRAN." That's the first language I learned back in the late '60s and I still have an active compiler on an old laptop that I still program on ... LOL!!

Keep attacking the certificate situation! The current situation is a disgrace, and I fully support the concept presented by Barnes: let's hurry it up!

Regards,

—Paul Dover
pdover@centeriem.com

The author responds:

Hi Paul,

Thanks for those recollections. I too spent a massive amount of time starting as a protocol analyzer, trying to make an IBM PC look enough like a Uniscope to allow file transfer between the PC and the Univac mainframe—no doubt it was a character-forming experience, but all I can say now is thank goodness for *tcpdump* and *wireshark*!

Thanks for your note—I truly appreciate the feedback!

Warm regards,

—*Geoff, Chief Scientist, APNIC*
gih@apnic.net

Dear Ole,

Congratulations on your 25-year anniversary!

You can tell how well people enjoy their professions by how great their products are, and yours is in the “excellent” category.

Regards,

—*Paul Dover*
pdover@centeriem.com

Ole,

Congratulations on your reaching a major milestone: 25 years of technology publishing! We are glad that you are continuing this service through *The Internet Protocol Journal* and look forward to many more years in this field.

Best,

—*T. Sridhar*
tsridhar@ieee.org

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ contains standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2012 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

September 2012

Volume 15, Number 3

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
Leaping Seconds	2
The Internet of Things	10
The Demise of Web 2.0	20
Binary Floor Control Protocol	25
Fragments	30

Internet devices use various forms of timers and timestamps to determine everything from when a given e-mail message arrives to the number of seconds since a particular device was rebooted. Most systems use the *Network Time Protocol* (NTP) to obtain the current time from a large network of Internet time servers. NTP will be the subject of a future article in this journal. This time we will focus our attention on the *Leap Second*, which is occasionally applied to *Coordinated Universal Time* (UTC) in order to keep its time of day close to the *Mean Solar Time*. Geoff Huston explains the mechanism and describes what happened to some Internet systems on July 1, 2012, as a result of a leap second addition.

The Internet of Things (IoT) is a phrase used to describe networks where not only computers, smartphones, and tablets are Internet-aware, but also autonomous sensors, control systems, light switches, and thousands of other embedded devices. In our second article, David Lake, Ammar Rayes, and Monique Morrow give an overview of this emerging field which already has its own conferences and journals.

The *World Wide Web* became a reality in the early 1990s, thanks mostly to the efforts of Tim Berners Lee and Robert Cailliau. The web has been a wonderful breeding ground for new protocols and technologies associated with access to and presentation of all kinds of media. The phrase *Web 2.0*, coined in 1999, has, per Wikipedia, "...been used to describe web sites that use technology beyond the static pages of earlier web sites." David Strom argues that the term is no longer appropriate and that we have moved on to a new phase of the web, dominated by mobile devices and Social Networking.

The last few years have seen great advances in Internet-based collaboration tools. Sometimes referred to as *Telepresence*, these systems allow not only high-quality audio and videoconferencing, but also the use of shared whiteboards and other presentation material. In our final article, Pat Jensen describes one important component of such systems, namely the *Binary Floor Control Protocol* (BFCP), which the IETF's XCON Centralized Conferencing working group has developed.

As always we welcome your feedback on anything you read in this journal. Contact us by e-mail at ipj@cisco.com

—Ole J. Jacobsen, Editor and Publisher

ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

Leaping Seconds

by Geoff Huston, APNIC

The tabloid press is never lost for a good headline, but in July 2012 this one in particular caught my eye: “Global Chaos as Moment in Time Kills the Interwebs.”^[1] I am pretty sure that “global chaos” is somewhat “over the top,” but a problem did happen on July 1 this year, and yes, it affected the Internet in various ways, as well as affecting many other enterprises that rely on IT systems. And yes, the problem had a lot to do with time and how we measure it. In this article I will examine the cause of this problem in a little more detail.

What Is a Second?

I would like to start with a rather innocent question: What exactly is a *second*? Obviously it is a unit of time, but what defines a second? Well, there are 60 seconds in a minute, 60 minutes in an hour, and 24 hours in a day. That information would infer that a “second” is 1/86,400 of a day, or 1/86,400 of the length of time it takes for the Earth to rotate about its own axis. Yes?

Almost, but this definition is still a little imprecise. What is the frame of reference that defines a unit of rotation of the Earth? As was established in the work a century ago in attempting to establish a frame of reference for the measurement of the speed of light, these frame-of-reference questions can be quite tricky!

What is the frame of reference to calibrate the Earth’s rotation about its own axis? A set of distant stars? The Sun? These days we use the Sun, a choice that seems logical in the first instance. But cosmology is far from perfect, and far from being a stable measurement, the length of time it takes for the Earth to rotate once about its axis relative to the Sun varies month by month by up to some 30 seconds from its mean value. This variation in the Earth’s rotational period is an outcome of both the Earth’s elliptical orbit around the Sun and the Earth’s axial tilt. These variations mean that by the time of the March equinox the *Solar Day* is some 18 seconds shorter than the mean, at the time of the June solstice it is some 13 seconds longer, at the September equinox it is some 21 seconds shorter, and in December it is some 29 seconds longer.

This variation in the rotational period of the Earth is unhelpful if you are looking for a stable way to measure time. To keep this unit of time at a constant value, then the definition of a second is based on an ideal version of the Earth’s rotational period, and we have chosen to base the unit of measurement of time on *Mean Solar Time*. This mean solar time is the average time for the Earth to rotate about its own axis, relative to the Sun.

This value is relatively constant, because the variations in solar time work to cancel out each other in the course of a full year. So a second is defined as $1/86,400$ of mean solar time, or in other words $1/86,400$ of the average time it takes for the Earth to rotate on its axis. And how do we measure this mean solar time? Well, in our search for precision and accuracy the measurement of mean solar time is not, in fact, based on measurements of the sun, but instead is derived from baseline interferometry from numerous distant radio sources. However, the measurement still reflects the average duration of the Earth's rotation about its own axis relative to the Sun.

So now we have a second as a unit of the measurement of time, based on the Earth's rotation about its own axis, and this definition allows us not only to construct a uniform time system to measure intervals of time, but also to all agree on a uniform value of absolute time. From this analysis we can make calendars that are not only "stable," in that the calendar does not drift forward or backward in time from year to year, but also accurate in that we can agree on absolute time down to units of minute fractions of a second. Well, so one would have thought, but the imperfections of cosmology intrude once again.

The Earth has the Moon, and the Earth generates a tidal acceleration of the Moon, and, in turn the Moon decelerates the Earth's rotational speed. In addition to this long-term factor arising from the gravitational interaction between the Earth and the Moon, the Earth's rotational period is affected by climatic and geological events that occur on and within the Earth^[2]. Thus it is possible for the Earth's rotation to both slow down and speed up at times. So the two requirements of a second—namely that it is a constant unit of time and it is defined as $1/86,400$ of the mean time taken for the Earth to rotate on its axis—cannot be maintained. Either one or the other has to go.

In 1955 we went down the route of a standard definition of a second, which was defined by the *International Astronomical Union* as $1/31,556,925.9747$ of the 1900.0 *Mean Tropical Year*. This definition was also adopted in 1956 by the *International Committee for Weights and Measures* and in 1960 by the *General Conference on Weights and Measures*, becoming a part of the *International System of Units* (SI). This definition addressed the problem of the drift in the value of the mean solar year by specifying a particular year as the baseline for the definition.

However, by the mid-1960s this definition was also found to be inadequate for precise time measurements, so in 1967 the SI second was again redefined, this time in experimental terms as a repeatable measurement. The new definition of a second was 9,192,631,770 periods of the radiation emitted by a Caesium-133 atom in the transition between the two hyperfine levels of its ground state.

Leaping Seconds

So we have the concept of a second as a fixed unit of time, but how does this relate to the astronomical measurement of time? For the past several centuries the length of the *Mean Solar Day* has been increasing by an average of some 1.7 milliseconds per century. Given that the solar day was fixed on the Mean Solar Day of the year 1900, by 1961 it was around a millisecond longer than 86,400 SI seconds. Therefore, absolute time standards that change the date after precisely 86,400 SI seconds, such as the *International Atomic Time* (TAI), get increasingly ahead of the time standards that are rigorously tied to the Mean Solar Day, such as *Greenwich Mean Time* (GMT).

When the *Coordinated Universal Time* (UTC) standard was instituted in 1961, based on atomic clocks, it was felt necessary that this time standard maintain agreement with the GMT time of day, which until then had been the reference for broadcast time services. Thus, from 1961 to 1971 the rate of broadcast time from the UTC atomic clock source had to be constantly slowed to remain synchronized with GMT. During that period, therefore, the “seconds” of broadcast services were actually slightly longer than the SI second and closer to the GMT seconds.

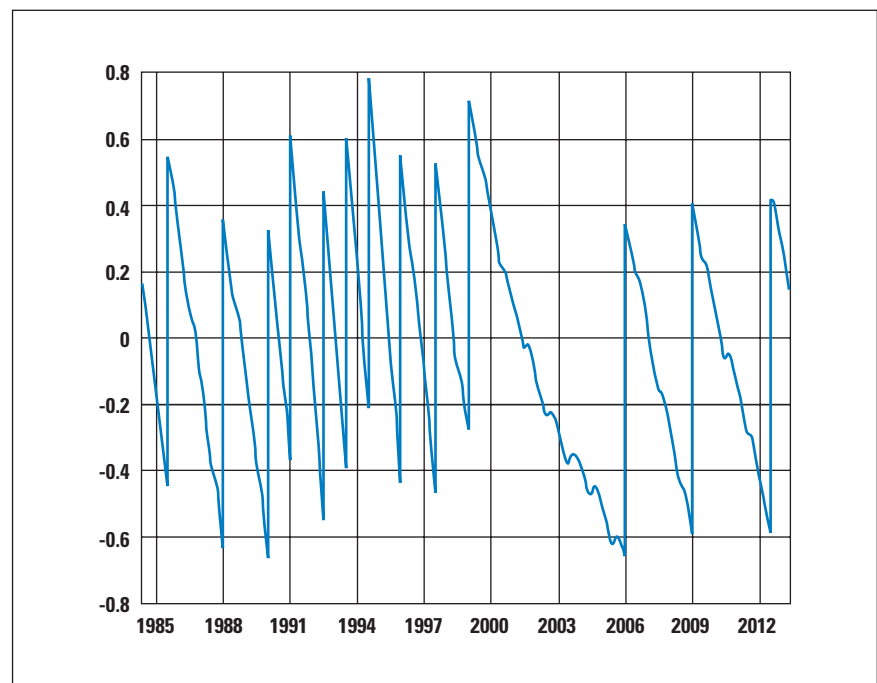
In 1972 the *Leap Second* system was introduced, so that the broadcast UTC seconds could be made exactly equal to the standard SI second, while still maintaining the UTC time of day and changes of UTC date synchronized with those of *UT1* (the solar time standard that superseded GMT). Reassuringly, a second is now a SI second in both the UTC and TAI standards, and the precise time when time transitions from one second to the next is synchronized in both of these reference frameworks. But this fixing of the two time standards to a common unit of exactly 1 second means that for the standard second to also track the time of day it is necessary to periodically add or remove entire standard seconds from the UTC time-of-day clock. Hence the use of so-called leap seconds. By 1972 the UTC clock was already 10 seconds behind TAI, which had been synchronized with UT1 in 1958 but had been counting true SI seconds since then. After 1972, both clocks have been ticking in SI seconds, so the difference between their readouts at any time is 10 seconds plus the total number of leap seconds that have been applied to UTC.

Since January 1, 1988, the role of coordinating the insertion of these leap-second corrections to the UTC time of day has been the responsibility of the *International Earth Rotation and Reference Systems Service* (IERS). IERS usually decides to apply a leap second whenever the difference between UTC and UT1 approaches 0.6 second in order to keep the absolute difference between UTC and the mean solar UT1 broadcast time from exceeding 0.9 second.

The UTC standard allows leap seconds to be applied at the end of any UTC month, but since 1972 all of these leap seconds have been inserted either at the end of June 30 or December 31, making the final minute of the month in UTC, either 1 second longer or 1 second shorter when the leap second is applied. IERS publishes announcements in its *Bulletin C* every 6 months as to whether leap seconds are to occur or not. Such announcements are typically published well in advance of each possible leap-second date—usually in early January for a June 30 scheduled leap second and in early July for a December 31 leap second. Greater levels of advance notice are not possible because of the degree of uncertainty in predicting the precise value of the cumulative effect of fluctuations of the deviation of the Earth’s rotational period from the value of the Mean Solar Day. Or, in other words, the Earth is unpredictably wobbly!

Between 1972 and 2012 some 25 leap seconds have been added to UTC. On average this number implies that a leap second has been inserted about every 19 months. However, the spacing of these leap seconds is quite irregular: there were no leap seconds in the 7-year interval between January 1, 1999, and December 31, 2005, but there were 9 leap seconds in the 13 years between 1985 and 1997, as shown in Figure 1. Since December 31, 1998, there have been only 3 leap seconds, on December 31, 2005, December 31, 2008, and June 30, 2012, each of which has added 1 second to that final minute of the month, at the UTC time of day.

Figure 1: The difference between UT1 and UTC 1984–2012



Leaping Seconds and Computer Systems

The June 30, 2012 leap second did not pass without a hitch, as reported by the tabloid press. The side effect of this particular leap second appeared to include computer system outages and crashes—an outcome that was unexpected and surprising. This leap second managed to crash some servers used in the Amadeus airline management system, throwing the Qantas airline into a flurry of confusion on Sunday morning on July 1 in Australia. But not just the airlines were affected, because LinkedIn, Foursquare, Yelp, and Opera were among numerous online service operators that had their servers stumble in some fashion. This event managed to also affect some *Internet Service Providers* and data center operators. One Australian service provider has reported that a large number of its Ethernet switches seized up over a 2-hour period following the leap second.

It appears that one common element here was the use of the Linux operating system. But Linux is not exactly a new operating system, and the use of the *Leap Second Option* in the *Network Time Protocol* (NTP) [7–10] is not exactly novel either. Why didn't we see the same problems in early 2009, following the leap second that occurred on December 31, 2008?

Ah, but there *were* problems then, but perhaps they were blotted out in the post new year celebratory hangover! Some folks noticed something wrong with their servers on January 1, 2009. Problems with the leap second were recorded with Red Hat Linux following the December 2008 leap second, where kernel versions of the system prior to Version 2.6.9 could encounter a deadlock condition in the kernel while processing the leap second.^[3]

“[...] the leap second code is called from the timer interrupt handler, which holds *xtime_lock*. The leap second code does a *printk* to notify about the leap second. The *printk* code tries to wake up *klogd* (I assume to prioritize kernel messages), and (under some conditions), the scheduler attempts to get the current time, which tries to get *xtime_lock* => *deadlock*.”^[4]

The advice in January 2009 to sysadmins was to upgrade the systems to Version 2.6.9 or later, which contained a patch that avoided this kernel-level deadlock. This time it is a different problem, where the server CPU encountered a 100-percent usage level:

“The problem is caused by a bug in the kernel code for high resolution timers (*hrtimers*). Since they are configured using the *CONFIG_HIGH_RES_TIMERS* option and most systems manufactured in recent years include the *High Precision Event Timers* (HPET) supported by this code, these timers are active in the kernels in many recent distributions.

“The kernel bug means that the *hrtimer* code fails to set the system time when the leap second is added. The result is that the *hrtimer* representation of the time taken from the kernel is a second ahead of the system time. If an application then calls a kernel function with a timeout of less than a second, the kernel assumes that the timeout has elapsed immediately after setting the timer, and so returns to the program code immediately. In the event of a timeout, many programs simply repeat the requested operation and immediately set a new timer. This results in an endless loop, leading to 100% CPU utilisation.”^[5]

Leap Smearing

Following a close monitoring of its systems in the earlier 2005 leap second, Google engineers were aware of problems in their operating system when processing this leap second. They had noticed that some clustered systems stopped accepting work during the leap second of December 31, 2005, and they wanted to ensure that this situation did not recur in 2008. Their approach was subtly different to that used by the Linux kernel maintainers.

Rather than attempt to hunt for bugs in the time management code streams in the system kernel, they noted that the intentional side effect of NTP was to continually perform slight time adjustments in the systems that are synchronizing their time according to the NTP signal. If the quantum of an entire second in a single time update was a problem to their systems, then what about an approach that allowed the 1-second time adjustment to be smeared across numerous minutes or even many hours? That way the leap second would be represented as a larger number of very small time adjustments that, in NTP terms, was nothing exceptional. The result of these changes was that NTP itself would start slowing down the time-of-day clock on these systems some time in advance of the leap second by very slight amounts, so that at the time of the applied leap second, at 23:59:59 UTC, the adjusted NTP time would have already been wound back to 23:59:58. The leap second, which would normally be recorded as 23:59:60 was now a “normal” time of 23:59:59, and whatever bugs that remained in the leap second time code of the system were not exercised.^[6]

More Leaping?

The topic of leap seconds remains a contentious one. In 2005 the United States made a proposal to the *ITU Radiocommunication Sector* (ITU-R) Study Group 7’s Working Party 7-A to eliminate leap seconds. It is not entirely clear whether these leap seconds would be replaced by a less frequent *Leap Hour*, or whether the entire concept of attempting to link UTC and the Mean Solar Day would be allowed to drift, and over time we would see UTC time shifting away from the UT1 concept of solar day time.

This proposal was most recently considered by the ITU-R in January 2012, and there was evidently no clear consensus on this topic. France, Italy, Japan, Mexico, and the United States were reported to be in favor of abandoning leap seconds, whereas Canada, China, Germany, and the United Kingdom were reportedly against these changes to UTC. At present a decision on this topic, or at the least a discussion on this topic, is scheduled for the 2015 *World Radio Conference*.

Although these computing problems with processing leap seconds are annoying and for some folks extremely frustrating and sometimes expensive, I am not sure this factor alone should affect the decision process about whether to drop leap seconds from the UTC time framework. With our increasing dependence on highly available systems, and the criticality of accurate time-of-day clocks as part of the basic mechanisms of system security and integrity, it would be good to think that we have managed to debug this processing of leap seconds.

It is often the case in systems maintenance that the more a bug is exercised the more likely it is that the bug will be isolated and corrected. However, with leap seconds, this task is a tough one because the occurrence of leap seconds is not easily predicted. The next time we have to leap a second in time, about the best we can do is hope that we are ready for it.

For Further Reading

The story of calendars, time, time of day, and time reference standards is a fascinating one. It includes ancient stellar observatories, the medieval quest to predict the date of Easter, the quest to construct an accurate clock that would allow the calculation of longitude, and the current constellations of time and location reference satellites. These days much of this material can be found on the Internet.

- [0] Wikipedia, “Leap Second,”
http://en.wikipedia.org/wiki/Leap_second
- [1] Herald Sun online,
<http://www.heraldsun.com.au/news/leap-second-crashes-qantas-and-leaves-passengers-stranded/story-e6frf7jo-1226413961235>
- [2] “The deviation of the Mean Solar Day from the SI-based day, 1962–2010,” graph in the Wikipedia article referenced earlier^[0],
http://upload.wikimedia.org/wikipedia/commons/thumb/2/28/Deviation_of_day_length_from_SI_day_.svg/1000px-Deviation_of_day_length_from_SI_day_.svg.png

- [3] Red Hat Bugzilla - Bug 479765, “Leap second message can hang the kernel,”
https://bugzilla.redhat.com/show_bug.cgi?id=479765
- [4] “Re: Bug: Status/Summary of slashdot leap-second crash on new years 2008–2009,”
<http://lkml.org/lkml/2009/1/2/373>
- [5] “Leap second bug in Linux wastes electricity,” *The H Open*, July 3, 2012,
<http://www.h-online.com/open/news/item/Leap-second-bug-in-Linux-wastes-electricity-1631462.html>
- [6] “Time, technology and leaping seconds,” Google Official Blog, September 15, 2011,
<http://googleblog.blogspot.de/2011/09/time-technology-and-leaping-seconds.html>
- [7] Burbank, J., Kasch, W., and D. Mills, “Network Time Protocol Version 4: Protocol and Algorithms Specification,” RFC 5905, June 2010.
- [8] Mills, D. and B. Haberman, “Network Time Protocol Version 4: Autokey Specification,” RFC 5906, June 2010.
- [9] Elliott, C., Haberman, B., and H. Gerstung, “Definitions of Managed Objects for Network Time Protocol Version 4 (NTPv4),” RFC 5907, June 2010.
- [10] Lourdelet, B. and R. Gayraud, “Network Time Protocol (NTP) Server Option for DHCPv6,” RFC 5908, June 2010.

Disclaimer

The views expressed are the author’s and not those of APNIC, unless APNIC is specifically identified as the author of the communication. APNIC will not be legally responsible in contract, tort, or otherwise for any statement made in this publication.

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001.
E-mail: gih@apnic.net

The Internet of Things

by David Lake, Ammar Rayes, and Monique Morrow, Cisco Systems

Until a point in time around 2008 or 2009, there were more human beings in the world than devices connected to the Internet. That is no longer the case.

In 2010, the global average of connected devices per person was 1.84. Taking only those people who use the Internet (around 2 billion in 2010), that figure becomes 6 devices per person.^[1] Chip makers such as ARM have targeted developments of low-power CPUs and predicts up to 50 billion devices connected by 2020.^[2]

Today, most of these devices are entities that the user interacts directly with—a PC or Mac, smartphone, tablet, etc. But what is changing is that other devices used every day to orchestrate and manage the world we live in are becoming connected entities in their own right.

They consist not just of users interacting with the end devices—the source and treatment of the information garnered will now occur autonomously, potentially linking to other networks of similarly interconnected entities.

Growing to an estimated 25 billion connected devices by 2015, the rapid explosion of devices on the Internet presents some new and interesting challenges.^[3]

A Definition of the Internet of Things

The *Internet of Things* (IoT) consists of networks of sensors attached to objects and communications devices, providing data that can be analyzed and used to initiate automated actions. The attributes of this world of things may be characterized by low energy consumption, auto-configuration, embeddable objects, etc. The data also generates vital intelligence for planning, management, policy, and decision making. In essence, the five properties that characterize the Internet of Things are as follows:

- *A Unique Internet Address* by which each connected physical object and device will be identified, and therefore be able to communicate with one another.
- *A Unique Location—can be fixed or mobile—within a network or system* (for example, a smart electricity grid) that makes sense of the function and purpose of the object in its specified environment, generating intelligence to enable autonomous actions in line with that purpose.
- *An Increase in Machine-Generated and Machine-Processed Information* that will surpass human-processed information, potentially linking in with other systems to create what some have called “the nervous system of the planet.”

- *Complex New Capabilities in Security, Analytics, and Management*, achievable through more powerful software and processing devices, that enable a network of connected devices and systems to cluster and interoperate transparently in a “network of networks.”
- *Time and Location Achieve New Levels of Importance* in information processing as Internet-connected objects work to generate ambient intelligence; for example, on the *Heating, Ventilation, and Air Conditioning* (HVAC) efficiency of a building, or to study soil samples and climatic change in relation to crop growth.

The concepts and technologies that have led to the IoT, or the interconnectivity of real-world objects, have existed for some time. Many people have referred to *Machine-to-Machine* (M2M) communications and IoT interchangeably and think they are the same. In reality, M2M is only a subset; IoT is a more encompassing phenomenon because it also includes *Machine-to-Human* communication (M2H). *Radio Frequency Identification* (RFID), *Location-Based Services* (LBS), *Lab-on-a-Chip* (LOC) sensors, *Augmented Reality* (AR), robotics, and vehicle telematics are some of the technology innovations that employ both M2M and M2H communications within the IoT as it exists today. They were spun off from earlier military and industrial supply chain applications; their common feature is to combine embedded sensory objects with communication intelligence, running data over a mix of wired and wireless networks.

What has really helped IoT gain traction outside these specific application areas is the greater commoditization of IP as a standard communication protocol, and the advent of IPv6 to allow for a unique IP address for each connected device and object. Researchers and early adopters have been further encouraged by advancements in wireless technologies, including radio and satellite; miniaturization of devices and industrialization; and increasing bandwidth, computing, and storage power.

All these factors have played a part in pushing the boundaries toward generating more context from data capture, communication, and analytics through various devices, objects, and machines in order to better understand our natural and man-made worlds. In exploring the relationship between the IoT and *Information-Centric Networking* (ICN), embedded distributed intelligence will be an important attribute for ICN. Context that is distributed as opposed to centralized is a core architectural component of the IoT for three main reasons:

- *Data Collection*: Centralized data collection and smart object management do not provide the scalability required by the Internet. Managing several hundreds of millions of sensors and actuators in a *Smart Grid* network, for example, cannot be done using a centralized approach.

- *Network Resource Preservation:* Network bandwidth is scarce and some smart objects are not mains-powered, meaning that collecting environmental data from a central point in the network unavoidably leads to using a large amount of the network capacity.
- *Closed-Loop Functioning:* The IoT needs reduced reaction times. For instance, sending an alarm via multiple hops from a sensor to a centralized system, which runs analytics before sending an order to an actuator, would entail unacceptable delays.

Service Management Systems (SMS) (also known as Management Systems, Network Management Systems, or back-end systems) are the brain in the IoT. SMS interacts with intelligent databases that contain *Intellectual Capital* (IC) information, contract information, and manufacturing and historical data. SMS also supports image-recognition technologies to identify objects, people, buildings, places, logos, and anything else that has value to consumers and enterprises. Smartphones and tablets equipped with cameras have pushed this technology from mainly industrial applications to broad consumer and enterprise applications.

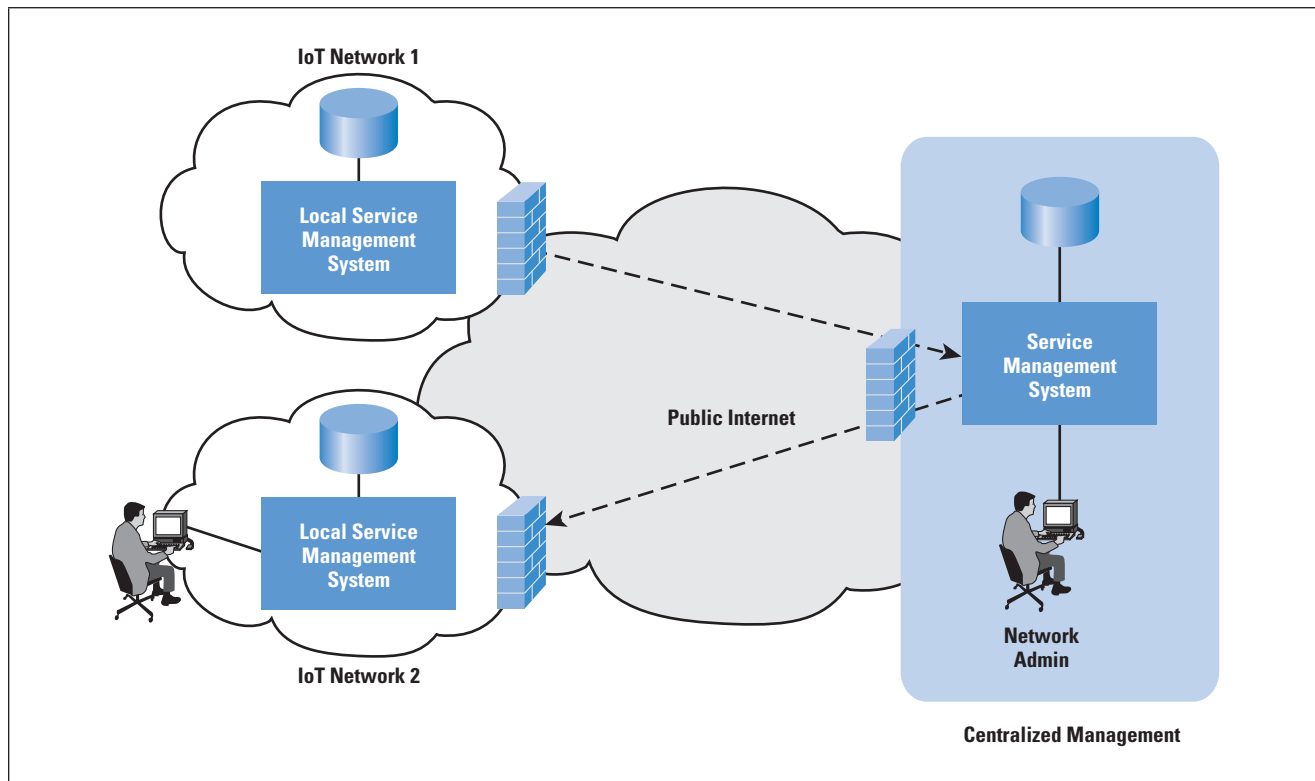
IC information includes intelligence of the vendor's (for example, Cisco) databases and systems such as contract DB, Manufacturing DB, and more importantly thousands of specific roles that are captured over the years by analyzing software bugs, technical support cases, etc.; that is, Cisco knows which devices were manufactured for which customers and with what features. Data collected by the collector is analyzed and correlated with the repository of proprietary Intellectual Capital, turning it into actionable intelligence to help network planners and administrators increase IT value, simplify IT infrastructure, reduce cost, and streamline processes.

Secure communications allow collected data to be sent securely from the agents or collection system to the SMS. SMS includes a database that stores the collected data and algorithms to correlate the collected data with Intellectual Capital information, turning the data into actionable intelligence that network planners and administrators can use with advanced analytics to determine the optimal solution for a problem (or potential problem) after the data is analyzed and corrected. More importantly, a secure mechanism allows the vendor to connect to the network remotely and take action. Secure communications also allows the SMS (automatically or via a network administer) to communicate back with the device to take action when needed.

However, centralized SMS for a large number of entities is very challenging given the near-real-time requirements and the effect on the network performance (see Figure 1). At the same time, centralized intelligence will be required for many IoT networks to interact with back-end centralized databases that are very difficult to distribute (for example, supplier Intellectual Capital databases).

This centralization is more demanding than the traditional multitier environments, servers, and back-end database types of applications where database caching was an effective approach to achieve high scalability and performance. Solution architects need to consider an optimal hybrid model that supports centralized and distributed systems at the same time. Distributed SMS may need to make sub-optimal decisions by using only narrow information to address real-time (or near-real-time) performance problems.

Figure 1: Typical Deployment of an IoT Network



Device and Data Security

The IoT will comprise many small devices, with varying operating systems, CPU types, memory, etc. Many of these devices will be inexpensive, single-function devices—for example, a temperature or pressure sensor—and could have rudimentary network connectivity. In addition, these devices could be in remote or inaccessible locations where human intervention or configuration is impossible.

The nature of sensors is such that they are embedded in what they are sensing—one can envisage a new workplace, hospital, or school construction project where the technology is introduced during the construction phase as part of the final fit rather than after completion as is common today. This paradigm in itself creates new challenges because the means of connectivity may exist only after the installation teams have left the site.

Additionally, methods must be taken to ensure that the authenticity of the data, the path from the sensor to the collector, and the connectivity authentication parameters cannot be compromised between the initial installation or configuration of the device and its eventual presence on the IoT infrastructure.

The challenges of designing and building IoT devices can be summarized as follows:

- IoT devices are typically small, inexpensive devices.
- They are designed to operate autonomously in the field.
- They may be installed prior to network availability.
- After deployment, these devices may require secure remote management.
- The computing platform may not support traditional security algorithms.

Because the IoT will not be a single-use, single-ownership “solution” with sources and the platform on which data may be consumed could be in different ownership, managerial, and connectivity domains, devices will be required to have equal and open access to numerous data consumers concurrently, while still retaining privacy and exclusivity of data where that is required between those consumers.

This requirement was neatly summarized by the IETF Security Area Directors as follows: “A house only needs one toaster even if it serves a family of four!”^[4]

So we have seemingly competing, complex security requirements to be deployed on a platform with limited resources:

- Authenticate to multiple networks securely.
- Ensure that data is available to multiple endpoints.
- Manage the contention between that data access.
- Manage privacy concerns among multiple consumers.
- Provide strong authentication and data protection that cannot be compromised.

And we have to manage existing challenges that all network-attached devices have to contend with such as *Denial of Service* (DoS) attacks, transaction replays, compromised identity through subscriber theft, device theft, or compromised encryption.

These problems have particular relevance in the IoT, where the availability of data is of paramount importance. For example, a critical industrial process may rely on accurate and timely temperature measurement—if that sensor is undergoing a DoS attack, the process collection agent must understand that, and be able to either source data from another location or take evasive action.

It must also be able to distinguish between loss of data because of an ongoing DoS attack and loss of the device because of a catastrophic event in the plant. This ability could mean the difference between a safe shut-down and a major incident.

Authentication and authorization will require reengineering to be appropriate for the IoT. Today's strong encryption and authentication schemes are based on cryptographic suites such as *Advanced Encryption Standard* (AES), *Rivest-Shamir-Adelman* (RSA) for digital signatures and key transport, and *Diffie-Hellman* (DH) for key agreement. Although the protocols are robust, they make very high demands of the compute platform—resources that may not exist in all IoT-attached devices.

These authentication and authorization protocols also require a degree of user intervention in terms of configuration. However, many IoT devices will have limited access; initial configuration needs to be protected from tampering, stealing, and other forms of compromise between device build and install, and also for its usable life, which could be many years.

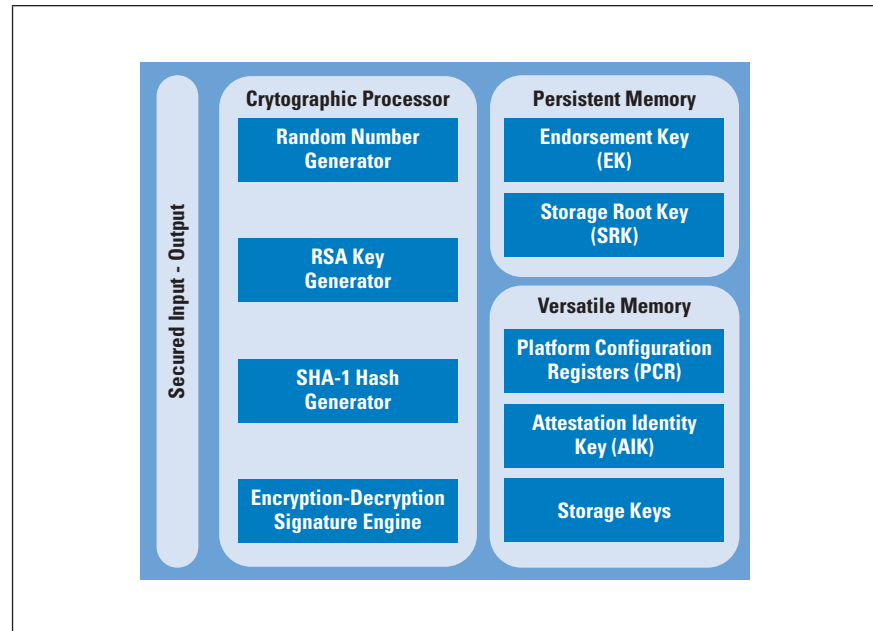
In order to overcome these difficulties, new authentication schemes that allow for strong authentication to many domains while building on the experience of today's strong encryption and authentication algorithms are required.

One possible approach could be to extend methodologies used in the PC industry such as the *Trusted Computing Group's Trusted Platform Model* (TPM).^[5,6]

TPM-enabled devices are fitted at build time with a highly secure hardware device containing a variety of cryptographic elements. Keys and other factors known from this device by trusted third parties are then used in an attestation—a request to validate the authenticity of one device from known parameters.

Because the cryptographic keys are burned into the device during build and the signatures are known to a controlled, trusted third party, a high degree of confidence in the authenticity of the device being queried can be obtained. A typical TPM-compliant cryptographic chip is shown in Figure 2.

Figure 2: Trusted Platform Module



TPM has traditionally been limited by requiring access not only between the devices, but also to a trusted third party. In the IoT, where connectivity may be transient, this requirement is obviously a limitation. Extensions to the TPM to allow for high-confidence attestation between devices without involving a third party have been built; for example, *Direct Anonymous Attestation* (DAA).^[7]

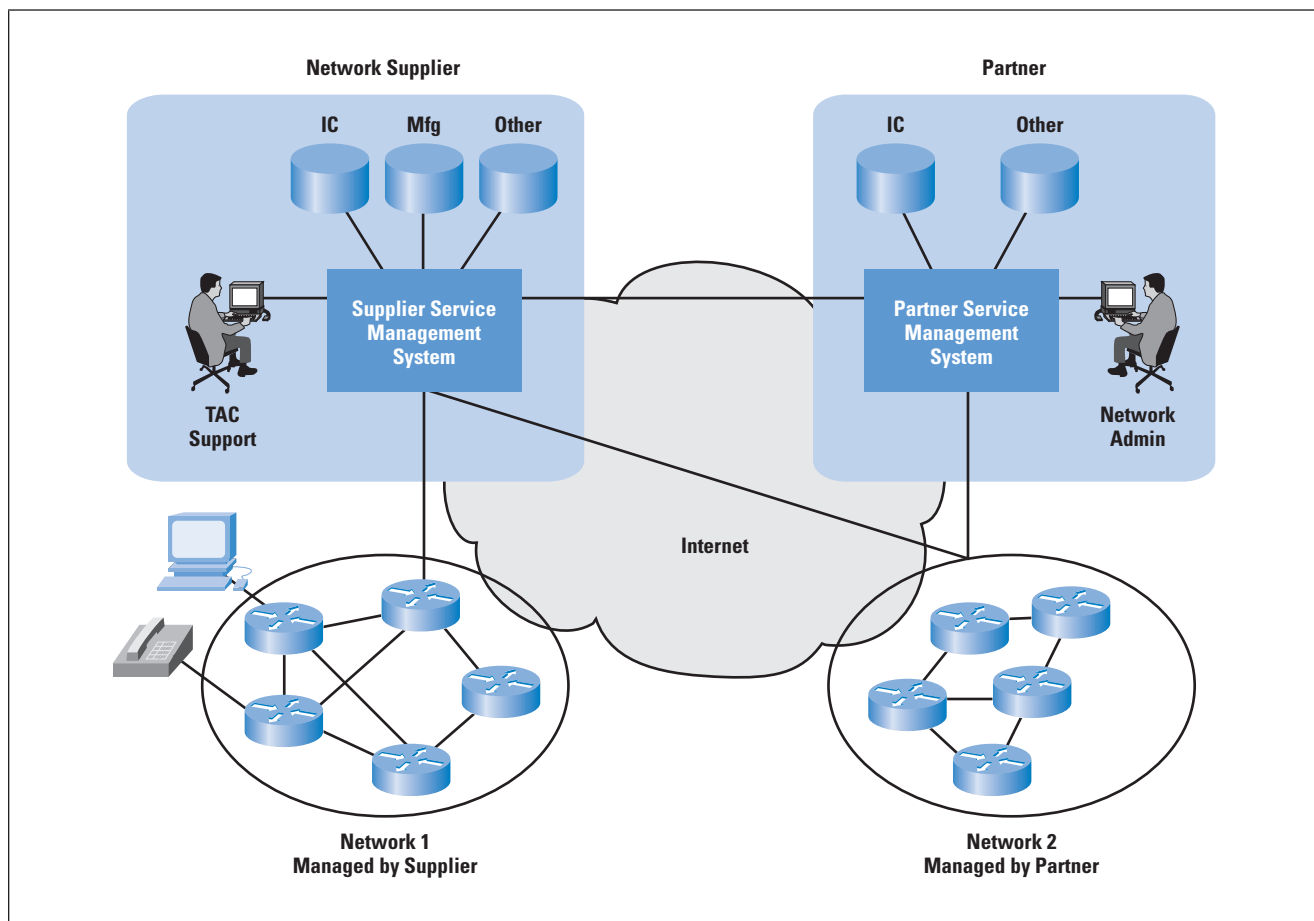
Other elements in security that could be considered include strong authentication between the device and the network attachment point (such as through electrical signatures at the *Media Access Control* [MAC] layer), application of geographic location and privacy levels to data, strengthening of other network-centric methods such as the *Domain Name System* (DNS) and the *Dynamic Host Configuration Protocol* (DHCP) to prevent attacks, and adoption of other protocols that are more tolerant to delay or transient connectivity (such as *Delay Tolerant Networks*).^[8]

An IoT Case Study

The concepts behind the IoT allow management of assets within an enterprise with responsibility shared among customer, partner, and manufacturer in a manner that would previously have been difficult to control.

A typical IT network consists of routers, switches, IP phones, telepresence systems, network management systems such as call managers, data center managers, and many other entities (also known as “machines”) with unique identification (for example, serial number, MAC address, or other address (for example, IP address)). Such a solution is depicted in Figure 3.

Figure 3: Example of Smart Services



The system has the following components:

- **The IT IP-based network:** The network typically is owned by a business customer or an end customer (for example, a small business network). It includes IP devices that may be managed either by the supplier (via service contract), by a third party, Partner 2, or by the customer network administrator.
- **Smart agent or collection system (or sensor):** An external collection system (for example, a server) or smart agent or collection systems on the managed devices gather the device and network information via numerous methods including *Simple Network Management Protocol* (SNMP) requests, *Command-Line Interface* (CLI) commands, syslog, etc. Collected information includes inventory, security data, performance data such as service-level agreement parameters, fault messages, etc.
- **Supplier or partner back-end service management system:** A service management system collects data from various devices and networks, correlates the collected data against intelligent Intellectual Capital rules and important databases (for example, Manufacturing database or Contact Management database), analyzes the results, and produces actionable and trending reports that examine the network and predict the performance.

- Two-way connectivity: Connectivity allows the front-end system (that is, smart agents and collection systems) to send data securely to the supplier or partner service management systems. It also allows the service management system to access the device or network securely to take action when required.
- Secure entitlement and data-transfer capability to register and entitle customer networks and communicate securely (via encryption and security keys) with service providers or network vendors: Such capability is typically deployed on the collector and back-end systems.

A Smart Service provides a proactive intelligence-based solution addressing the installed-based lifecycle and *Fault, Configuration, Accounting, Performance, and Security* (FCAPS) management with the unique benefit of correlating data with the supplier's Intellectual Capital and recognized best practices. Using smart agents, Smart Services collects basic inventory information from the network in order to establish Install Base context.

Conclusions

The implications of the IoT on today's Internet are vast. With such a large number of devices and highly constrained network environments, provisioning and management of the IoT needs to be a part of the architecture. It is both unwise and impractical to provision each active device in the network manually throughout its lifecycle. Earlier technologies, including IP phones, wireless access points, or service provider *Customer Premises Equipment* (CPE), have demonstrated that provisioning can be carried out securely over the network.

The IoT encompasses heterogeneous types of devices that can be on public or private IP networks: from low-powered, low-cost sensors, to fully functioning multipurpose computers with commercial operating systems. For this reason, there can be no "one-size-fits-all" approach to IoT security. What is required is a series of architectural approaches that are dictated by specific IoT use cases. In certain industry solutions, most notably healthcare, security is not just important; information privacy is specifically mandated in many countries.

The challenges of designing, deploying, and supporting billions of IP-enabled endpoints, each producing data that needs to be analyzed and acted on, present exciting opportunities for the next generation of the Internet.

References

- [1] Dave Evans, "The Internet of Things: How the Next Evolution of the Internet Is Changing Everything," April 2011, http://www.cisco.com/web/about/ac79/docs/innov/IoT_IBSG_0411FINAL.pdf

- [2] “ARM targets Internet of Things with new low-power chip,”
Institute of Nanotechnology,
<http://www.instituteofnanotechnology.co.uk/arm-targets-internet-of-things-with-new-low-power-chip/>
- [3] <http://share.cisco.com/internet-of-things.html>
- [4] Tim Polk and Sean Turner, “Security Challenges for The Internet of Things,” IETF Security Area Directors, Feb. 14, 2011,
<http://www.iab.org/wp-content/IAB-uploads/2011/03/Turner.pdf>
- [5] Guillaume Piolle and Yves Demazeau, “Une architecture pour la protection étendue des données personnelles,” 2010,
<http://guillaume.piolle.fr/doc/piolle10b.pdf>
- [6] “Trusted Platform Module,” ISO/IEC 11889,
http://www.iso.org/iso/catalogue_detail.htm?csnumber=50970
- [7] Jan Camenisch, “Direct Anonymous Attestation: Achieving Privacy in Remote Authentication,” June 15, 2004.
<http://www.zurich.ibm.com/security/daa/daa-slides-ZISC.pdf>
- [8] Delay Tolerant Networking Research Group:
<http://www.dtnrg.org/wiki>

DAVID LAKE, B.Sc., is a Consulting Engineer in the Research and Advanced Development Group at Cisco. He has more than 20 years of network design and deployment experience, ranging from X.25 and SNA, through the era of multiprotocol routing to IP, covering a wide range of networking technologies. He has extensive experience in transporting rich-media technologies across complex enterprise and service provider networks. David has worked as customer, network integrator, and manufacturer, and he understands the unique positioning of each of these areas in the IT industry. E-mail: dlake@cisco.com

AMMAR RAYES is a Distinguished Service Engineer at Cisco Systems focusing on Smart Service Technology Strategy. He has authored and co-authored over a hundred papers, patents and books on advances in telecommunications-related technologies. He is the President of the *International Society of Service Innovation Professionals* www.issip.org, an Associate Editor of *ACM Transactions on Internet Technology* and Editor-in-Chief of *Advances of Internet of Things Journal*. Dr. Rayes received his BS and MS Degrees in EE from the University of Illinois at Urbana and his Ph.D. degree in EE from Washington University in St. Louis, Missouri.
E-mail: rayes@cisco.com

MONIQUE MORROW is a Distinguished Engineer in the Research and Advanced Development Group at Cisco. She has over 20 years experience in IP internetworking that includes design, implementation of complex customer projects and service development for service providers. She has presented in various conferences on the topic of RFID, Grid Networking and Cloud Computing; and, she co-authored several books and publications. She is currently focused on the Internet of Things/ Machine-to-Machine Communications in eHealth and security and is active in various SDOS and forums such as the IETF, ITU-T and the FTTH Council Asia-Pacific, holding leadership positions in these organizations. Monique is a Senior Member of the IEEE and a Life Member of the ACM. Monique has a Masters of Science Degree in Telecommunications Management and an MBA.
E-mail: mmorrow@cisco.com

The Demise of Web 2.0 and Why You Should Care

by David Strom

The term Web 2.0 has been around for about a decade^[1], but we are finally seeing its disuse. No, the web itself is not going away, but the notion that an interactive layer of applications, protocols, programming languages, and tools has become subsumed into a new kind of web—one where everything is a *service*, mobile browsing is more important, and social networking has helped discover and promote new content. As a result, we do not really need the term anymore, because it is so much of what the web has become.

Think of this concept as going beyond the 2.0 label of the web: now we have a richer world of interactions that is just the beginning of how we use that tired old TCP port 80. All these developments mean that the readers of *The Internet Protocol Journal* are well poised to help others take advantage of this new complex web environment, because it has become the norm rather than some fancy address in the better part of town. Understanding its new structure and purpose is critical to building the next generation of websites and interactive applications.

Back in the early days of the web in the mid-1990s, it was largely static content that a browser would access from a web server. The notion of having dynamic pages that would automatically update from a database server was exciting and difficult to accomplish without a lot of programming help.

But then came Web 2.0, where the interactive web was born. We had blogging tools such as Google's *Blogger* and Automattic's *Wordpress*, and anyone could create a website that could be easily changed and instantly updated. Web and database servers became better connected, and new protocols were invented to better marry the two.

Everything as a Service

The past few years have seen the rise of *Software as a Service*, *Infrastructure as a Service*, and even *Platforms as a Service*.^[2] The coming of *Cloud Computing* has meant that just about anything can be virtualized and moved into a far-away data center, where it can be managed and replicated easily, obviating the need for any physical infrastructure in the traditional enterprise data center.

Why is this change relevant for the modern web era? Four reasons:

- The web browser is still used as the main remote-access tool to configure and manage a wide variety of applications, network equipment, and servers, including all kinds of cloud-based infrastructures.

- Most of these “as-a-service” entities still run over ports 80 and 443 and piggyback on top of web protocols, for better or worse. We have gotten used to having these ports carry all sorts of traffic that has nothing to do with ordinary web browsing, and we have to do a better job of sorting out the ways apps use the traditional web ports too.
- We do not need to buy any software or install it on our own desktops; everything is available in the cloud at a moment’s notice. What is more, we have gotten used to having the web as the go-to place to get new tools, software drivers, and programs. Software repositories such as *GitHub* and open source projects such as *Apache* have blossomed into places that corporate developers use daily for building their own apps. And why not? They have large support communities and hundreds of projects that are as well tended as something out of Oracle or Microsoft (and some would argue better, too).
- The days of a simple web server serving up pages is ever more complex, with typical commercial websites having ad servers, built-in analytics to track page views and visitors, discussion forums to moderate comments, connections to share the post on *Twitter* and *Facebook* (more on these in a moment), and videos embedded in various ways. All of these websites require coordinated applications and add-ons to the basic web server that require various cloud services. For example, the sites that I run for *ReadWriteWeb* use *Moveable Type* for our content, *Google Analytics*, *Disqus* discussions, interactive polls from **PollDaddy.com**, and custom-built advertising servers, just to name a few of the numerous add-ons. The ever increasing numbers of add-ons means maintaining this system is not easy, and it requires a lot of detailed adjustments on a too-frequent basis.

The Rise of Mobile Browsers

According to the research firm NetApplications^[3], the share of web browsing originating from mobile devices has more than doubled in the past year. Although desktops still account for more than 90 percent of the data accessed from browsers, mobile devices are consuming the web at an increasing rate.

Part of this trend is that we are using more devices and they have become more capable. Android-based phones constitute the largest market share, and they have the fastest-growing consumer mobile phone adoption rate.^[4] Certainly, more and more of us are browsing more webpages from mobile devices these days.

Another part of the trend of increased roaming on mobile devices is that more people are creating and using more mobile apps, too. Hundreds of new mobile apps with a wide variety of content are created every day. Professors at major universities teach computer science students how to code mobile apps, and you can even take online courses on *Java* programming.

But mobile browsing poses a conundrum for web designers. One school of thought is to build custom tablet applications for your website, to show off the features of the tablet interface and to make it easier for tablet users to interact with your content. The U.K. *Guardian*, for example, is leading the way in this area.^[5]

Another school of thought is to improve the mobile experience, by either building a separate site that is optimized for smaller screens and lower bandwidth connections or allowing the site to work automatically under the constraints of the mobile browser itself.^[6]

One real challenge for the mobile web browsing experience is the role of Adobe Flash and the newest of the *Hypertext Markup Language* (HTML) standards, HTMLv5. Apple decided when it released its first iPads to not support Flash, and since then there has been additional effort and movement to migrate many Flash-based sites, such as **YouTube.com**, toward HTMLv5, which is supported by Apple's tablets and can be more efficient for lower-bandwidth connections. Although this topic could easily be the subject of an entire article for this journal, our point in mentioning it here is that displaying video and similar content is still a problem for the web, even today.

Our mobile traffic at *ReadWriteWeb* has increased tremendously in the past year, and I suspect our site is typical of other sites. But this increase in traffic presents challenges for content creators: is it better to sell ad units around the content, even ads that have sub-par browsing experiences on mobile devices? Or code up your own iPad app (or use Verve's tools [<http://www.vervewireless.com/>] or something equivalent)? Certainly the level of engagement with the custom mobile app is greater, but it amazes me that sites with just static pages still are not optimized for mobile browsers yet, with large image downloads or multiple included links, for example.

Let's consider the site **Remodelista.com** as a case study of how to properly optimize a site for mobile browsing. The owners have implemented tricks to adjust its layout for different screen sizes. As you make your browsing window smaller (or as you run it on a mobile device with a small screen), the integrity of the site content remains intact, meaning that font sizes change and ad blocks appear on wider, higher-resolution screens and disappear on smaller ones, but the overall content stream remains the same, no matter what device is used to view it. This consistency is achieved by adding a lot of special coding to the webpages, as the following snippet shows:

```
<!--[if IEMobile 7]> <html class="no-js iem7 oldie" itemscope itemtype="http://schema.org/"><![endif]-->
<!--[if lt IE 7]> <html lang="en" class="no-js ie6 oldie" xmlns="http://www.w3.org/1999/xhtml"
xmlns:nectar="http://saymedia.com/2011/swml" itemscope itemtype="http://schema.org/"><![endif]-->
<!--[if (IE 7)&!(IEMobile)]> <html lang="en" class="no-js ie7 oldie" xmlns="http://www.w3.org/1999/xhtml"
xmlns:nectar="http://saymedia.com/2011/swml" itemscope itemtype="http://schema.org/"><![endif]-->
<!--[if (IE 8)&!(IEMobile)]> <html lang="en" class="no-js ie8 oldie" xmlns="http://www.w3.org/1999/xhtml"
xmlns:nectar="http://saymedia.com/2011/swml" itemscope itemtype="http://schema.org/"><![endif]-->
<!--[if gt IE 8]> <html class="no-js" lang="en" itemscope itemtype="http://schema.org/"><![endif]-->
<!--[if (gte IE 9)|(gt IEMobile 7)]> <html class="no-js" lang="en" itemscope itemtype="http://schema.org/">
<![endif]-->
```


The Social Web Is Now Everywhere

It used to be the odd person in your professional circle who did not have or use an Internet e-mail account. Now the odd person is the one who does not have an account on Facebook or some other social networking site. What began in a Harvard dorm room in this decade has turned into a juggernaut of more than a billion users—and it is growing rapidly.

But the social web is more than a bunch of college kids swapping photos of their party pictures. A recent study from the University of Massachusetts at Dartmouth^[7] shows that nearly 75 percent of the Inc. 500 (the fastest-growing 500 American private companies) are using *Facebook* or *LinkedIn*, a level that is about twice the percentage that are using corporate blogs. “Ninety percent of responding executives report that social media tools are important for brand awareness and company reputation. Eighty-eight percent see these tools as important for generating web traffic while 81% find them important for lead generation. Seventy-three percent say that social media tools are important for customer support programs.” Clearly, these tools have become the accepted corporate intranet, the mainstream mechanism for communications among distributed work teams, and the way that many of us share events in our professional lives as well.

The social web means more than a “Like” button on a particular page of content; it is a way to curate and disseminate that content quickly and easily. It has replaced the Usenet *news groups* that many of us remember with a certain fondness for their arcane and complex structure. Or maybe that is just nostalgia talking.

In the presocial web past, even in the days when Web 2.0 was the rage, sharing and curation was not easy. If you wanted to share something you found online, more than likely you would e-mail your colleagues a URL. Now you can *Tweet*, post on *Facebook* and *Google+*, add an update to your *LinkedIn* account, put up a page on your corporate **Yammer.com** or **tibbr.com** server, or use one of dozens more services that will stream your likes and notable sites to the world at large. Or you likely have to do all of these tasks.

Back in the days of yore (say 2000), when I wrote a freelance article, it was sufficient to post a link to the story on my own personal website, in addition to perhaps sending an e-mail message or two to the people I thought might be interested in reading the content. Those days seem so quaint. Today, the process of writing the article is actually just the beginning, not the end. When the article appears online, a whole series of promotional activities must take place, including monitoring online discussions and adding my own comments, posting on the various social media sites, and re-Tweeting a link to my article several times over the next several days—all to ensure generation of lots of traffic.

There are even services such as **Ping.fm** and **Graspr.com** that can coordinate batch updates to numerous services, so that at the push of a button all of your social media will get your news at once. Or services such as **Nimble.com** that attempt to coordinate your entire social graph (as it is called) of friends and admirers so you can track what is going out across all your various networks.

Where We Go from Here

I have just tried to touch on a few topics to show that the days of the simple static web are “so over,” as Generation Y says. Clearly, we have a long and rich future ahead of us for more interesting web applications.

References

- [1] http://en.wikipedia.org/wiki/Web_2.0
- [2] See “Alphabet Soup in the Cloud”:
<http://www.readwriteweb.com/cloud/2011/10/alphabet-soup-in-the-cloud-und.php>
- [3] NetApplications research cited in this September 2011 article in *Computerworld*:
http://www.computerworld.com/s/article/9219696/Apples_rules_phone_tablet_browsing_market
- [4] Nielsen’s statistics are typical:
<http://blog.nielsen.com/nielsenwire/?p=29786>
- [5] See <http://m.guardian.co.uk/>, but you really need to view it on an iPad or other tablet device to understand what they are trying to do with their content. See also:
http://www.readwriteweb.com/archives/the_guardian_ipad_edition_hits_ios_5_newsstands.php
- [6] See Thomas Husson’s May 2011 Forrester Research blog post here:
http://blogs.forrester.com/thomas_husson/11-05-03-why_the_web_versus_application_debate_is_irrelevant

Also see my own January 2012 article in ReadWriteWeb here:
<http://www.readwriteweb.com/hack/2012/01/do-you-really-need-your-own-mo.php>
- [7] See their January 2012 study here:
<http://www.umassd.edu/cmr/studiesandresearch/2011inc500socialmediaupdate/>

DAVID STROM has created dozens of editorial-rich websites for publications such as *ReadWriteWeb*, *Tom’sHardware.com*, *eeTimes*, and others, as well as written thousands of articles for numerous IT magazines. He was the founding editor-in-chief of *Network Computing* magazine and author of two books on computer networking. He lives in St. Louis, Mo. and can be found at **strominator.com**, on Twitter **@dstrom**, and **david@strom.com** for those that still prefer e-mail.

Binary Floor Control Protocol

by Pat Jensen, Cisco Systems

Over the last decade, communication technologies have evolved to encompass new modalities of collaboration across IP networks—from instant messaging on a personal computer, to being able to make *Voice-over-IP* (VoIP) calls and also now including the growing adoption of *High-Definition* (HD) videoconferencing.

Operating systems, device types, and physical locations now are less affected as continued growth in networking has evolved to promote high bandwidth across wireless and wired networks. An example is the emergence of growing network-access technologies such as *Multiprotocol Label Switching* (MPLS), *Very-High-Speed Digital Subscriber Line 2* (VDSL2), *Long Term Evolution* (LTE), and *Data over Cable Service Interface Specification* (DOCSIS). With both availability of bandwidth and broadband user penetration increasing, the user's expectation of delivering immersive collaboration now becomes more apparent.

This evolution includes modern use cases accelerating the adoption of videoconferencing, such as enabling telemedicine for remote surgeries and diagnostic procedures as well as distance learning applications being used to connect educators with students across the globe.

This article introduces the *Binary Floor Control Protocol* (BFCP) as a standard for managing floor control during collaboration sessions across dedicated video endpoints, mobile devices, and personal computers running collaboration software. These capabilities can be delivered using an enabled *Session Initiation Protocol* (SIP) standards-based endpoint or as a software implementation in a collaboration application stack.

History

BFCP is a deliverable developed as part of the *Internet Engineering Task Force* (IETF) XCON Centralized Conferencing working group. The IETF XCON working group was formed to focus on delivering a standards-based approach to managing IP conferencing while promoting broad interoperability between software and equipment vendors.^[1]

This mandate includes defining the objects, mechanisms, and provisions to assist in scheduling conferencing resources. These resources could be consumed as a conference enabled in a web browser, via an audio conference call or during a videoconference.

As defined, privacy, security, and authorization are considered integral in protecting the ability to join, participate in, and manage each conference session. The IETF XCON working group's initial focus was on unicast media conferences.

The IETF XCON working group was proposed in August 2003, with work starting early in October of that year.^[2] Early requirements for BFCP were defined in RFC 4376, which describes important concepts, including a model for floor control and how it should be integrated in a conferencing platform.^[3] Other important aspects such as security, including using authentication and encryption to provide protection against man-in-the-middle attacks, were also outlined.

In November 2006, Gonzalo Camarillo, Joerg Ott, and Keith Drage authored RFC 4582, which defined the Binary Floor Control Protocol.^[4]

Besides BFCP, other standardization efforts around conference role and content management also were defined, including the ITU-T H.239 recommendation.^[5] Unlike BFCP, H.239 applies specifically to H.323-enabled *Integrated Services Digital Network* (ISDN) and IP conferencing endpoints, whereas BFCP is designed to be agnostic of the underlying signaling protocol.

Protocol Details

The basic concept of floor control is analogous to managing a live in-person presentation, where you want to control who is presenting, manage and transition your presenters, and maintain a feedback loop. Also important is the ability to allow a presenter to show slides and share with your audience a white board or transparency projector.

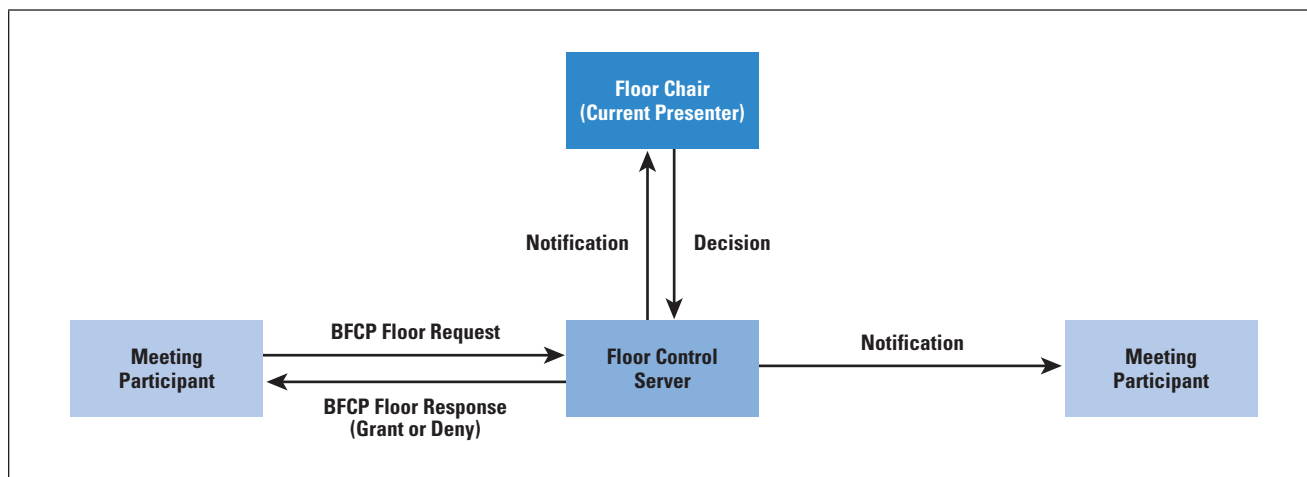
During an active collaboration session, a presenter may choose to present material to a remote user, or optionally to an audience on a call with multiple endpoints through a *Multipoint Control Unit* (MCU). This session could include many additional sources; for example, using a secondary video camera to show zoomed-in content (that is, an optical examination camera used in telemedicine) or any external video source.

This floor-control mechanism can also encompass functions available in a collaboration application stack, such as the ability to share the content of the presenter's desktop, application, or web browser.

BFCP provides the ability to manage multiple streams being presented during a collaboration session using floor control. BFCP accomplishes this management using a token-based mechanism where a single presenter can request control of the floor from the floor-control server.

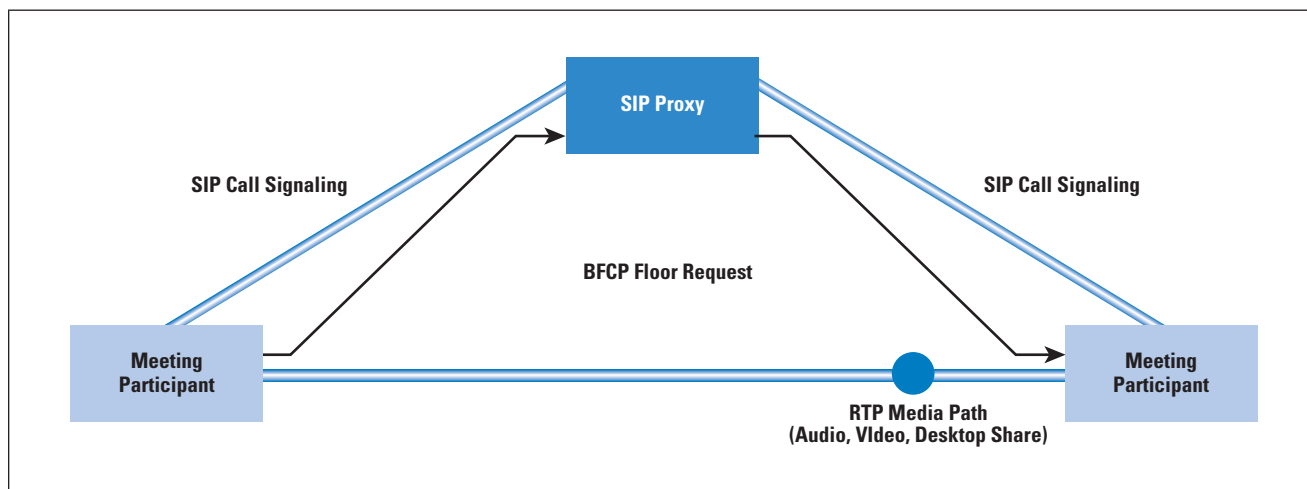
When this request is granted, the presenter holds the token and has the ability to open an additional stream to provide presentation data. Figure 1 examines this process in detail, with a meeting attendee requesting the token from the floor-control server to become an active presenter during the session.

Figure 1: BFCP Floor Request from Floor-Control Server



This same interaction can also take place during a point-to-point audio or video call with only two parties. In this case, a token can be used to signify which party will be presenting an additional stream, such as a secondary camera or application providing a desktop sharing session. Figure 2 shows an overview of this process. One of the critical differences here is that in a point-to-point call, the floor-control server capability is being provided by the user's device or application instead of using a multipoint control unit or conference server.

Figure 2: BFCP Floor Request in a Point-to-Point SIP Call



For instance, as a presenter, you can choose to present auxiliary streams via your application or endpoint and determine whether it is your primary, secondary, or tertiary stream. As a conference participant, you can also choose which stream you are currently viewing, also including the definition and quality of the secondary stream. In this case, current network conditions such as bandwidth and latency will also dictate the quality of additional streams.

BFCP is designed to be signaling protocol-agnostic, in that it is relying on the capabilities of the underlying signaling and transport protocols to set up each stream that is being managed, including whether voice, video, or content is being provided in the *Real-Time Transport Protocol* (RTP) stream.

For example, using a standards-based endpoint and *Session Initiation Protocol* (SIP), a SIP INVITE message is sent with the media capabilities line specifying the session description information about the stream. This data provides relevant information about the underlying video codec being used and the bit rate that is required to support the video and presentation streams.

In this case as multiple RTP media streams are transported across the network carrying audio and video traffic, *Call Admission Control* (CAC) and *Quality of Service* (QoS) tagging can be applied and enforced by the call-control platform, providing the ability to limit bandwidth usage and helping ensure that bandwidth is available on the network after the additional media stream is added.

Also important to note, BFCP can use *Transport Layer Security* (TLS) to provide encryption of floor information pertaining to each resource that is being controlled as well as the participants using and viewing them. BFCP provides the ability to support anonymous users as well for sessions where you may have a large audience or where anonymity is desired. An example of where this feature could be used is hosting a large web conferencing event where you have external attendees who may be outside of your organization.

One use case for BFCP includes the ability to focus on the presenter while the presenter is sharing a desktop application. With the ability to control the presenter's media stream, this feature adds additional immersion in a collaboration session, allowing you to both identify the presenter's visual cues and posture as well as focus on relevant content the presenter supplies.

Summary

The Binary Floor Control Protocol plays a very important role in helping manage diverse types of content being shared across multiple parties in a conference session. Today's modern implementations of BFCP span web conferencing applications as well as video and audio conferencing solutions across a wide array of vendors.

While these vendors are focused on delivering these capabilities across screen-led PC-centric types of devices, because of its inherent transport-agnostic capabilities, it is likely we will see BFCP being used to enable new modalities of content sharing across collaboration applications in the future.

Industry efforts are focusing on promoting collaboration applications across new arrays of devices, including using touchscreen technology on handheld computers and stationary LCD televisions to manipulate and visualize data in new ways.

Concepts such as manipulating session content using cognitive mapping as an evolution of electronic whiteboarding and transitioning an active conference from a tablet device to another type of room-based video-enabled endpoint during a collaboration session are two powerful examples of ways BFCP could be used in the future. On the horizon, touchscreen-enabled tablet and smartphone devices and HTML5-enabled web browsers also provide yet another avenue to enable rich standards-based multimedia conferencing with advanced content management.

Disclaimer

The views of this article do not necessarily represent the views or positions of Cisco Systems.

For Further Reading

- [1] <http://datatracker.ietf.org/wg/xcon/charter/>
- [2] <http://datatracker.ietf.org/wg/xcon/history/>
- [3] Petri Koskelainen, Joerg Ott, Henning Schulzrinne, and Xiaotao Wu, "Requirements for Floor Control Protocols," RFC 4376, February 2006.
- [4] Gonzalo Camarillo, Joerg Ott, and Keith Drage, "The Binary Floor Control Protocol," RFC 4582, November 2006.
- [5] "Role management and additional media channels for H.300-series terminals," International Telecommunication Union Standard H.239, September 2005.

PAT JENSEN is a member of the Unified Communications consulting systems engineering team at Cisco Systems. Since 2010, he has designed collaboration architectures for Cisco's customers across the western United States. Prior to joining Cisco Systems, he served as an IP telephony design engineer designing and implementing unified communications and telepresence technologies at AT&T and SBC DataComm. E-mail, XMPP, SIP: patjense@cisco.com



© Stonehouse Photography/Internet Society

Pierre Ouedraogo Receives 2012 Jonathan B. Postel Service Award

The Internet Society recently announced that its prestigious *Jonathan B. Postel Service Award* was presented to Pierre Ouedraogo for his exceptional contributions to the growth and vitality of the Internet in Africa. The international award committee, comprised of former Jonathan B. Postel award winners, noted that Mr. Ouedraogo played a significant role in the growth of the Internet in Africa and demonstrated an extraordinary commitment to training young engineers and participating in regional Internet organizations.

Mr. Ouedraogo is the Director of Digital Francophonie at *Organisation Internationale de la Francophonie* (OIF) based in Paris, France. Over the years, he has established networks of IT experts to coordinate African efforts to develop IT and use it as a tool for development. Mr. Ouedraogo initiated many IT technical workshops in Africa and is a founding member of numerous African regional organizations, including AfriNIC (the African Internet Registry for IP addresses); AfTLD (*African Internet Top Level Domain Names Association*); AFNOG (*African Network Operators Group*); AfCERT (*African CERT network*), and AfrICANN (*African network of participants to the ICANN process*).

“Pierre Ouedraogo is a highly-regarded technical leader in Africa, and he has been instrumental in bringing the Internet to Burkina Faso as well as other French-speaking African countries,” said Lynn St. Amour, President and Chief Executive Officer of the Internet Society.”

“His commitment to the expansion of the Internet and encouragement of young engineers to help them build their skills through training workshops has had a profound impact on the growth of the Internet across Africa.”

The Postel Award was established by the Internet Society to honour individuals or organisations that, like Jon Postel, have made outstanding contributions in service to the data communications community. The committee places particular emphasis on candidates who have supported and enabled others in addition to their own specific actions. The award is focused on sustained and substantial technical contributions, service to the community, and leadership.

For more information about the Internet Society and the Postel award, see: <http://www.internetsociety.org/>

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



Atsushi Seike (L) with Vint Cerf and Jun Murai.

Vint Cerf Awarded Honorary Doctorate by Keio University

Keio University in Tokyo recently awarded Dr. Vinton Gray Cerf an honorary doctorate in Media and Governance for his work in the creation and governance of our modern Internet over the last forty years. On the recommendation of Professor Jun Murai, dean of the Faculty of Environment and Information Studies, Keio University president Atsushi Seike presented Dr. Cerf with the degree. The ceremony was held in the Enzetsu-kan, the historic public speaking hall on Keio's Mita Campus in Tokyo, and streamed live via the Internet to viewers around the world.

Professor Murai's recommendation for the degree, read during the ceremony, said that not only is Dr. Cerf the founding father of internetworking technology, "he is the global leader in many ways of the largest innovation for the 21st century, the Internet itself, which has become the core of today's information-based society." In addition to his work on TCP/IP with Robert Khan, Dr. Cerf's work in establishing the Internet Society and his stewardship of ICANN as its chairman were highlighted. Also mentioned was his role in *Delay/Disruption-Tolerant Networking* (DTN) and the first experiments connecting a space probe twenty million miles away using Internet protocols.

In his remarks, President Seike mentioned Dr. Cerf's forty-year commitment to advancing the role of networks in creating our global society, from the earliest days of the ARPANET through today's Internet. "[Dr. Cerf] understood quickly and clearly the international nature of the Internet and its potential for having a positive impact on the lives of not just the technical elite, but for all of the people of the world, as a tool for education, commerce, and the advance of democracy," he noted. Professor Seike compared Dr. Cerf's role in using technology to make the world a better place to the efforts of Yukichi Fukuzawa, the founder of Keio University, who in the mid-19th century was instrumental in bringing knowledge to Japan from the outside world, not as an academic exercise but in order to improve society.

Following the ceremony, Dr. Cerf gave an invited technical talk titled "Re-Inventing the Internet." He discussed the potential of DTN and *Mobile Ad Hoc Networks* as tools for disaster recovery. He presented his view of urgent technical problems, including the need for strong authentication and digital forensics. He also outlined society's need for preserving data, the programs that create and manipulate that data, and even the systems that are used to run those programs. Without such an effort, we will fail to preserve our own technical and cultural history for the thousands of years we have come to expect, he noted.

Dr. Cerf left behind the inscription, "I cannot imagine a greater honor than to be brought into this august and highly regarded university where contrary thinking is rewarded! I am most grateful to my good friend, Jun Murai, for his decades long commitment to the Internet."



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2012 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

December 2012

Volume 15, Number 4

A Quarterly Technical Publication for
Internet and Intranet Professionals

In This Issue

From the Editor	1
Network Time Protocol.....	2
Packet Classification	12
Fragments	23
Call for Papers	31

FROM THE EDITOR

Accurate timekeeping has long been an engineering challenge if not obsession in some circles. Take for example the iconic Swiss *chronometer* watch or the pendulum-controlled clock mechanism in London's Palace of Westminster, often referred to as "Big Ben." Such mechanical systems—accurate as they may be—are no match for the clocks we use in telecommunication and computer networks. In our last issue, Geoff Huston described the glitches encountered last June when a *Leap Second* was applied to *Coordinated Universal Time* (UTC). In this issue he explains the operation of the *Network Time Protocol* (NTP). The article is another installment in our series "Protocol Basics."

It is difficult to believe that it has been more than 25 years since the first publication of Douglas Comer's book series *Internetworking With TCP/IP*. Volume 1 of this series will soon be available in its sixth edition, and we asked the author to write an article about *Packet Classification* based on material in the book.

The recent *World Conference on International Telecommunications* (WCIT) did not have the outcome with respect to the Internet that many had hoped for. We plan to publish an analysis of this event in our next issue. This time—in our "Fragments" section—we have some reactions from the *Number Resource Organization* (NRO) and the Internet Society, as well as pointers to further information about WCIT.

January 1, 2013, marked the 30th anniversary of the *Transmission Control Protocol/Internet Protocol* (TCP/IP). A transition from the earlier *Network Control Program* (NCP) took place on January 1, 1983, also known as "Flag Day." Such an instant technology change would have been desirable for the transition from IPv4 to IPv6, but sadly this isn't possible. Instead we are happy to honor those who dedicate their careers to IPv6 deployment with an *Itojun Service Award*. See page 25 for more details.

On page 30 you will find some frequently asked questions about subscriptions to this journal. If you have other questions or comments, please contact us at ipj@cisco.com

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

Protocol Basics: The Network Time Protocol

by Geoff Huston, APNIC

Back at the end of June 2012^[0] there was a brief IT hiccup as the world adjusted the *Coordinated Universal Time* (UTC) standard by adding an extra second to the last minute of the 31st of June. Normally such an adjustment would pass unnoticed by all but a small dedicated collection of time keepers, but this time the story spread out into the popular media as numerous Linux systems hiccupped over this additional second, and they supported some high-profile services, including a major air carrier's reservation and ticketing backend system. The entire topic of time, time standards, and the difficulty of keeping a highly stable and regular clock standard in sync with a slightly wobbly rotating Earth has been a longstanding debate in the *International Telecommunication Union Radiocommunication Sector* (ITU-R) standards body that oversees this coordinated time standard. However, I am not sure that anyone would argue that the challenges of synchronizing a strict time signal with a less than perfectly rotating planet is sufficient reason to discard the concept of a coordinated time standard and just let each computer system drift away on its own concept of time. These days we have become used to a world that operates on a consistent time standard, and we have become used to our computers operating at sub-second accuracy. But how do they do so? In this article I will look at how a consistent time standard is spread across the Internet, and examine the operation of the *Network Time Protocol* (NTP).

Some communications protocols in the IP protocol suite are quite recent, whereas others have a long and rich history that extends back to the start of the Internet. The ARPANET switched over to use the TCP/IP protocol suite in January 1983, and by 1985 NTP was in operation on the network. Indeed it has been asserted that NTP is the longest running, continuously operating, distributed application on the Internet^[1].

The objective of NTP is simple: to allow a client to synchronize its clock with UTC time, and to do so with a high degree of accuracy and a high degree of stability. Within the scope of a WAN, NTP will provide an accuracy of small numbers of milliseconds. As the network scope gets finer, the accuracy of NTP can increase, allowing for sub-millisecond accuracy on LANs and sub-microsecond accuracy when using a precision time source such as a *Global Positioning System* (GPS) receiver or a caesium oscillator.

If a collection of clients all use NTP, then this set of clients can operate with a synchronized clock signal. A shared data model, where the modification time of the data is of critical importance, is one example of the use of NTP in a networked context.

(I have relied on NTP timer accuracy at the microsecond level when trying to combine numerous discrete data sources, such as a web log on a server combined with a *Domain Name System* (DNS) query log from DNS resolvers and a packet trace.)

NTP, Time, and Timekeeping

To consider NTP, it is necessary to consider the topic of timekeeping itself. It is useful to introduce some timekeeping terms at this juncture:

<i>Stability</i>	How well a clock can maintain a constant frequency
<i>Accuracy</i>	How well the frequency and absolute value of the clock compares with a standard reference time
<i>Precision</i>	How well the accuracy of a clock can be maintained within a particular timekeeping system
<i>Offset</i>	The time difference in the absolute time of two clocks
<i>Skew</i>	The variation of offset over time (first-order derivative of offset over time)
<i>Drift</i>	The variation of skew over time (second-order derivative of offset over time)

NTP is designed to allow a computer to be aware of three critical metrics for timekeeping: the *offset* of the local clock to a selected reference clock, the *round-trip delay* of the network path between the local computer and the selected reference clock server, and the *dispersion* of the local clock, which is a measure of the maximum error of the local clock relative to the reference clock. Each of these components is maintained separately in NTP. They provide not only precision measurements of offset and delay, to allow the local clock to be adjusted to synchronize with a reference clock signal, but also definitive maximum error bounds of the synchronization process, so that the user interface can determine not only the time, but the quality of the time as well.

Universal Time Standards

It would be reasonable to expect that the time is just the time, but that is not the case. The Universal Time reference standard has several versions, but these two standards are of interest to network timekeeping.

UT1 is the principal form of Universal Time. Although conceptually it is *Mean Solar Time* at 0° longitude, precise measurements of the Sun are difficult. Hence, it is computed from observations of distant quasars using long baseline interferometry, laser ranging of the Moon and artificial satellites, as well as the determination of GPS satellite orbits. *UT1* is the same everywhere on Earth, and is proportional to the rotation angle of the Earth with respect to distant quasars, specifically the *International Celestial Reference Frame* (ICRF), neglecting some small adjustments.

The observations allow the determination of a measure of the Earth's angle with respect to the ICRF, called the *Earth Rotation Angle* (ERA), which serves as a modern replacement for *Greenwich Mean Sidereal Time*). UT1 is required to follow the relationship

$$\text{ERA} = 2\pi(0.7790572732640 + 1.00273781191135448T_u) \text{ radians}$$

where $T_u = (\text{Julian UT1 date} - 2451545.0)$

Coordinated Universal Time (UTC) is an atomic timescale that approximates UT1. It is the international standard on which civil time is based. It ticks SI seconds, in step with *International Atomic Time* (TAI). It usually has 86,400 SI seconds per day, but is kept within 0.9 seconds of UT1 by the introduction of occasional intercalary leap seconds. As of 2012 these leaps have always been positive, with a day of 86,401 seconds.^[9]

NTP uses UTC, as distinct from the *Greenwich Mean Time* (GMT), as the reference clock standard. UTC uses the TAI time standard, based on the measurement of 1 second as 9,192,631,770 periods of the radiation emitted by a caesium-133 atom in the transition between the two hyperfine levels of its ground state, implying that, like UTC itself, NTP has to incorporate leap second adjustments from time to time.

NTP is an “absolute” time protocol, so that local time zones—and conversion of the absolute time to a calendar date and time with reference to a particular location on the Earth's surface—are not an intrinsic part of the NTP protocol. This conversion from UTC to the wall-clock time, namely the local date and time, is left to the local host.

Servers and Clients

NTP uses the concepts of *server* and *client*. A server is a source of time information, and a client is a system that is attempting to synchronize its clock to a server.

Servers can be either a *primary server* or a *secondary server*. A primary server (sometimes also referred to as a *stratum 1* server using terminology borrowed from the time reference architecture of the telephone network) is a server that receives a UTC time signal directly from an authoritative clock source, such as a configured atomic clock or—very commonly these days—a GPS signal source. A *secondary server* receives its time signal from one or more *upstream servers*, and distributes its time signal to one or more *downstream servers* and *clients*. Secondary servers can be thought of as clock signal repeaters, and their role is to relieve the client query load from the primary servers while still being able to provide their clients with a clock signal of comparable quality to that of the primary servers. The secondary servers need to be arranged in a strict hierarchy in terms of upstream and downstream, and the stratum terminology is often used to assist in this process.

As noted previously, a stratum 1 server receives its time signal from a UTC reference source. A stratum 2 server receives its time signal from a stratum 1 server, a stratum 3 server from stratum 2 servers, and so on. A stratum n server can peer with many stratum $n - 1$ servers in order to maintain a reference clock signal. This stratum framework is used to avoid synchronization loops within a set of time servers.

Clients peer with servers in order to synchronize their internal clocks to the NTP time signal.

The NTP Protocol

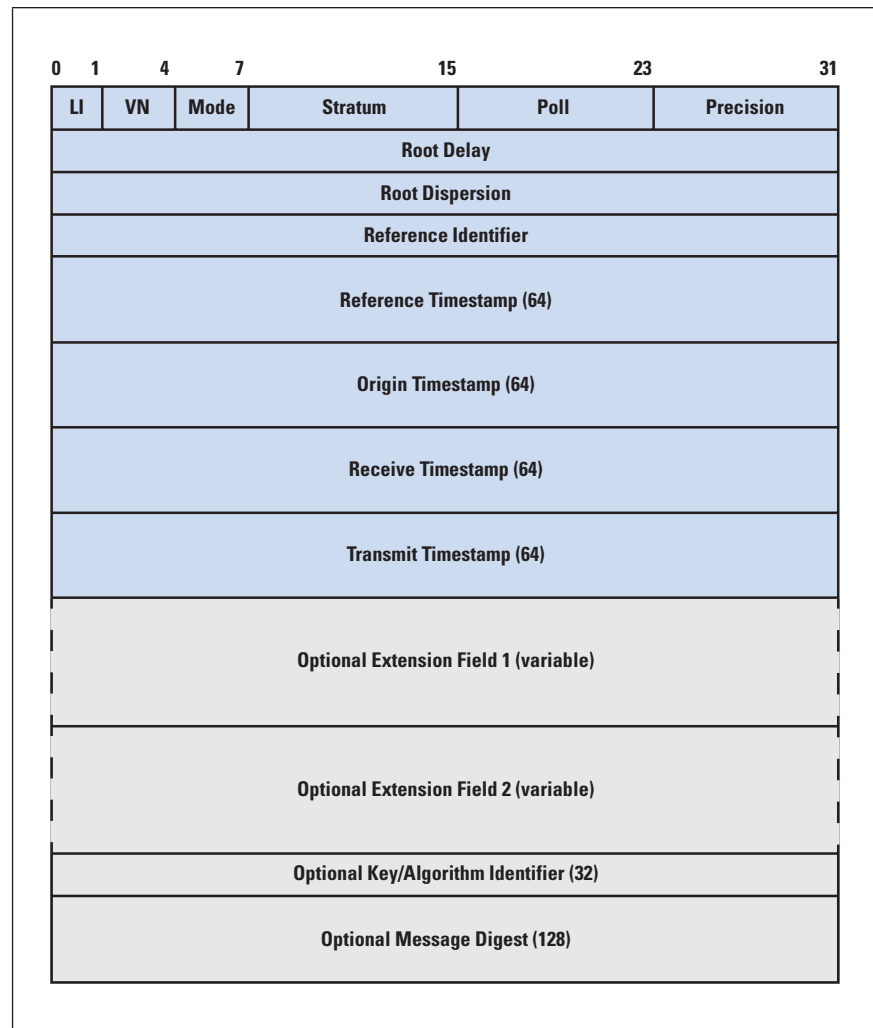
At its most basic, the NTP protocol is a clock request transaction, where a client requests the current time from a server, passing its own time with the request. The server adds its time to the data packet and passes the packet back to the client. When the client receives the packet, the client can derive two essential pieces of information: the reference time at the server and the elapsed time, as measured by the local clock, for a signal to pass from the client to the server and back again. Repeated iterations of this procedure allow the local client to remove the effects of network jitter and thereby gain a stable value for the delay between the local clock and the reference clock standard at the server. This value can then be used to adjust the local clock so that it is synchronized with the server. Further iterations of this protocol exchange can allow the local client to continuously correct the local clock to address local clock skew.

NTP operates over the *User Datagram Protocol* (UDP). An NTP server listens for client NTP packets on port 123. The NTP server is stateless and responds to each received client NTP packet in a simple transactional manner by adding fields to the received packet and passing the packet back to the original sender, without reference to preceding NTP transactions.

Upon receipt of a client NTP packet, the receiver time-stamps receipt of the packet as soon as possible within the packet assembly logic of the server. The packet is then passed to the NTP server process. This process interchanges the IP Header Address and Port fields in the packet, overwrites numerous fields in the NTP packet with local clock values, time-stamps the egress of the packet, recalculates the checksum, and sends the packet back to the client.

The NTP packets sent by the client to the server and the responses from the server to the client use a common format, as shown in Figure 1.

Figure 1: NTP Message Format



The header fields of the NTP message are as follows:

- LI** Leap Indicator (2 bits)
This field indicates whether the last minute of the current day is to have a leap second applied. The field values follow:
0: No leap second adjustment
1: Last minute of the day has 61 seconds
2: Last minute of the day has 59 seconds
3: Clock is unsynchronized
- VN** NTP Version Number (3 bits) (current version is 4).

<i>Mode</i>	<p>NTP packet mode (3 bits)</p> <p>The values of the Mode field follow:</p> <ul style="list-style-type: none"> 0: Reserved 1: Symmetric active 2: Symmetric passive 3: Client 4: Server 5: Broadcast 6: NTP control message 7: Reserved for private use
<i>Stratum</i>	<p>Stratum level of the time source (8 bits)</p> <p>The values of the Stratum field follow:</p> <ul style="list-style-type: none"> 0: Unspecified or invalid 1: Primary server 2–15: Secondary server 16: Unsynchronized 17–255: Reserved
<i>Poll</i>	<p>Poll interval (8-bit signed integer)</p> <p>The \log_2 value of the maximum interval between successive NTP messages, in seconds.</p>
<i>Precision</i>	<p>Clock precision (8-bit signed integer)</p> <p>The precision of the system clock, in \log_2 seconds.</p>
<i>Root Delay</i>	<p>The total round-trip delay from the server to the primary reference sourced. The value is a 32-bit signed fixed-point number in units of seconds, with the fraction point between bits 15 and 16. This field is significant only in server messages.</p>
<i>Root Dispersion</i>	<p>The maximum error due to clock frequency tolerance. The value is a 32-bit signed fixed-point number in units of seconds, with the fraction point between bits 15 and 16. This field is significant only in server messages.</p>
<i>Reference Identifier</i>	<p>For stratum 1 servers this value is a four-character ASCII code that describes the external reference source (refer to Figure 2). For secondary servers this value is the 32-bit IPv4 address of the synchronization source, or the first 32 bits of the <i>Message Digest Algorithm 5</i> (MD5) hash of the IPv6 address of the synchronization source.</p>

Figure 2: Reference Identifier Codes
(from RFC 4330)

Code	External Reference Source
LOCL	uncalibrated local clock
CESM	calibrated Cesium clock
RBDM	calibrated Rubidium clock
PPS	calibrated quartz clock or other pulse-per-second source
IRIG	Inter-Range Instrumentation Group
ACTS	NIST telephone modem service
USNO	USNO telephone modem service
PTB	PTB (Germany) telephone modem service
TDF	Allouis (France) Radio 164 kHz
DCF	Mainflingen (Germany) Radio 77.5 kHz
MSF	Rugby (UK) Radio 60 kHz
WWV	Ft. Collins (US) Radio 2.5, 5, 10, 15, 20 MHz
WWVB	Boulder (US) Radio 60 kHz
WWVH	Kauai Hawaii (US) Radio 2.5, 5, 10, 15 MHz
CHU	Ottawa (Canada) Radio 3330, 7335, 14670 kHz
LORC	LORAN-C radionavigation system
OMEG	OMEGA radionavigation system
GPS	Global Positioning Service

The next four fields use a 64-bit time-stamp value. This value is an unsigned 32-bit seconds value, and a 32-bit fractional part. In this notation the value 2.5 would be represented by the 64-bit string:

0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0010 . | 1000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000

The unit of time is in seconds, and the epoch is 1 January 1900, meaning that the NTP time will cycle in the year 2036 (two years before the 32-bit Unix time cycle event in 2038).

The smallest time fraction that can be represented in this format is 232 picoseconds.

<i>Reference Timestamp</i>	This field is the time the system clock was last set or corrected, in 64-bit time-stamp format.
<i>Originate Timestamp</i>	This value is the time at which the request departed the client for the server, in 64-bit time-stamp format.
<i>Receive Timestamp</i>	This value is the time at which the client request arrived at the server in 64-bit time-stamp format.
<i>Transmit Timestamp</i>	This value is the time at which the server reply departed the server, in 64-bit time-stamp format.

The basic operation of the protocol is that a client sends a packet to a server and records the time the packet left the client in the *Origin Timestamp* field (T1). The server records the time the packet was received (T2). A response packet is then assembled with the original Origin Timestamp and the *Receive Timestamp* equal to the packet receive time, and then the *Transmit Timestamp* is set to the time that the message is passed back toward the client (T3). The client then records the time the packet arrived (T4), giving the client four time measurements, as shown in Figure 3.

Figure 3: NTP Transaction
Timestamps (from RFC 4330)

Timestamp Name	ID	When Generated
Originate Timestamp	T1	time request sent by client
Receive Timestamp	T2	time request received by server
Transmit Timestamp	T3	time reply sent by server
Destination Timestamp	T4	time reply received by client

These four parameters are passed into the client timekeeping function to drive the clock synchronization function, which we will look at in the next section.

The optional Key and Message Digest fields allow a client and a server to share a secret 128-bit key, and use this shared secret to generate a 128-bit MD5 hash of the key and the NTP message fields. This construct allows a client to detect attempts to inject false responses from a man-in-the-middle attack.

The final part of this overview of the protocol operation is the polling frequency algorithm. A NTP client will send a message at regular intervals to a NTP server. This regular interval is commonly set to be 16 seconds. If the server is unreachable, NTP will back off from this polling rate, doubling the back-off time at each unsuccessful poll attempt to a minimum poll rate of 1 poll attempt every 36 hours. When NTP is attempting to resynchronize with a server, it will increase its polling frequency and send a burst of eight packets spaced at 2-second intervals.

When the client clock is operating within a sufficient small offset from the server clock, NTP lengthens the polling interval and sends the eight-packet burst every 4 to 8 minutes (or 256 to 512 seconds).

Timekeeping on the Client

The next part of the operation of NTP is how an NTP process on a client uses the information generated by the periodic polls to a server to moderate the local clock.

From an NTP poll transaction, the client can estimate the delay between the client and the server. Using the time fields described in Figure 3, the transmission delay can be calculated as the total time from transmission of the poll to reception of the response minus the recorded time for the server to process the poll and generate a response:

$$\delta = (T4 - T1) - (T3 - T2)$$

The offset of the client clock from the server clock can also be estimated by the following:

$$\Theta = \frac{1}{2} [(T2 - T1) + (T3 - T4)]$$

It should be noted that this calculation assumes that the network path delay from the client to the server is the same as the path delay from the server to the client.

NTP uses the minimum of the last eight delay measurements as δ_0 . The selected offset, Θ_0 , is one measured at the lowest delay. The values (Θ_0, δ_0) become the NTP update value.

When a client is configured with a single server, the client clock is adjusted by a slew operation to bring the offset with the server clock to zero, as long as the server offset value is within an acceptable range.

When a client is configured with numerous servers, the client will use a selection algorithm to select the preferred server to synchronize against from among the candidate servers. Clustering of the time signals is performed to reject outlier servers, and then the algorithm selects the server with the lowest stratum with minimal offset and jitter values. The algorithm used by NTP to perform this operation is *Marzullo's Algorithm*^[2].

When NTP is configured on a client, it attempts to keep the client clock synchronized against the reference time standard. To do this task NTP conventionally adjusts the local time by small offsets (larger offsets may cause side effects on running applications, as has been found when processing leap seconds). This small adjustment is undertaken by an *adjtime()* system call, which slews the clock by altering the frequency of the software clock until the time correction is achieved. Slewing the clock is a slow process for large time offsets; a typical slew rate is 0.5 ms per second.

Obviously this informal description has taken a rather complex algorithm and some rather detailed math formulas without addressing the details. If you are interested in how NTP operates at a more detailed level, consult the references that follow, which will take you far deeper into the algorithms and the underlying models of clock selection and synchronization than I have done here.

Conclusion

NTP is in essence an extremely simple stateless transaction protocol that provides a quite surprising outcome. From a regular exchange of simple clock readings between a client and a server, it is possible for the client to train its clock to maintain a high degree of precision despite the possibility of potential problems in the stability and accuracy of the local clock and despite the fact that this time synchronization is occurring over network paths that impose a noise element in the form of jitter in the packet exchange between client and server. Much of today's distributed Internet service infrastructure relies on a common time base, and this base is provided by the common use of the Network Time Protocol.

References and Further Reading

- [0] Geoff Huston, “Leaping Seconds,” *The Internet Protocol Journal*, Volume 15, No. 3, September 2012.
- [1] David L. Mills, “A Brief History of NTP Time: Confessions of an Internet Timekeeper,” ACM SIGCOMM, *Computer Communication Review*, Vol. 33, No. 2, pp. 9–12, April 2003, <http://www.eecis.udel.edu/~mills/database/papers/history.pdf>
- [2] K. A. Marzullo, “Maintaining the Time in a Distributed System: An Example of a Loosely-Coupled Distributed Service,” Ph.D. dissertation, Stanford University, Department of Electrical Engineering, February 1984, http://en.wikipedia.org/wiki/Marzullo%27s_algorithm
- [3] David L. Mills, “NTP Architecture, Protocol and Algorithms,” University of Delaware, www.eecis.udel.edu/~mills/database/brief/arch/arch.ppt
- [4] Jack Burbank, William Kasch, and David Mills, “Network Time Protocol Version 4: Protocol and Algorithms Specification,” RFC 5905, June 2010.
- [5] David L. Mills, “Simple Network Time Protocol (SNTP) Version 4 for IPv4, IPv6 and OSI,” RFC 4330, January 2006.
- [6] <http://www.ntp.org>
- [7] <http://www.eecis.udel.edu/~mills/ntp.html>
- [8] David Mills, *Computer Network Time Synchronization: the Network Time Protocol on Earth and in Space*, Second Edition, CRC Press, 2011.
- [9] http://en.wikipedia.org/wiki/Universal_Time

Disclaimer

The views expressed are the author’s and not those of APNIC, unless APNIC is specifically identified as the author of the communication. APNIC will not be legally responsible in contract, tort or otherwise for any statement made in this publication.

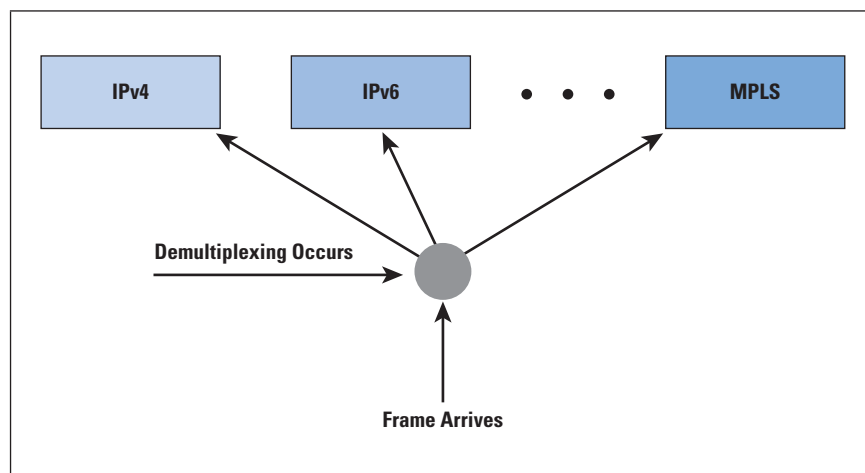
GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001.
E-mail: gih@apnic.net

Packet Classification: A Faster, More General Alternative to Demultiplexing

by Douglas Comer, Purdue University

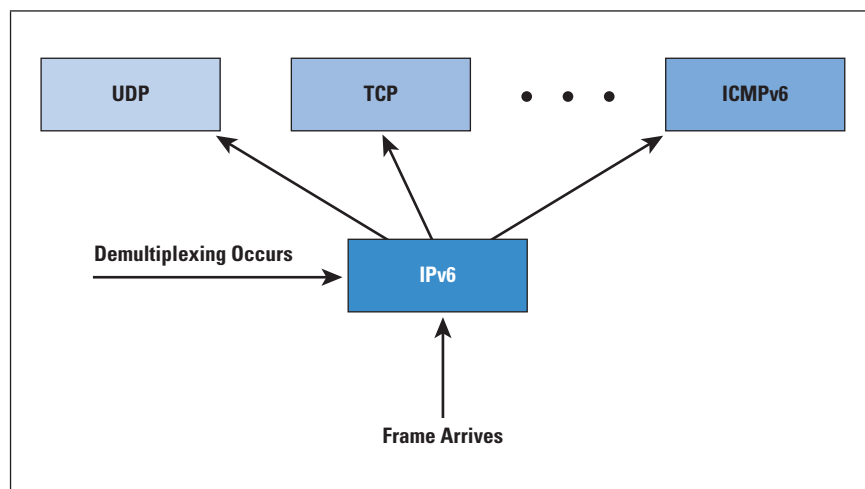
Traditional packet-processing systems use an approach known as *demultiplexing* to handle incoming packets (refer to [1] for details). When a packet arrives, protocol software uses the contents of a *Type Field* in a protocol header to decide how to process the payload in the packet. For example, the Type field in a frame is used to select a Layer 3 module to handle the frame, as Figure 1 illustrates.

Figure 1: Frame Demultiplexing



Demultiplexing is repeated at each level of the protocol stack. For example, IPv6 uses the *Next Header* field to select the correct transport layer protocol module, as Figure 2 illustrates.

Figure 2: Demultiplexing at Layer 3



Modern, high-speed network systems take an entirely different view of packet processing. In place of demultiplexing, they use a technique known as *classification*^[2]. Instead of assuming that a packet proceeds through a protocol stack one layer at a time, they allow processing to cross layers. (In addition to being used by companies such as Cisco and Juniper, classification has been used in Linux^[3] and with network processors by companies such as Intel and Netronome^[4].)

Packet classification is especially pertinent to three key network technologies. First, Ethernet switches use classification instead of demultiplexing when they choose how to forward packets. Second, a router that sends incoming packets over *Multiprotocol Label Switching* (MPLS) tunnels uses classification to choose the appropriate tunnel. Third, classification provides the basis for *Software-Defined Networking* (SDN) and the *OpenFlow* protocol.

Motivation for Classification

To understand the motivation for classification, consider a network system that has protocol software arranged in a traditional layered stack. Packet processing relies on demultiplexing at each layer of the protocol stack. When a frame arrives, protocol software looks at the Type field to learn about the contents of the frame payload. If the frame carries an IP datagram, the payload is sent to the IP protocol module for processing. IP uses the destination address to select a next-hop address. If the datagram is in *transit* (that is, passing through the router on its way to a destination), IP forwards the datagram by sending it back out one of the interfaces. A datagram reaches TCP only if the datagram is destined for the router itself. TCP then uses the protocol port numbers in the TCP segment to further demultiplex the incoming datagram among multiple application programs.

To understand why traditional layering does not solve all problems, consider MPLS processing. In particular, consider a router at the border between a traditional internet and an MPLS core. Such a router must accept packets that arrive from the traditional internet and choose an MPLS path over which to send the packet. Why is layering pertinent to path selection? In many cases, network managers use transport layer protocol port numbers when choosing a path. For example, suppose a manager wants to send all web traffic down a specific MPLS path. All the web traffic will use TCP port 80, meaning that the selection must examine TCP port numbers.

Unfortunately, in a traditional demultiplexing scheme, a datagram does not reach the transport layer unless the datagram is destined for the local network system. Therefore, protocol software must be reorganized to handle MPLS path selection. We can summarize:

A traditional protocol stack is insufficient for the task of MPLS path selection because path selection often involves transport layer information and a traditional stack will not send transit datagrams to the transport layer.

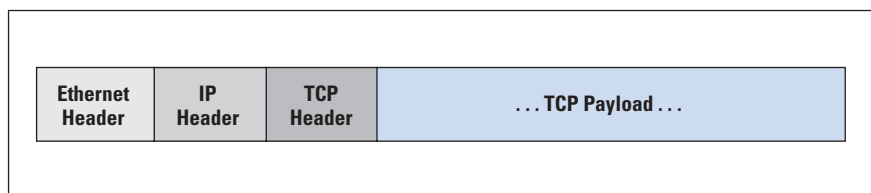
Classification Instead of Demultiplexing

How should protocol software be structured to handle tasks such as MPLS path selection? The answer lies in the use of *classification*. A classification system differs from conventional demultiplexing in two ways:

- Ability to cross multiple layers
- Higher speed than demultiplexing

To understand classification, imagine a packet that has been received at a router and placed in memory. *Encapsulation* means that the packet will have a set of contiguous protocol headers at the beginning. For example, Figure 3 illustrates the headers in a TCP packet (for example, a request sent to a web server) that has arrived over an Ethernet.

Figure 3: Layout of a Packet in Memory



Given a packet in memory, how can we quickly determine whether the packet is destined to the web? A simplistic approach simply looks at one field in the headers: the TCP destination port number. However, it could be that the packet is not a TCP packet at all. Maybe the frame is carrying *Address Resolution Protocol* (ARP) data instead of IP. Or maybe the frame does indeed contain an IP datagram, but instead of TCP the transport layer protocol is the *User Datagram Protocol* (UDP). To make certain that it is destined for the web, software needs to verify each of the headers: the frame contains an IP datagram, the IP datagram contains a TCP segment, and the TCP segment is destined for the web.

Instead of parsing protocol headers, think of the packet as an array of octets in memory. Consider IPv4 as an example. To be an IPv4 datagram, the Ethernet Type field (located in array positions 12 and 13) must contain **0x0800**. The IPv4 Protocol field, located at position 23, must contain **6** (the protocol number for TCP). The Destination Port field in the TCP header must contain **80**. To know the exact position of the TCP header, we must know the size of the IP header. Therefore, we check the header length octet of the IPv4 header. If the octet contains **0x45**, the TCP destination port number will be found in array positions 36 and 37.

As another example, consider classifying *Voice over IP* (VoIP) traffic that uses the *Real-Time Transport Protocol* (RTP). Because RTP is not assigned a specific UDP port, vendors use a heuristic to determine whether a given packet carries RTP traffic: check the Ethernet and IP headers to verify that the packet carries UDP, and then examine the octets at a known offset in the RTP packet to verify that the value matches the value used by a known codec.

Observe that all the checks described in the preceding paragraphs require only array lookup. That is, the lookup mechanism treats the packet as an array of octets and merely checks to verify that location *X* contains value *Y*, location *Z* contains value *W*, and so on—the mechanism does not need to understand any of the protocol headers or the meaning of values. Furthermore, observe that the lookup scheme crosses multiple layers of the protocol stack.

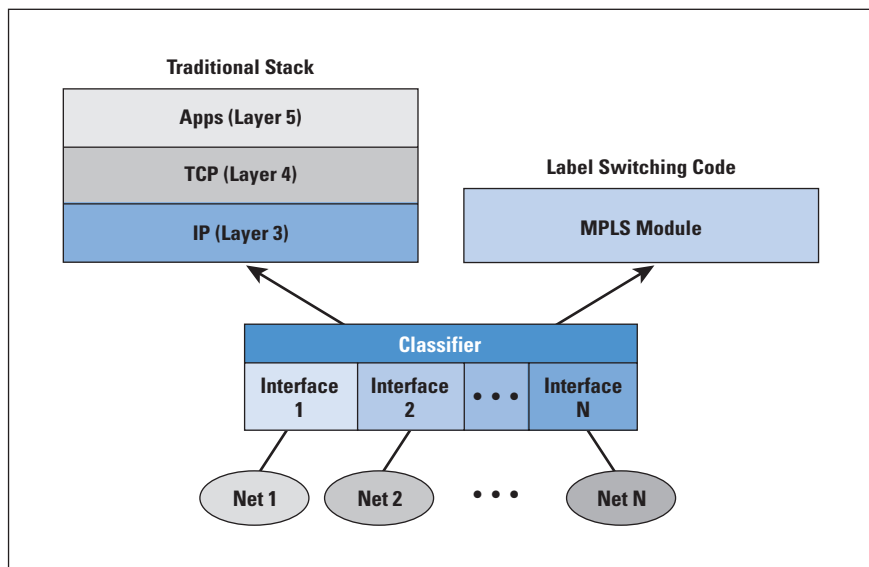
We use the term *classifier* to describe a mechanism that uses the lookup approach described previously, and we say that the result is a packet *classification*. In practice, a classification mechanism usually takes a list of classification *rules* and applies them until a match is found. For example, a manager might specify three rules: send all web traffic to MPLS path 1, send all FTP traffic to MPLS path 2, and send all VPN traffic to MPLS path 3.

Layering When Classification Is Used

If classification crosses protocol layers, how does it relate to traditional layering diagrams? We can think of classification as an extra layer that has been squeezed between Layer 2 and Layer 3. When a packet arrives, the packet passes from a Layer 2 module to the classification module. All packets proceed to the classifier; no demultiplexing occurs before classification. If any of the classification rules matches the packet, the classification layer follows the rule. Otherwise, the packet proceeds up the traditional protocol stack. For example, Figure 4 illustrates layering when classification is used to send some packets across MPLS paths.

Interestingly, a classification layer can subsume all demultiplexing. That is, instead of classifying packets only for MPLS paths, the classifier can be configured with additional rules that check the Type field in a frame for IPv4, IPv6, ARP, *Reverse ARP* (RARP), and so on.

Figure 4: Layering in a Router that Uses Classification to Select MPLS Paths



Classification Hardware and Network Switches

The text in the previous section describes a classification mechanism that is implemented in software—an extra layer is added to a software protocol stack that classifies frames after they arrive at a router. Classification can also be implemented in hardware. In particular, Ethernet switches and other packet-processing hardware devices contain classification hardware that allows packet classification and forwarding to proceed at high speed. The next sections explain hardware classification mechanisms.

We think of network devices, such as switches, as being divided into broad categories by the level of protocol headers they examine and the consequent level of functions they provide:

- Layer 2 Switching
- Layer 2 *Virtual Local-Area Network* (VLAN) Switching
- Layer 3 Switching
- Layer 4 Switching

A *Layer 2 Switch* examines the *Media Access Control* (MAC) source address in each incoming frame to learn the MAC address of the computer that is attached to each port. When a switch learns the MAC addresses of all the attached computers, the switch can use the destination MAC address in each frame to make a forwarding decision. If the frame is unicast, the switch sends only one copy of the frame on the port to which the specified computer is attached. For a frame destined to the broadcast or a multicast address, the switch delivers a copy of the frame to all ports.

A *VLAN Switch* adds one level of virtualization by permitting a manager to assign each port to a specific VLAN. Internally, VLAN switches extend forwarding in a minor way: instead of sending broadcasts and multicasts to all ports on the switch, a VLAN switch consults the VLAN configuration and sends them only to ports on the same VLAN as the source.

A *Layer 3 Switch* acts like a combination of a VLAN switch and a router. Instead of using only the Ethernet header when forwarding a frame, the switch can look at fields in the IP header. In particular, the switch watches the source IP address in incoming packets to learn the IP address of the computer attached to each switch port. The switch can then use the IP destination address in a packet to forward the packet to its correct destination.

A *Layer 4 Device* extends the examination of a packet to the transport layer. That is, the device can include the TCP or UDP Source and Destination Port fields when making a forwarding decision.

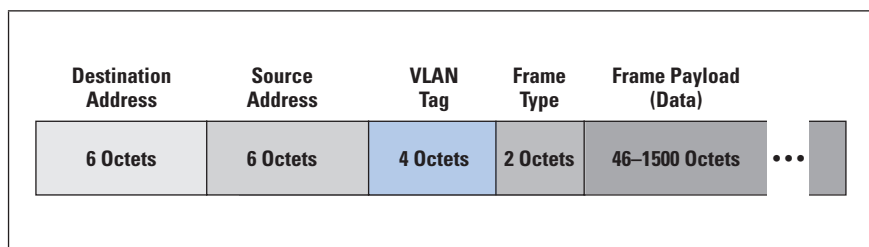
Switching Decisions and VLAN Tags

All types of switching hardware described previously use classification. That is, switches operate on packets as if a packet is merely an array of octets, and individual fields in the packet are specified by giving offsets in the array. Thus, instead of demultiplexing packets, a switch treats a packet syntactically by applying a set of classification rules similar to the rules described previously.

Surprisingly, even VLAN processing is handled in a syntactic manner. Instead of merely keeping VLAN information in a separate data structure that holds meta information, the switch inserts an extra field in an incoming packet and places the VLAN number of the packet in the extra field. Because it is just another field, the classifier can reference the VLAN number just like any other header field.

We use the term *VLAN Tag* to refer to the extra field inserted in a packet. The tag contains the VLAN number that the manager assigned to the port over which the frame arrived. For Ethernet, IEEE standard 802.1Q specifies placing the VLAN Tag field after the MAC Source Address field. Figure 5 illustrates the format.

Figure 5: An Ethernet Frame with a VLAN Tag Inserted



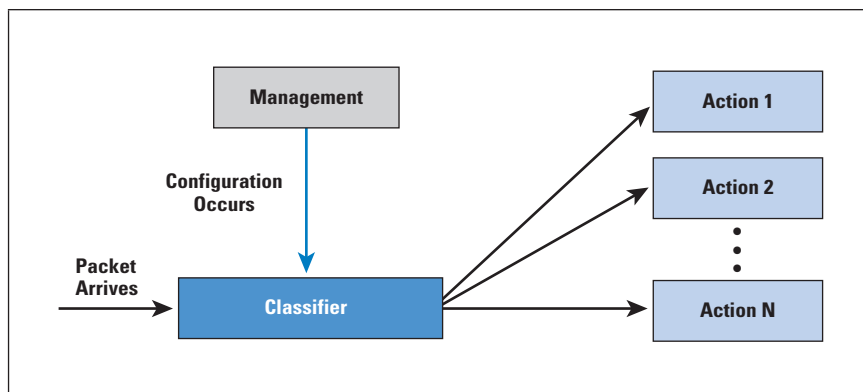
A VLAN tag is used only internally—after the switch has selected an output port and is ready to transmit the frame, the tag is removed. Thus, when computers send and receive frames, the frames do not contain a VLAN tag.

An exception can be made to the rule: a manager can configure one or more ports on a switch to leave VLAN tags in frames when sending the frame. The purpose is to allow two or more switches to be configured to operate as a single, large switch. That is, the switches can share a set of VLANs—a manager can configure each VLAN to include ports on one or both of the switches.

Classification Hardware

We can think of hardware in a switch as being divided into three main components: a classifier, a set of units that perform actions, and a management component that controls the overall operation. Figure 6 illustrates the overall organization and the flow of packets.

Figure 6: Hardware Components
Used for Classification



As black arrows in the figure indicate, the classifier provides the high-speed data path that packets follow. When a packet arrives, the classifier uses the rules that have been configured to choose an action. The management module usually consists of a general-purpose processor that runs management software. A network administrator can interact with the management module to configure the switch, in which case the management module can create or modify the set of rules the classifier follows.

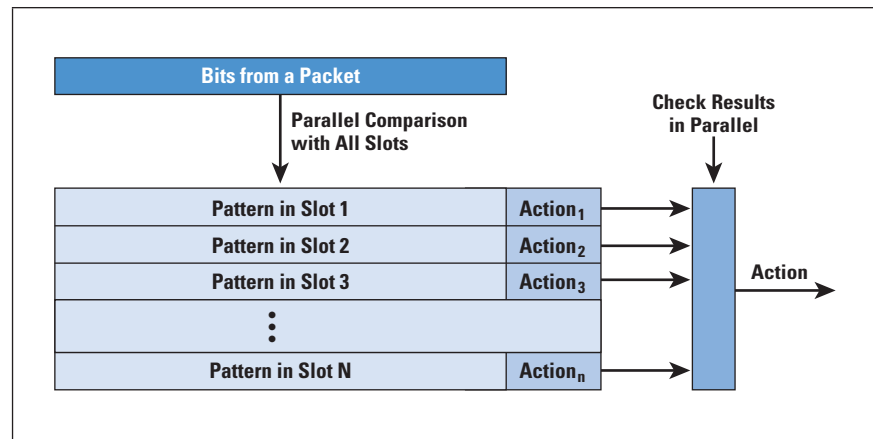
A network system, such as a switch, must be able to handle two types of traffic: transit traffic and traffic destined for the switch itself. For example, to provide management or routing functions, a switch may have a local TCP/IP protocol stack and packets destined for the switch must be passed to the local stack. Therefore, one of the actions a classifier takes may be “*pass packet to the local stack for Demultiplexing*”.

High-Speed Classification and TCAM

Modern switches can allow each interface to operate at 10 Gbps. At 10 Gbps, a frame takes only 1.2 microseconds to arrive, and a switch usually has many interfaces. A conventional processor cannot handle classification at such speeds, so a question arises: how can a hardware classifier achieve high speed? The answer lies in a hardware technology known as *Ternary Content Addressable Memory* (TCAM).

TCAM uses parallelism to achieve high speed—instead of testing one field of a packet at a given time, TCAM checks all fields simultaneously. Furthermore, TCAM performs multiple checks at the same time. To understand how TCAM works, think of a packet as a string of bits. We imagine TCAM hardware as having two parts: one part holds the bits from a packet and the other part is an array of values that will be compared to the packet. Entries in the array are known as *slots*. Figure 7 illustrates the idea.

Figure 7: The Conceptual Organization of TCAM



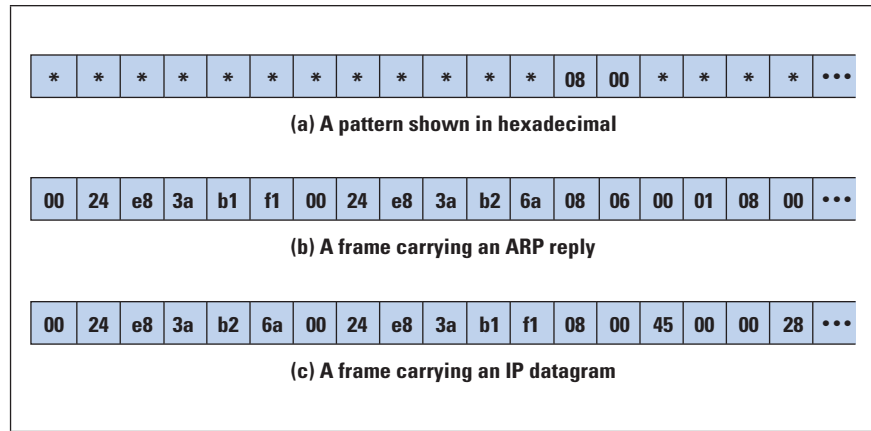
In the figure, each slot contains two parts. The first part consists of hardware that compares the bits from the packet to the pattern stored in the slot. The second part stores a value that specifies an action to be taken if the pattern matches the packet. If a match occurs, the slot hardware passes the action to the component that checks all the results and announces an answer.

One of the most important details concerns the way TCAM handles multiple matches. In essence, the output circuitry selects one match and ignores the others. That is, if multiple slots each pass an action to the output circuit, the circuit accepts only one and passes the action as the output of the classification. For example, the hardware may choose the lowest slot that matches. In any case, the action that the TCAM announces corresponds to the action from one of the matching slots.

The figure indicates that a slot holds a *pattern* rather than an exact value. Instead of merely comparing each bit in the pattern to the corresponding bit in the packet, the hardware performs a pattern match. The adjective *ternary* is used because each bit position in a pattern can have three possible values: a one, a zero, or a “don’t care”. When a slot compares its pattern to the packet, the hardware checks only the one and zero bits in the pattern—the hardware ignores pattern bits that contain “don’t care”. Thus, a pattern can specify exact values for some fields in a packet header and omit other fields.

To understand TCAM pattern matching, consider a pattern that identifies IP packets. Identifying such packets is easy because an Ethernet frame that carries an IPv4 datagram will have the value **0x0800** in the Ethernet Type field. Furthermore, the Type field occupies a fixed position in the frame: bits 96 through 111. Thus, we can create a pattern that starts with 96 “don’t care” bits (to cover the Ethernet destination and source MAC addresses) followed by 16 bits with the binary value **0000100000000000** (the binary equivalent of **0x0800**) to cover the Type field. All remaining bit positions in the pattern will be “don’t care”. Figure 8 illustrates the pattern and example packets.

Figure 8: A TCAM Pattern and Example Packets



Although a TCAM hardware slot has one position for each bit, the figure does not display individual bits. Instead, each box corresponds to one octet, and the value in a box is a hexadecimal value that corresponds to 8 bits. We use hexadecimal simply because binary strings are too long to fit into a figure comfortably.

The Size of a TCAM

A question arises: how large is a TCAM? The question can be divided into two important aspects:

- *The number of bits in a slot:* The number of bits per slot depends on the type of Ethernet switch. A basic switch uses the destination MAC address to classify a packet. Because a MAC address is 48 bits, TCAM in a basic switch needs only 48 bit positions. A VLAN switch needs 128 bit positions to cover the VLAN tag as well as source and destination MAC addresses. A Layer 3 switch must have sufficient bit positions to cover the IP header as well as the Ethernet header. For IPv6, the header size is large and variable—in most cases, a pattern will need to cover extension headers as well as the base header.
- *The total number of slots:* The total number of TCAM slots determines the maximum number of patterns a classifier can hold. When a switch learns the MAC address of a computer that has been plugged into a port, the switch can store a pattern for the address. For example, if a computer with MAC address X is plugged into port 29, the switch can create a pattern in which destination address bits match X and the action is “*send packet to output port 29*”.

A switch can also use patterns to control broadcasting. When a manager configures a VLAN, the switch can add an entry for the VLAN broadcast. For example, if a manager configures VLAN 9, an entry can be added in which the destination address bits are all 1s (that is, the Ethernet broadcast address) and the VLAN tag is 9. The action associated with the entry is “*broadcast on VLAN 9*”.

A Layer 3 switch can learn the IP source address of computers attached to the switch, and can use TCAM to store an entry for each IP address. Similarly, it is possible to create entries that match Layer 4 protocol port numbers (for example, to direct all web traffic to a specific output). SDN technologies allow a manager to place patterns in the classifier to establish paths through a network and direct traffic along the paths. Because such classification rules cross multiple layers of the protocol stack, the potential number of items stored in a TCAM can be large.

TCAM seems like an ideal mechanism because it is both extremely fast and versatile. However, TCAM has two significant drawbacks: cost and heat. The cost is high because TCAM has parallel hardware for each slot and the overall system is designed to operate at high speed. In addition, because it operates in parallel, TCAM consumes much more energy than conventional memory (and generates more heat). Therefore, designers minimize the amount of TCAM to keep costs and power consumption low. A typical switch has 32,000 entries.

Classification-Enabled Generalized Forwarding

Perhaps the most significant advantage of a classification mechanism arises from the generalizations it enables. Because classification examines arbitrary fields in a packet before any demultiplexing occurs, cross-layer combinations are possible. For example, classification can specify that all packets from a given MAC address should be forwarded to a specific output port regardless of the packet contents. In addition, classification can make forwarding decisions depend on combinations of source and destination. An *Internet Service Provider* (ISP) can choose to forward all packets with IP source address X that are destined for web server W along one path while forwarding packets with IP source address Y that are destined to the same web server along another path.

ISPs need the generality that classification offers to handle traffic engineering that is not usually available in a conventional protocol stack. In particular, classification allows an ISP to offer tiered services in which the path a packet follows depends on a combination of the type of traffic and how much the customer pays.

Summary

Classification is a fundamental performance optimization that allows a packet-processing system to cross layers of the protocol stack without demultiplexing. A classifier treats each packet as an array of bits and checks the contents of fields at specific locations in the array.

Classification offers high-speed forwarding for network systems such as Ethernet switches and routers that send packets across MPLS tunnels. To achieve the highest speed, classification can be implemented in hardware; a hardware technology known as TCAM is especially useful because it employs parallelism to perform classification at extremely high speed.

The generalized forwarding capabilities that classification provides allow ISPs to perform traffic engineering. When making a forwarding decision, a classification mechanism can use the source of a packet as well as the destination (for example, to choose a path based on the tier of service to which a customer subscribes).

Acknowledgment

Material in this article has been taken with permission from Douglas E. Comer, *Internetworking With TCP/IP Volume 1: Principles, Protocols, and Architecture*, Sixth edition, 2013.

References

- [1] Douglas E. Comer and David L. Stevens, *Internetworking With TCP/IP Volume 2: Design, Implementation, and Internals*, Prentice-Hall, Upper Saddle River, NJ, Third edition, 1999.
- [2] yuba.stanford.edu/~nickm/papers/classification_tutorial_01.pdf
- [3] Patrick McHardy, “nfttables: A Successor to iptables, ip6tables, ebtables and arptables,” *Netfilter Workshop 2008*, Paris, 2008.
- [4] Douglas E. Comer, *Network Systems Design Using Network Processors*, Intel IXP 2xxx version, Prentice-Hall, Upper Saddle River, NJ, 2006.

DOUGLAS E. COMER is a Distinguished Professor of Computer Science at Purdue University. Formerly, he served as VP of Research and Research Collaboration at Cisco Systems. As a member of the original IAB, he participated in early work on the Internet, and is internationally recognized as an authority on TCP/IP protocols and Internet technologies. He has written a series of best-selling technical books, and his three-volume *Internetworking* series is cited as an authoritative work on Internet technologies. His books, which have been translated into 16 languages, are used in industry and academia in many countries. Comer consults for industry, and has lectured to thousands of professional engineers and students around the world. For 20 years he was editor-in-chief of the journal *Software—Practice and Experience*. He is a Fellow of the ACM and the recipient of numerous teaching awards. E-mail: comer@cs.purdue.edu

Internet Society Disappointed over Fundamental Divides at WCIT-12

On December 14, 2012, The Internet Society released the following statement from President and CEO Lynn St. Amour:

“The Internet Society, like other participants at the *World Conference on International Telecommunications* (WCIT), came to this conference looking for a successful outcome. We were hopeful that it would result in a treaty that would enable growth, further innovation, and advance interoperability in international telecommunications. It was extremely important that this treaty not extend to content, or implicitly or explicitly undermine the principles that have made the Internet so beneficial.

While progress was made in some areas such as transparency in international roaming fees, fundamental divides were exposed leaving a significant number of countries unable to sign the *International Telecommunication Regulations* (ITRs). Statements made by a host of delegations today made it very clear that Internet issues did not belong in the ITRs and that they would not support a treaty that is inconsistent with the multi-stakeholder model of Internet Governance.

We are disappointed that the conference has not been successful in reaching consensus. The Internet Society is dedicated to working with all stakeholders around the world to create the environment that will allow the Internet to grow for the betterment of all people.”

For more information, see:

<http://www.internetsociety.org/wcit>

See also:

[0] Geoff Huston, “December in Dubai,” *The Internet Protocol Journal*, Volume 15, No. 2, June 2012.

[1] World Conference on International Telecommunications (WCIT-12), <http://www.itu.int/en/wcit-12/Pages/default.aspx>

[2] “NRO contribution to the WCIT Public Consultation Process,” <http://www.nro.net/wp-content/uploads/2012/joint-submission-WCIT-RIR.pdf>

[3] “Stop the Net Grab”: NRO Shares Concerns About the WCIT Process,” <http://www.nro.net/news/nro-shares-concerns-aboutwcit-process>

[4] WCITLeaks.org “Bringing transparency to the ITU,” <http://wcitleaks.org>

NRO Observations on WCIT-12 Process

The *Number Resource Organization* (NRO), representing the world's five *Regional Internet address Registries* (RIRs), issued the following statement from Dubai, the site of the recent *World Conference on International Telecommunications* (WCIT):

The conference has clearly not met expectations of many *International Telecommunications Union* (ITU) Member States, and with this unfortunate outcome now clear, we feel compelled to put the following observations on record.

The NRO is concerned about aspects of the WCIT-12 meetings, which have just ended in Dubai, particularly with events in the last days of the conference. Neither the content of this conference, nor its conduct during this critical final period, have met community expectations or satisfied public assurances given prior to the event.

Internet stakeholders around the world watched the WCIT preparations closely, and were hopeful, throughout those processes, of two things: that WCIT would have no bearing on the Internet, its governance or its content; and that the event would allow all voices to be heard. The ITU Secretary General himself made these assurances on multiple occasions, and reiterated them in his opening remarks to the conference.

Regrettably, expected WCIT discussions on traditional telecommunication issues were eclipsed by debates about Internet-related issues. The intensity and length of these debates revealed clearly the depth of genuine concern about the proposals, and also the determination of those who brought them to the meeting.

Perhaps more importantly, an open multi-stakeholder conduct of the WCIT conference did not eventuate. Plenary sessions of the conference were webcast, but contributions were allowed only from official Government delegates and ITU officials, relegating all other stakeholders to an observer role.

Furthermore, an important number of critical negotiations occurred in small groups accessible only to Member States; and key experts and other stakeholders were unable even to observe them.

The NRO strongly supports the principles established in 2005 by the *World Summit on the Information Society*, which call for Internet Governance to be carried out in a multi-stakeholder manner, and we note that these represent the view of the global community as expressed through the United Nations system itself.

The NRO has also participated in many ITU conferences and study groups over the years, at very substantial cost, in genuine efforts to build relationships between our communities and to demonstrate the value of multi-stakeholder cooperation and collaboration. The NRO will continue to participate in the ITU, itself a member of the UN system, in expectation that its processes can evolve visibly, and much more rapidly, towards these accepted principles.

John Jason Brzozowski, Donn Lee, and Paul Saab win 2012 Itojun Awards

The fourth annual *Itojun Service Awards* were recently presented to John Jason Brzozowski for his tireless efforts in providing IPv6 connectivity to cable broadband users across North America and evangelizing the importance of IPv6 deployment globally, and to Donn Lee and Paul Saab for their efforts in making high-profile online content available over IPv6 and for their key contributions to *World IPv6 Day* and *World IPv6 Launch*. The awardees were recognized at the *Internet Engineering Task Force* (IETF) 85 meeting in November 2012 in Atlanta, Georgia.

First awarded in 2009, the award honors the memory of Dr. Jun-ichiro “Itojun” Hagino, who passed away in 2007 at the age of 37. The award, established by the friends of Itojun and administered by the Internet Society, recognizes and commemorates the extraordinary dedication exercised by Itojun over the course of IPv6 development. IPv6, the next-generation Internet protocol developed within the IETF, provides more than 340 trillion, trillion, trillion addresses, enabling billions of people and a huge range of devices to connect with one another, and helping ensure the Internet continues its current growth rate indefinitely.

“The combined work of John, Donn, and Paul has made IPv6 a technology used every day by people around the world as they access some of the most popular websites from their homes and offices,” said Jun Murai of the Itojun Service Award committee and founder of the WIDE Project.

“On behalf of the Itojun Service Award committee, I am extremely pleased to present this award to them for their ongoing efforts that have made IPv6 a mainstream technology for global web companies looking to ensure their continued growth.”

The Itojun Service Award is focused on pragmatic contributions to developing and deploying IPv6 in the spirit of serving the Internet. With respect to the spirit, the selection committee seeks contributors to the Internet as a whole; open source developers are a common example of such contributors, although this is not a requirement for expected nominees.

While the committee primarily considers practical contributions such as software development or network operation, higher level efforts that help those direct contributions will also be appreciated in this regard. The contribution should be substantial, but could be at an immature stage or be ongoing; this award aims to encourage the contributor to continue their efforts, rather than just recognizing well established work. Finally, contributions of a group of individuals will be accepted, as deployment work is often done by a large project, not just a single outstanding individual.

The award includes a presentation crystal, a US\$3,000 honorarium, and a travel grant.

John Jason Brzozowski said, “It is truly humbling to be a recipient of the Itojun Service Award, being recognized with others that have worked tirelessly to make IPv6 a reality is rewarding personally and professionally. I would like to thank the award committee and the Internet Society as well as my family and co-workers for their support. As many are aware, the IPv6 journey at Comcast has been unfolding since 2005. It is an honor and pleasure to provide the technical and strategic leadership for IPv6 that has led to the success of our program and the widespread adoption of IPv6.”

Donn Lee said, “Deploying IPv6 continues to be an amazing experience. I’m thankful to be sharing this award with my colleagues Paul and John, whom I have worked alongside through the challenging and exciting milestones of World IPv6 Day 2011 and World IPv6 Launch 2012. I especially want to thank the award committee for this honor that remembers Itojun, a truly inspirational IPv6 scientist, leader, and visionary.”

Paul Saab said, “I’m honored to be sharing the Itojun Service Award with Donn and John. We should never forget that we would not be here today if it were not for Itojun’s trailblazing work and passion for IPv6. To be recognized is extremely humbling, as Facebook’s participation could not have been done without our amazing co-workers and their own hard work to bring IPv6 to our users. Thank you for recognizing us and remember that this journey is only 2% complete.”

For more information about the Itojun Service Award see:
<http://www.internetsociety.org/what-we-do/grants-and-awards/awards/itojun-service-award>



Left to right: Jun Murai, John Jason Brzozowski, Paul Saab and Don Lee

Leading Global Standards Organizations Endorse “OpenStand” Principles

Five leading global organizations—the *Institute for Electrical and Electronics Engineers* (IEEE), the *Internet Architecture Board* (IAB), the *Internet Engineering Task Force* (IETF), the *Internet Society* and the *World Wide Web Consortium* (W3C)—recently announced that they have signed a statement affirming the importance of a jointly developed set of principles establishing a modern paradigm for global, open standards. The shared “OpenStand” principles—based on the effective and efficient standardization processes that have made the Internet and Web the premiere platforms for innovation and borderless commerce—are proven in their ability to foster competition and cooperation, support innovation and interoperability and drive market success.

The IEEE, IAB, IETF, Internet Society and W3C invite other standards organizations, governments, corporations and technology innovators globally to endorse the principles, available at open-stand.org

The OpenStand principles strive to encapsulate that successful standardization model and make it extendable across the contemporary, global economy’s gamut of technology spaces and markets. The principles comprise a modern paradigm in which the economics of global markets—fueled by technological innovation—drive global deployment of standards, regardless of their formal status within traditional bodies of national representation. The OpenStand principles demand:

- Cooperation among standards organizations;
- Adherence to due process, broad consensus, transparency, balance and openness in standards development;
- Commitment to technical merit, interoperability, competition, innovation and benefit to humanity;
- Availability of standards to all; and
- Voluntary adoption.

“New dynamics and pressures on global industry have driven changes in the ways that standards are developed and adopted around the world,” said Steve Mills, president of the IEEE Standards Association.

“Increasing globalization of markets, the rapid advancement of technology and intensifying time-to-market demands have forced industry to seek more efficient ways to define the global standards that help expand global markets. The OpenStand principles foster the more efficient international standardization paradigm that the world needs.”

Added Leslie Daigle, chief Internet technology officer with the Internet Society: “International standards development for borderless economics is not ad hoc; rather, it has a paradigm—one that has demonstrated agility and is driven by technical merit.

The OpenStand principles convey the power of bottom-up collaboration in harnessing global creativity and expertise to the standards of any technology space that will underpin the modern economy moving forward.”

Standards developed and adopted via the OpenStand principles include IEEE standards for the Internet’s physical connectivity, IETF standards for end-to-end global Internet interoperability and the W3C standards for the World Wide Web.

“The Internet and World Wide Web have fueled an economic and social transformation, touching billions of lives. Efficient standardization of so many technologies has been key to the success of the global Internet,” said Russ Housley, IETF chair. “These global standards were developed with a focus toward technical excellence and deployed through collaboration of many participants from all around the world. The results have literally changed the world, surpassing anything that has ever been achieved through any other standards-development model.”

Globally adopted design-automation standards, which have paved the way for a giant leap forward in industry’s ability to define complex electronic solutions, provide another example of standards developed in the spirit of the OpenStand principles. Another technology space that figures to demand such standards over the next decades is the global smart-grid effort, which seeks to augment regional facilities for electricity generation, distribution, delivery and consumption with a two-way, end-to-end network for communications and control.

“Think about all that the Internet and Web have enabled over the past 30 years, completely transforming society, government and commerce,” said W3C chief executive officer Jeff Jaffe. “It is remarkable that a small number of organizations following a small number of principles have had such a huge impact on humanity, innovation and competition in global markets.”

Bernard Aboba, chair of the IAB said: “The Internet has been built on specifications adopted voluntarily across the globe. By valuing running code, interoperability and deployment above formal status, the Internet has democratized the development of standards, enabling specifications originally developed outside of standards organizations to gain recognition based on their technical merit and adoption, contributing to the creation of global communities benefiting humanity. We now invite standards organizations, as well as governments, companies and individuals to join us at open-stand.org in order to affirm the principles that have nurtured the Internet and underpin many other important standards—and will continue to do so.”

New Year's Day 2013 Marks 30th Anniversary of Major Milestone for the Internet

On January 1, 1983, the ARPANET, a direct predecessor of today's Internet, implemented the *Transmission Control Protocol/Internet Protocol* (TCP/IP) in a transition that required all connected computers to convert to the protocol simultaneously. The open TCP/IP protocol is now a foundational technology for the networks around the world that make up the global Internet and interconnect billions of devices. The transition, which was carefully planned over several years before it actually took place, is documented in RFC 801^[1] authored by Jon Postel^[2].

Throughout its history, the Internet has continued to evolve. Today, deploying IPv6, the latest generation of the IP protocol, is critical to ensuring the Internet's continued growth and to connect the billions of people not yet online. Thousands of major *Internet Service Providers* (ISPs), home networking equipment manufacturers, and web companies around the world are coming together to permanently enable IPv6 for their products and services through efforts such as *World IPv6 Launch*^[3] organized by the Internet Society.

For more information about the Internet Society's work to facilitate the open development of standards, protocols, and administration, and to ensure a robust, secure technical infrastructure, see the *Internet Technology Matters* blog^[4] and the *Deploy360 Programme*^[5]. For further details about the Internet's history and development, see [6].

[1] Jon Postel, "NCP/TCP transition plan," RFC 801, November 1981.

[2] <http://www.internethalloffame.org/inductees/jon-postel>

[3] <http://www.worldipv6launch.org/>

[4] <http://www.internetsociety.org/what-we-do/internet-technology-matters>

[5] <http://www.internetsociety.org/deploy360/>

[6] Barry M. Leiner, Vinton G. Cerf, David D. Clark, Robert E. Kahn, Leonard Kleinrock, Daniel C. Lynch, Jon Postel, Larry G. Roberts, and Stephen Wolff, "Brief History of the Internet," <http://www.internetsociety.org/internet/what-internet/history-internet/brief-history-internet>

What is my “Subscription ID” for The Internet Protocol Journal (IPJ) and where do I find it?

IPJ Subscription FAQ

Your Subscription ID is a unique combination of letters and numbers used to locate your subscription in our database. It is printed on the back of your IPJ issue or on the envelope. You will also find information about your subscription expiration date near your Subscription ID. Here is an example:



How do I renew or update my subscription?

From the IPJ homepage (www.cisco.com/ipj) click “Subscriber Service” and then enter your Subscription ID and your e-mail address in the boxes. After you click “Login” the system will send you an e-mail message with a unique URL that allows access to your subscription record. You can then update your postal and e-mail details, change delivery options, and of course *renew* your subscription.

What will you use my e-mail address and postal address for?

This information is used *only* to communicate with you regarding your subscription. You will receive renewal reminders as well as other information about your subscription. We will never use your address for any form of marketing or unsolicited e-mail.

I didn’t receive the special URL that allows me to renew or update my Subscription. Why?

This is likely due to some form of spam filtering. Just send an e-mail message to ipj@cisco.com with your Subscription ID and any necessary changes and we will make the changes for you.

Do I need my Subscription ID to read IPJ online? What is my username and password?

Your Subscription ID is used *only* for access to your subscription record. No username or password is required to read IPJ. All back issues are available for online browsing or for download at www.cisco.com/ipj

I can’t find my Subscription ID and I have since changed e-mail address anyway; what do I do now?

Just send a message to ipj@cisco.com and we will take care of it for you.

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ contains standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2012 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

March 2013

Volume 16, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
SDN and OpenFlow	2
Address Authentication	15
WCIT Report	21
Letters to the Editor.....	34
Book Review.....	36
Fragments	38
Call for Papers.....	39

FROM THE EDITOR

This is the 60th edition of *The Internet Protocol Journal*, and in June we will celebrate our 15th anniversary. Fifteen years is not a long time in absolute terms, but when it comes to networking technology a lot can happen in a short time.

Throughout this 15-year period we have published numerous articles on “emerging technologies,” and in this issue we present yet another. *Software-Defined Networks* (SDNs) have become a mainstream topic for research, development, and standardization. We asked William Stallings to give us an overview of SDNs, and we plan further articles on this topic in the future.

A recurring theme in this journal has been Internet *security* at all levels of the protocol stack. We have covered security in routing, securing the *Domain Name System* (DNS), secure wireless networks, secure HTTP, and much more. This time, Scott Hogg discusses the advantages and disadvantages of using IPv4 or IPv6 addresses as a form of user authentication.

In our previous issue we published some reactions to the outcomes of the *World Conference on International Telecommunications* (WCIT) held in Dubai in December 2012. In this edition, Robert Pepper and Chip Sharp provide analysis and background on this conference and discuss how the revised *International Telecommunication Regulations* (ITRs) might affect the future of the Internet.

It has been some time since we have published a book review, but we are happy to bring you one in this issue. For the first time in history, we are reviewing a book that exists only in electronic form, another sign of a rapidly changing technology landscape. We are always looking for new book reviews. Please send your reviews, letters to the editor, or any subscription questions to ipj@cisco.com

If you want to look back at 15 years of IPJ, visit our website at www.cisco.com/ipj where you will find all of our back issues (as a single PDF file, as a collection of individual PDF files, or in HTML format), as well as an index of all IPJ articles.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

Software-Defined Networks and OpenFlow

by William Stallings

A network organizing technique that has come to recent prominence is the *Software-Defined Network* (SDN)^[1]. In essence, an SDN separates the data and control functions of networking devices, such as routers, packet switches, and LAN switches, with a well-defined *Application Programming Interface* (API) between the two. In contrast, in most large enterprise networks, routers and other network devices encompass both data and control functions, making it difficult to adjust the network infrastructure and operation to large-scale addition of end systems, virtual machines, and virtual networks. In this article we examine the characteristics of an SDN, and then describe the *OpenFlow* specification, which is becoming the standard way of implementing an SDN.

Evolving Network Requirements

Before looking in more detail at SDNs, let us examine the evolving network requirements that lead to a demand for a flexible, response approach to controlling traffic flows within a network or the Internet.

One key leading factor is the increasingly widespread use of *Server Virtualization*. In essence, server virtualization masks server resources, including the number and identity of individual physical servers, processors, and operating systems, from server users. This masking makes it possible to partition a single machine into multiple, independent servers, conserving hardware resources. It also makes it possible to migrate a server quickly from one machine to another for load balancing or for dynamic switchover in the case of machine failure. Server virtualization has become a central element in dealing with “big data” applications and in implementing cloud computing infrastructures. But it creates problems with traditional network architectures (for example, refer to [2]). One problem is configuring *Virtual LANs* (VLANs). Network managers need to make sure the VLAN used by the *Virtual Machine* is assigned to the same switch port as the physical server running the virtual machine. But with the virtual machine being movable, it is necessary to reconfigure the VLAN every time that a virtual server is moved. In general terms, to match the flexibility of server virtualization, the network manager needs to be able to dynamically add, drop, and change network resources and profiles. This process is difficult to do with conventional network switches, in which the control logic for each switch is co-located with the switching logic.

Another effect of server virtualization is that traffic flows differ substantially from the traditional client-server model. Typically, there is a considerable amount of traffic among virtual servers, for such purposes as maintaining consistent images of the database and invoking security functions such as access control. These server-to-server flows change in location and intensity over time, demanding a flexible approach to managing network resources.

Another factor leading to the need for rapid response in allocating network resources is the increasing use by employees of mobile devices such as smartphones, tablets, and notebooks to access enterprise resources. Network managers must be able to respond to rapidly changing resource, *Quality of Service* (QoS), and security requirements.

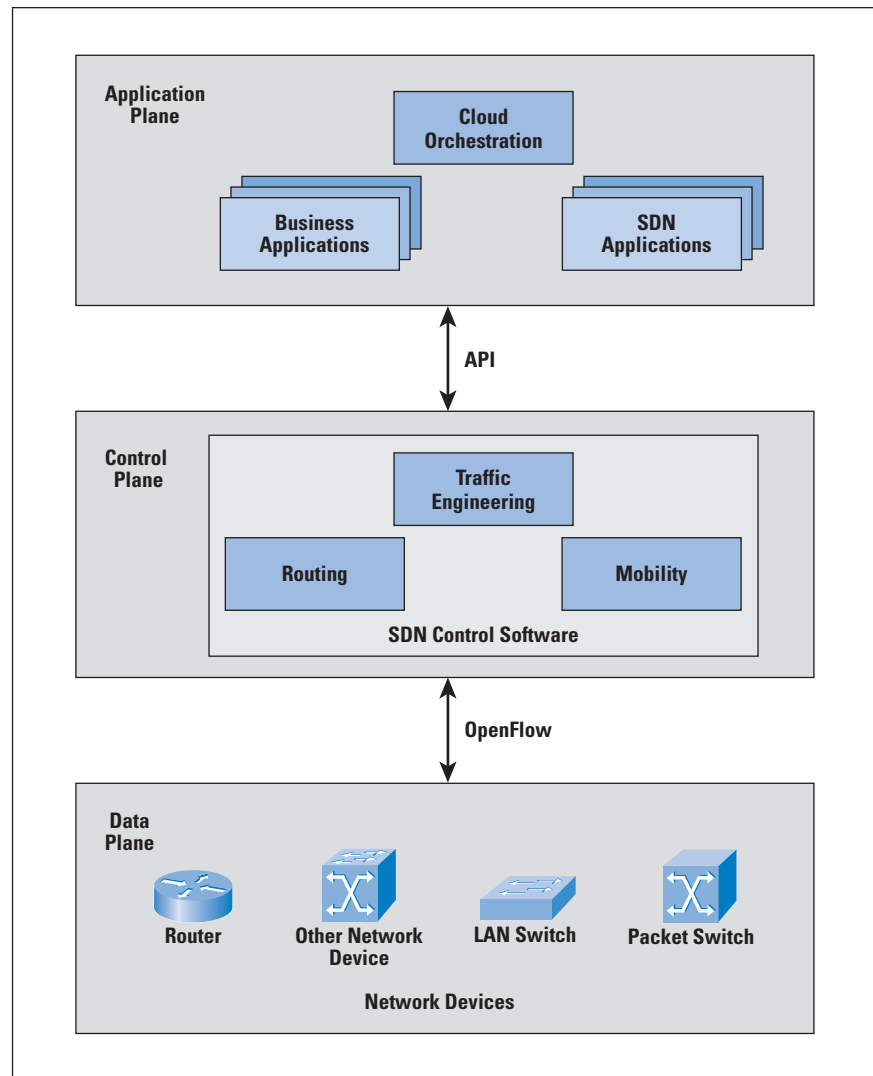
Existing network infrastructures can respond to changing requirements for the management of traffic flows, providing differentiated QoS levels and security levels for individual flows, but the process can be very time-consuming if the enterprise network is large and/or involves network devices from multiple vendors. The network manager must configure each vendor's equipment separately, and adjust performance and security parameters on a per-session, per-application basis. In a large enterprise, every time a new virtual machine is brought up, it can take hours or even days for network managers to do the necessary reconfiguration^[3].

This state of affairs has been compared to the mainframe era of computing^[4]. In the era of the mainframe, applications, the operating system, and the hardware were vertically integrated and provided by a single vendor. All of these ingredients were proprietary and closed, leading to slow innovation. Today, most computer platforms use the x86 instruction set, and a variety of operating systems (Windows, Linux, or Mac OS) run on top of the hardware. The OS provides APIs that enable outside providers to develop applications, leading to rapid innovation and deployment. In a similar fashion, commercial networking devices have proprietary features and specialized control planes and hardware, all vertically integrated on the switch. As will be seen, the SDN architecture and the OpenFlow standard provide an open architecture in which control functions are separated from the network device and placed in accessible control servers. This setup enables the underlying infrastructure to be abstracted for applications and network services, enabling the network to be treated as a logical entity.

SDN Architecture

Figure 1 illustrates the logical structure of an SDN. A central controller performs all complex functions, including routing, naming, policy declaration, and security checks. This plane constitutes the *SDN Control Plane*, and consists of one or more SDN servers.

Figure 1: SDN Logical Structure



The *SDN Controller* defines the data flows that occur in the *SDN Data Plane*. Each flow through the network must first get permission from the controller, which verifies that the communication is permissible by the network policy. If the controller allows a flow, it computes a route for the flow to take, and adds an entry for that flow in each of the switches along the path. With all complex functions subsumed by the controller, switches simply manage flow tables whose entries can be populated only by the controller. Communication between the controller and the switches uses a standardized protocol and API. Most commonly this interface is the OpenFlow specification, discussed subsequently.

The SDN architecture is remarkably flexible; it can operate with different types of switches and at different protocol layers. SDN controllers and switches can be implemented for Ethernet switches (Layer 2), Internet routers (Layer 3), transport (Layer 4) switching, or application layer switching and routing. SDN relies on the common functions found on networking devices, which essentially involve forwarding packets based on some form of flow definition.

In an SDN architecture, a switch performs the following functions:

- The switch encapsulates and forwards the first packet of a flow to an SDN controller, enabling the controller to decide whether the flow should be added to the switch flow table.
- The switch forwards incoming packets out the appropriate port based on the flow table. The flow table may include priority information dictated by the controller.
- The switch can drop packets on a particular flow, temporarily or permanently, as dictated by the controller. Packet dropping can be used for security purposes, curbing *Denial-of-Service* (DoS) attacks or traffic management requirements.

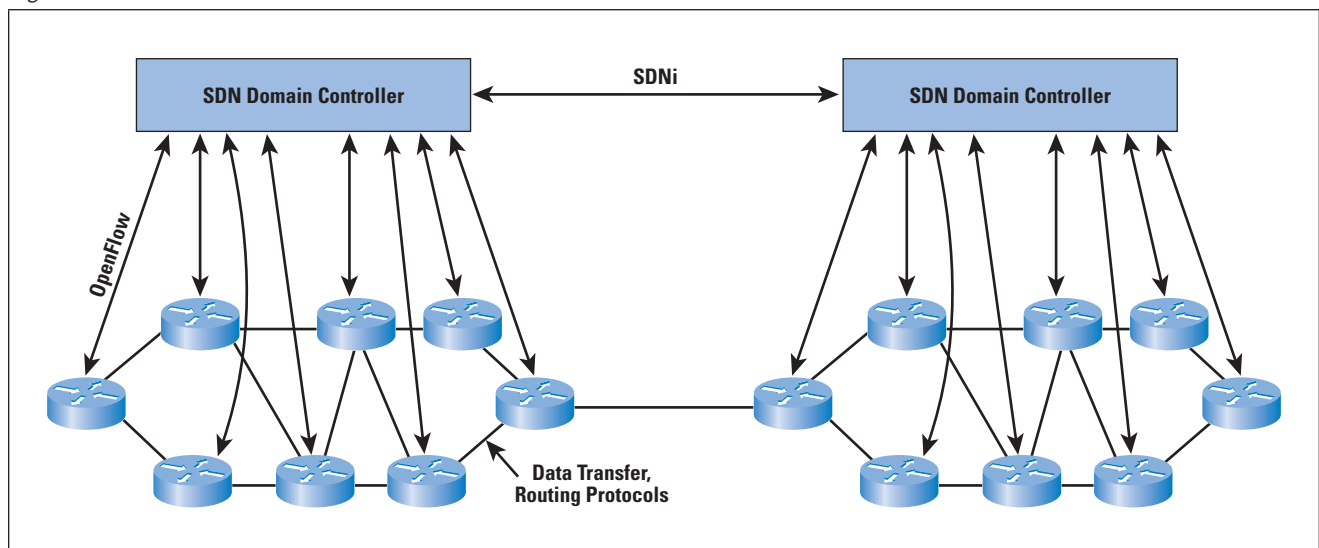
In simple terms, the SDN controller manages the forwarding state of the switches in the SDN. This management is done through a vendor-neutral API that allows the controller to address a wide variety of operator requirements without changing any of the lower-level aspects of the network, including topology.

With the decoupling of the control and data planes, SDN enables applications to deal with a single abstracted network device without concern for the details of how the device operates. Network applications see a single API to the controller. Thus it is possible to quickly create and deploy new applications to orchestrate network traffic flow to meet specific enterprise requirements for performance or security.

SDN Domains

In a large enterprise network, the deployment of a single controller to manage all network devices would prove unwieldy or undesirable. A more likely scenario is that the operator of a large enterprise or carrier network divides the whole network into numerous nonoverlapping SDN domains as shown in Figure 2.

Figure 2: SDN Domain Structure



Reasons for using SDN domains include the following:

- *Scalability*: The number of devices an SDN controller can feasibly manage is limited. Thus, a reasonably large network may need to deploy multiple SDN controllers.
- *Privacy*: A carrier may choose to implement different privacy policies in different SDN domains. For example, an SDN domain may be dedicated to a set of customers who implement their own highly customized privacy policies, requiring that some networking information in this domain (for example, network topology) not be disclosed to an external entity.
- *Incremental deployment*: A carrier's network may consist of portions of traditional and newer infrastructure. Dividing the network into multiple, individually manageable SDN domains allows for flexible incremental deployment.

The existence of multiple domains creates a requirement for individual controllers to communicate with each other via a standardized protocol to exchange routing information. The IETF is currently working on developing a protocol, called *SDNi*, for “interfacing SDN Domain Controllers”^[5]. SDNi functions include:

- Coordinate flow setup originated by applications containing information such as path requirement, QoS, and service-level agreements across multiple SDN domains.
- Exchange reachability information to facilitate inter-SDN routing. This information exchange will allow a single flow to traverse multiple SDNs and have each controller select the most appropriate path when multiple such paths are available.

The message types for SDNi tentatively include the following:

- Reachability update
- Flow setup/tear-down/update request (including application capability requirements such as QoS, data rate, latency etc.)
- Capability update (including network-related capabilities such as data rate and QoS, and system and software capabilities available inside the domain)

OpenFlow

To turn the concept of SDN into practical implementation, two requirements must be met. First, there must be a common logical architecture in all switches, routers, and other network devices to be managed by an SDN controller. This logical architecture may be implemented in different ways on different vendor equipment and in different types of network devices, so long as the SDN controller sees a uniform logical switch function. Second, a standard, secure protocol is needed between the SDN controller and the network device.

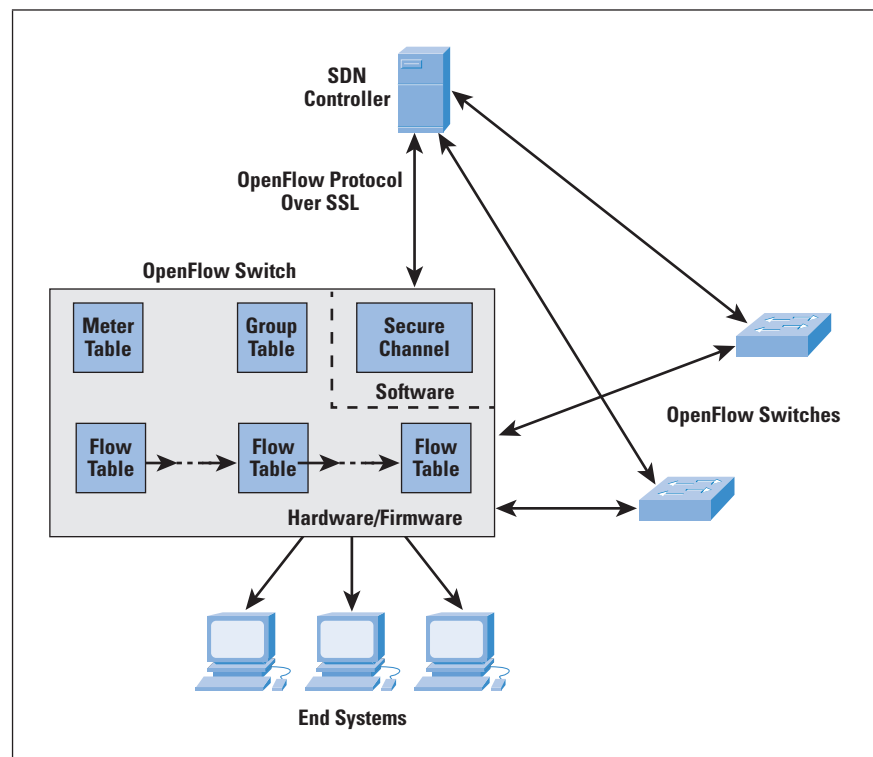
Both of these requirements are addressed by *OpenFlow*, which is both a protocol between SDN controllers and network devices, as well as a specification of the logical structure of the network switch functions^[6, 7]. OpenFlow is defined in the *OpenFlow Switch Specification*, published by the *Open Networking Foundation* (ONF). ONF is a consortium of software providers, content delivery networks, and networking equipment vendors whose purpose is to promote software-defined networking.

This discussion is based on the current OpenFlow specification, Version 1.3.0, June 25, 2012^[8]. The original specification, 1.0, was developed at Stanford University and was widely implemented. OpenFlow 1.2 was the first release from ONF after inheriting the project from Stanford. OpenFlow 1.3 significantly expands the functions of the specification. Version 1.3 is likely to become the stable base upon which future commercial implementations for OpenFlow will be built. ONF intends for this version to be a stable target for chip and software vendors, so little if any change is planned for the foreseeable future^[9].

Logical Switch Architecture

Figure 3 illustrates the basic structure of the OpenFlow environment. An SDN controller communicates with OpenFlow-compatible switches using the OpenFlow protocol running over the *Secure Sockets Layer* (SSL). Each switch connects to other OpenFlow switches and, possibly, to end-user devices that are the sources and destinations of packet flows. Within each switch, a series of tables—typically implemented in hardware or firmware—are used to manage the flows of packets through the switch.

Figure 3: OpenFlow Switch



The OpenFlow specification defines three types of tables in the logical switch architecture. A *Flow Table* matches incoming packets to a particular flow and specifies the functions that are to be performed on the packets. There may be multiple flow tables that operate in a pipeline fashion, as explained subsequently. A flow table may direct a flow to a *Group Table*, which may trigger a variety of actions that affect one or more flows. A *Meter Table* can trigger a variety of performance-related actions on a flow.

Before proceeding, it is helpful to define what the term *flow* means. Curiously, this term is not defined in the OpenFlow specification, nor is there an attempt to define it in virtually all of the literature on OpenFlow. In general terms, a flow is a sequence of packets traversing a network that share a set of header field values. For example, a flow could consist of all packets with the same source and destination IP addresses, or all packets with the same VLAN identifier. We provide a more specific definition subsequently.

Flow-Table Components

The basic building block of the logical switch architecture is the flow table. Each packet that enters a switch passes through one or more flow tables. Each flow table contains entries consisting of six components:

- *Match Fields*: Used to select packets that match the values in the fields.
- *Priority*: Relative priority of table entries.
- *Counters*: Updated for matching packets. The OpenFlow specification defines a variety of timers. Examples include the number of received bytes and packets per port, per flow table, and per flow-table entry; number of dropped packets; and duration of a flow.
- *Instructions*: Actions to be taken if a match occurs.
- *Timeouts*: Maximum amount of idle time before a flow is expired by the switch.
- *Cookie*: Opaque data value chosen by the controller. May be used by the controller to filter flow statistics, flow modification, and flow deletion; not used when processing packets.

A flow table may include a *table-miss* flow entry, which renders all Match Fields wildcards (every field is a match regardless of value) and has the lowest priority (priority 0). The Match Fields component of a table entry consists of the following required fields:

- *Ingress Port*: The identifier of the port on the switch where the packet arrived. It may be a physical port or a switch-defined virtual port.
- *Ethernet Source and Destination Addresses*: Each entry can be an exact address, a bitmasked value for which only some of the address bits are checked, or a wildcard value (match any value).

- *IPv4 or IPv6 Protocol Number*: A protocol number value, indicating the next header in the packet.
- *IPv4 or IPv6 Source Address and Destination Address*: Each entry can be an exact address, a bitmasked value, a subnet mask value, or a wildcard value.
- *TCP Source and Destination Ports*: Exact match or wildcard value.
- *User Datagram Protocol (UDP) Source and Destination Ports*: Exact match or wildcard value.

The preceding match fields must be supported by any OpenFlow-compliant switch. The following fields may be optionally supported:

- *Physical Port*: Used to designate underlying physical port when packet is received on a logical port.
- *Metadata*: Additional information that can be passed from one table to another during the processing of a packet. Its use is discussed subsequently.
- *Ethernet Type*: Ethernet Type field.
- *VLAN ID and VLAN User Priority*: Fields in the IEEE 802.1Q Virtual LAN header.
- *IPv4 or IPv6 DS and ECN*: Differentiated Services and Explicit Congestion Notification fields.
- *Stream Control Transmission Protocol (SCTP) Source and Destination Ports*: Exact match or wildcard value.
- *Internet Control Message Protocol (ICMP) Type and Code Fields*: Exact match or wildcard value.
- *Address Resolution Protocol (ARP) Opcode*: Exact match in Ethernet Type field.
- *Source and Target IPv4 Addresses in Address Resolution Protocol (ARP) Payload*: Can be an exact address, a bitmasked value, a subnet mask value, or a wildcard value.
- *IPv6 Flow Label*: Exact match or wildcard.
- *ICMPv6 Type and Code fields*: Exact match or wildcard value.
- *IPv6 Neighbor Discovery Target Address*: In an IPv6 Neighbor Discovery message.
- *IPv6 Neighbor Discovery Source and Target Addresses*: Link-layer address options in an IPv6 Neighbor Discovery message.
- *Multiprotocol Label Switching (MPLS) Label Value, Traffic Class, and Bottom of Stack (BoS)*: Fields in the top label of an MPLS label stack.

Thus, OpenFlow can be used with network traffic involving a variety of protocols and network services. Note that at the MAC/link layer, only Ethernet is supported. Thus, OpenFlow as currently defined cannot control Layer 2 traffic over wireless networks.

We can now offer a more precise definition of the term *flow*. From the point of view of an individual switch, a flow is a sequence of packets that matches a specific entry in a flow table. The definition is packet-oriented, in the sense that it is a function of the values of header fields of the packets that constitute the flow, and not a function of the path they follow through the network. A combination of flow entries on multiple switches defines a flow that is bound to a specific path.

The *instructions component* of a table entry consists of a set of instructions that are executed if the packet matches the entry. Before describing the types of instructions, we need to define the terms “Action” and “Action Set.” Actions describe packet forwarding, packet modification, and group table processing operations. The OpenFlow specification includes the following actions:

- *Output*: Forward packet to specified port.
- *Set-Queue*: Sets the queue ID for a packet. When the packet is forwarded to a port using the output action, the queue id determines which queue attached to this port is used for scheduling and forwarding the packet. Forwarding behavior is dictated by the configuration of the queue and is used to provide basic QoS support.
- *Group*: Process packet through specified group.
- *Push-Tag/Pop-Tag*: Push or pop a tag field for a VLAN or MPLS packet.
- *Set-Field*: The various Set-Field actions are identified by their field type; they modify the values of respective header fields in the packet.
- *Change-TTL*: The various Change-TTL actions modify the values of the IPv4 Time To Live (TTL), IPv6 Hop Limit, or MPLS TTL in the packet.

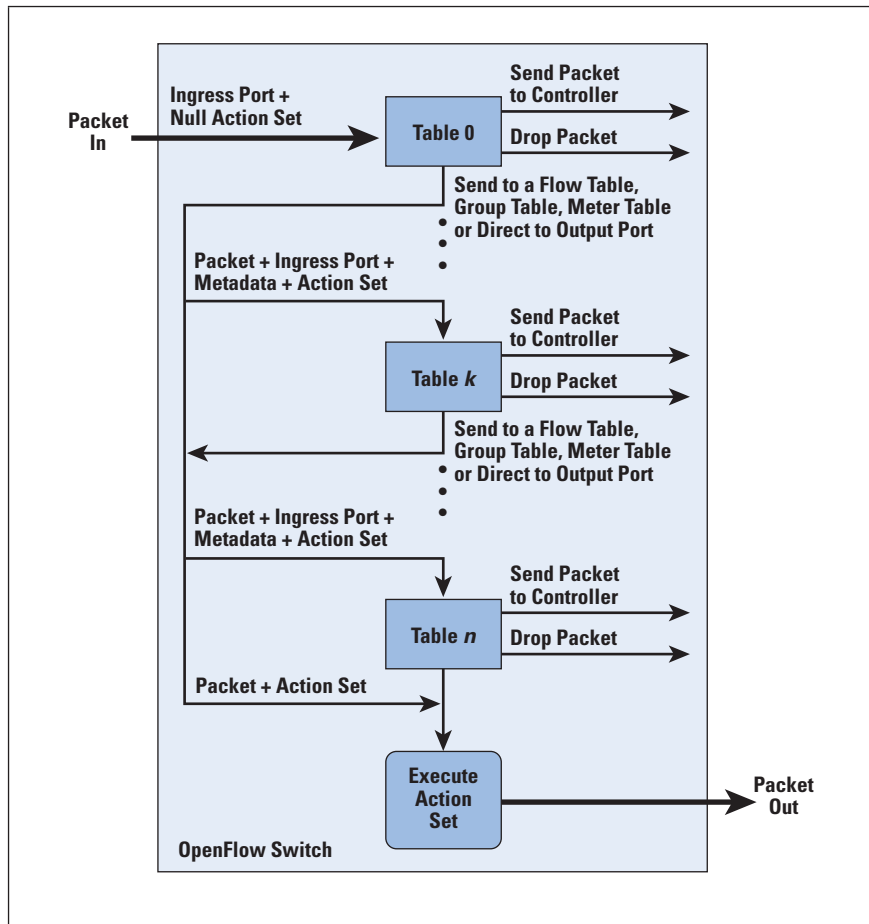
An *Action Set* is a list of actions associated with a packet that are accumulated while the packet is processed by each table and executed when the packet exits the processing pipeline. Instructions are of four types:

- *Direct packet through pipeline*: The Goto-Table instruction directs the packet to a table farther along in the pipeline. The Meter instruction directs the packet to a specified meter.
- *Perform action on packet*: Actions may be performed on the packet when it is matched to a table entry.
- *Update action set*: Merge specified actions into the current action set for this packet on this flow, or clear all the actions in the action set.
- *Update metadata*: A metadata value can be associated with a packet. It is used to carry information from one table to the next.

Flow-Table Pipeline

A switch includes one or more flow tables. If there is more than one flow table, they are organized as a pipeline as shown in Figure 4, with the tables labeled with increasing numbers starting with 0.

Figure 4: Packet Flow Through OpenFlow-Compliant Switch



When a packet is presented to a table for matching, the input consists of the packet, the identity of the ingress port, the associated metadata value, and the associated action set. For Table 0, the metadata value is blank and the action set is null. Processing proceeds as follows:

1. Find the highest-priority matching flow entry. If there is no match on any entry and there is no table-miss entry, then the packet is dropped. If there is a match only on a table-miss entry, then that entry specifies one of three actions:
 - a. Send packet to controller. This action will enable the controller to define a new flow for this and similar packets, or decide to drop the packet.
 - b. Direct packet to another flow table farther down the pipeline.
 - c. Drop the packet.

2. If there is a match on one or more entries other than the table-miss entry, then the match is defined to be with the highest-priority matching entry. The following actions may then be performed:
 - a. Update any counters associated with this entry.
 - b. Execute any instructions associated with this entry. These instructions may include updating the action set, updating the metadata value, and performing actions.
 - c. The packet is then forwarded to a flow table further down the pipeline, to the group table, or to the meter table, or it could be directed to an output port.

For the final table in the pipeline, forwarding to another flow table is not an option.

If and when a packet is finally directed to an output port, the accumulated action set is executed and then the packet is queued for output.

OpenFlow Protocol

The OpenFlow protocol describes message exchanges that take place between an OpenFlow controller and an OpenFlow switch. Typically, the protocol is implemented on top of SSL or *Transport Layer Security* (TLS), providing a secure OpenFlow channel.

The OpenFlow protocol enables the controller to perform add, update, and delete actions to the flow entries in the flow tables. It supports three types of messages, as shown in Table 1.

- *Controller-to-Switch*: These messages are initiated by the controller and, in some cases, require a response from the switch. This class of messages enables the controller to manage the logical state of the switch, including its configuration and details of flow- and group-table entries. Also included in this class is the Packet-out message. This message is used when a switch sends a packet to the controller and the controller decides not to drop the packet but to direct it to a switch output port.
- *Asynchronous*: These types of messages are sent without solicitation from the controller. This class includes various status messages to the controller. Also included is the Packet-in message, which may be used by the switch to send a packet to the controller when there is no flow-table match.
- *Symmetric*: These messages are sent without solicitation from either the controller or the switch. They are simple yet helpful. Hello messages are typically sent back and forth between the controller and switch when the connection is first established. Echo request and reply messages can be used by either the switch or controller to measure the latency or bandwidth of a controller-switch connection or just verify that the device is operating. The Experimenter message is used to stage features to be built into future versions of OpenFlow.

Table 1: OpenFlow Messages

Message		Description
Controller-to-Switch		
Features		Request the capabilities of a switch. Switch responds with a features reply that specifies its capabilities.
Configuration		Set and query configuration parameters. Switch responds with parameter settings.
Modify-State		Add, delete, and modify flow/group entries and set switch port properties.
Read-State		Collect information from switch, such as current configuration, statistics, and capabilities.
Packet-out		Direct packet to a specified port on the switch.
Barrier		Barrier request/reply messages are used by the controller to ensure message dependencies have been met or to receive notifications for completed operations.
Role-Request		Set or query role of the OpenFlow channel. Useful when switch connects to multiple controllers.
Asynchronous-Configuration		Set filter on asynchronous messages or query that filter. Useful when switch connects to multiple controllers.
Asynchronous		
Packet-in		Transfer packet to controller.
Flow-Removed		Inform the controller about the removal of a flow entry from a flow table.
Port-Status		Inform the controller of a change on a port.
Error		Notify controller of error or problem condition.
Symmetric		
Hello		Exchanged between the switch and controller upon connection startup.
Echo		Echo request/reply messages can be sent from either the switch or the controller, and they must return an echo reply.
Experimenter		For additional functions.

The OpenFlow protocol enables the controller to manage the logical structure of a switch, without regard to the details of how the switch implements the OpenFlow logical architecture.

Summary

SDNs, implemented using OpenFlow, provide a powerful, vendor-independent approach to managing complex networks with dynamic demands. The software-defined network can continue to use many of the useful network technologies already in place, such as virtual LANs and an MPLS infrastructure. SDNs and OpenFlow are likely to become commonplace in large carrier networks, cloud infrastructures, and other networks that support the use of big data.

References

- [1] Greg Goth, “Software-Defined Networking Could Shake Up More than Packets,” *IEEE Internet Computing*, July/August, 2011.
- [2] Robin Layland, “The Dark Side of Server Virtualization,” *Network World*, July 7, 2010.
- [3] Open Networking Foundation, “Software-Defined Networking: The New Norm for Networks,” ONF White Paper, April 12, 2012.
- [4] Dell, Inc., “Software Defined Networking: A Dell Point of View,” Dell White Paper, October 2012.
- [5] Hongtao Yin, Haiyong Xie, Tina Tsou, Diego Lopez, Pedro Aranda, and Ron Sidi, “SDNi: A Message Exchange Protocol for Software Defined Networks (SDNS) across Multiple Domains,” Internet Draft, work in progress, June 2012, **draft-yin-sdn-sdni-00.txt**
- [6] Steven Vaughan-Nichols, “OpenFlow: The Next Generation of the Network?” *Computer*, August 2011.
- [7] Thomas A. Limoncelli, “OpenFlow: A Radical New Idea in Networking,” *Communications of the ACM*, August 2012.
- [8] Open Networking Foundation, “OpenFlow Switch Specification Version 1.3.0,” June 25, 2012.
- [9] Sean Michael Kerner, “OpenFlow Protocol 1.3.0 Approved,” *Enterprise Networking Planet*, May 17, 2012.

WILLIAM STALLINGS is an independent consultant and author of many books on security, computer networking, and computer architecture. His latest book is *Data and Computer Communications* (Pearson, 2013). He maintains a computer science resource site for computer science students and professionals at **ComputerScienceStudent.com**. He has a Ph.D. in computer science from M.I.T. He can be reached at **ws@shore.net**

IPv4 and IPv6 Address Authentication

by Scott Hogg, GTRI

Some Internet services use the source address of the client's computer as a form of authentication. These systems keep track of the *Internet Protocol* (IP) address that an end user used the last time that user accessed the site and try to determine if the user is legitimate. When that same user accesses the site from a different source IP address, the site asks for further authentication to revalidate the client's computer. The theory is that a user's typical location computer has a somewhat persistent IP address, but when the user has a new address, that user may be mobile or using a less secure wireless media, and then require further authentication. For example, many organizations have firewall policies with objects named like "Bob's Laptop" with the single IP address of his computer. This technique is used by some banking sites, some online gaming sites, and Gmail (for example, *Google Authenticator*)^[1].

Some online retailers track the client IP address for Business Intelligence or fraud detection and forensics purposes. The retailer tracks the client IP address using the source address to analyze fraudulent purchases and to track down criminal activity. Some industries frequently use the customer's IP address as a form of authentication. Also, many sites that use *Server Load Balancers* (SLBs) and *Application Delivery Controllers* (ADCs) use *X-forwarded-for* (XFF)^[2] or *Hypertext Transfer Protocol* (HTTP) header insertion so that the back-end real servers are aware of the client's original IP address associated with the reverse proxy connection. The application can then use the IP address for tracking purposes or simply log the address with the transaction details.

Other applications try to validate the client's source IP address when the server receives an inbound connection. E-mail *Simple Mail Transfer Protocol* (SMTP) servers or *Internet Relay Chat* (IRC) servers can use the *Ident* protocol^[3] to try to validate the originating e-mail server or client computer validity. SMTP e-mail servers^[4] also use other protocols such as *SenderID*^[5], *Sender Policy Framework* (SPF)^[6], and *DomainKeys Identified Mail* (DKIM)^[7] in an effort to restrict spam. *Domain Name System* (DNS) *pointer* (PTR) records are sometimes used as a way to confirm that the client IP address is configured in DNS (for example, forward-confirmed reverse DNS^[8]).

Statically configured IP addresses are frequently used to signify some limited form of authentication. These addresses may not be used to authenticate a user, but authenticate IT systems to each other. Many manually configured systems rely on IP address to permit connectivity, including manually configured tunnels, *IP Security* (IPsec) peers, Apache *.htaccess*^[9], *.rhosts*^[10], SAMBA, and *Border Gateway Protocol* (BGP) peers, among many others.

The address is used as one part of the connection authentication. Obviously, IPsec connections are authenticated with certificates or preshared keys to strengthen their validation of the endpoints. Similarly, BGP peers use passwords (and/or *Time To Live* [TTL]^[11]) to help secure the peer beyond just IP address confirmation.

Identity-based firewalls police users' network behavior by IP address through *Windows Active Directory*, *Remote Authentication Dial-In User Service* (RADIUS), or *Lightweight Directory Access Protocol* (LDAP). Palo Alto firewalls championed the *UserID* concept as part of their analysis of connections to permit or deny authentication^[12]. The Cisco *Adaptive Security Appliance* (ASA) firewalls running Version 8.4 or later can be configured for Identify firewall functions^[13]. Firewalls have always used manually configured IP addresses as the fundamental element of their policies. The IP address is used in the policy as if that concretely defines a system and/or user. This process of adding rules based on IP address continues until the firewall is a pincushion full of pinholes.

Organizations that rely on using an IP address as a form of authentication run the risk of an attacker learning that IP address and attacking using that address. Attackers who know the addresses that are being used could perform a *Man-in-the-Middle* (MITM) attack or use TCP session hijacking. The attacker needs to know only the information about which IP addresses are used for the communications. The attacker might be able to ascertain the IP addresses the organization uses by guessing or by other means. The attacker could find the external IP address of the company's firewall and assume that IPv4 *Network Address Translation* (NAT)^[24] was being performed. The attacker could also suppose the business partner IP address. Organizations that use these techniques are relying on the secrecy of their IP addressing for the purposes of security.

Address Quality

The quality of the IP address is an important concept to consider. For example, a global address is of higher surety and authenticity than a private address. Many organizations use private addresses and overlap between private networks, whereas global addresses are unique and they are registered to a specific entity. Public addresses can reveal the client's *Internet Service Provider* (ISP), the organization that has registered the IP addresses, and some geolocation information. However, any IP packet can be spoofed and the source-address modified or crafted. Of course, if the source IP addresses is spoofed, the return packets will not necessarily be sent back to the attacker's source in these cases, but one-way blind attacks are still possible. Furthermore, systems such as *Tor*^[14] are intended to protect the identity of the end user.

Using the IP address as a form of authentication does not work if the client changes its location frequently. Today, many clients use mobile devices that can change their Layer 3 addresses often. The source IP address of the mobile device could change frequently and could even change during the transaction.

With increasing mobile device usage for business purposes, the ability to determine the typical IP address of the client becomes impossible. Increased scarcity of IPv4 addresses is leading service providers to use *Carrier-Grade NAT* (CGN) or *Large-Scale NAT* (LSN) and shorter and shorter *Dynamic Host Configuration Protocol* (DHCP) lease times, meaning that the client IP address is not static.

Many organizations and systems assume that a single computer with a single IP address represents a single user. The problem arises where IPv4 public addresses may not uniquely identify a single user. The industry may be trying to anticipate the implications of CGN/LSN and the effect of systems that rely on the uniqueness of a public IP address. Similar problems related to the mega-proxies of the late 1990s occurred (for example, AOL). With CGN/LSN systems in place, online retailers and banks will no longer be able to use the client IPv4 as the “real client IP.” Instead, the IP address observed on the retailer’s web servers will come from a pool of IPv4 addresses configured in the LSN system. In this situation, one bad actor could spoil that NAT pool IPv4 address for subsequent lawful users who follow. When a legitimate user tries to make an online purchase and that user’s system happens to use that IPv4 address of the bad actor, then the purchase attempt might be blocked. This situation would be bad for business on Cyber-Monday, or any day for that matter.

Table 1 compares IPv4 and IPv6 for their authentication purposes.

Table 1: IPv4 vs. IPv6 for Authentication

IPv4	IPv6
Extensive use of NAT	No motivation for NAT
End users use private addresses	End users use global addresses
Use of CGN/LSN starting	Abundance of IPv6 addresses
Robust geolocation	Geolocation needs improvement
Addresses could be spoofed	Addresses could be spoofed

Public Addresses

Public IPv4 addresses are becoming increasingly scarce^[15, 25], however, an abundance of global IPv6 addresses are available^[16]. Global IPv6 addresses can be obtained from *Regional Internet Registries* (RIRs) or from an IPv6-capable service provider. Residential broadband Internet users today use private IPv4 addresses on their internal computers, but these computers will soon start to use global IPv6 addresses as they upgrade to IPv6-capable *Customer Premises Equipment* (CPE). IPv6-enabled residential subscribers and employees of IPv6-enabled enterprises will be using global addresses when they access an IPv6-capable Internet service.

To online retailers, this situation may represent a change to their IP address authentication measures. As IPv4 residential users start to go through CGN/LSN systems, their IPv4 addresses will be useless for authentication.

However, their IPv6 addresses will be global addresses with no NAT taking place between the client and the server^[17]. It will be seemingly more accurate to use the IPv6 address to determine the validity of the source. IPv6 could potentially help to create an environment with more “trustworthiness” and less anonymity. For example, IPv6 IPsec connections could use the *Authentication Header* (AH) and *Encapsulating Security Payload* (ESP) together to create stronger connections, where IPv4 IPsec connections rely on NAT-Traversal and can use only ESP^[18].

As we head toward an increasingly dual-stack world, applications will need to do “dual-checking” of both the client’s IPv4 and IPv6 addresses. In a dual-stack world, there is more work to do^[19], and servers using IP address authentication will need to understand that a single user will have both an IPv4 address and an IPv6 address and keep track of both. The other consideration is that IPv6 nodes may have multiple global IPv6 addresses in some situations.

Authentication with Addresses

Security experts know that the secrecy of the encryption algorithm is not important, but the secrecy of the key is vitally important (Kerckhoffs’s Principle^[20]). The same concept should hold true for an IP address. Users should not rely on the secrecy of their IP addresses to be secure; the security of the individual node should be strong enough to defend against attacks. To the extreme, users should feel confident enough in their security posture that they feel comfortable widely publicizing their IP address. However, even if you are using *LifeLock*^[21], you should still keep your Social Security Number or government ID number private.

Security practitioners know that authentication should involve multiple factors. A combination of “something you are” (biometrics), “something you know” (username/password) and “something you have” (token, *Common Access Card*^[22]) forms a more solid foundation for identifying a user. Combining two factors provides more assurance than just one factor. We are all aware of the weaknesses of using username and password as a means of authentication^[23].

The systems mentioned so far in this article are three-factor systems (username, password, and IP address) which are presumably better than just username/password. However, we should acknowledge that an IP address is not a characteristic of a person. IP addresses have more to do with “somewhere you are,” because the IP address reflects location within a network topology by the prefix/subnet. The last few bits of an IPv4 address representing the point-of-attachment or an IPv6 *Interface Identifier* (IID) do not necessarily uniquely identify a user. Having authentication based on your location becomes difficult with mobile devices that roam widely. However, controlling authentication to users who are within the office subnet rather than outside the office may be useful.

An IP address is not something anyone really owns outright. Few organizations actually have complete ownership of their IP addresses. Organizations should read the fine print in the policies of their RIR. Organizations just pay RIR annual fees for their addresses, but if they stop paying those dues, the IP address allocation is revoked and the addresses go back into a pool for reallocation to another organization. Therefore, public IP addresses do not truly represent unequivocal ownership or legitimacy of a network.

Conclusion

Many different types of systems use the client's source address as a form of authentication. Systems that rely on IP address checking will need to do so for IPv4 and will need to be modified to use IPv6 addresses. IPv6 systems will use global addresses without NAT, so the security systems must stand on their own even though the IPv6 address is publicized. IPv4 and IPv6 addresses can be spoofed, and as CGN/LSN systems become widely deployed the validity of a public IPv4 address decreases. However, IPv6 addresses are not necessarily any more trustworthy than IPv4 addresses when used for authentication. Regardless, the IP address should not be the only factor used for authentication, and we should not be using IPv4 or IPv6 addresses as a form of authentication. The truth is that the IT industry needs to be aware of where IP addresses are used as a form of authentication and seek out better forms of authentication beyond just username, password, and IP address.

References

- [1] Google Authenticator,
<http://support.google.com/a/bin/answer.py?hl=en&answer=1037451>
- [2] <http://en.wikipedia.org/wiki/X-Forwarded-For>
- [3] Mike St. Johns, "Identification Protocol," RFC 1413, February 1993.
- [4] http://en.wikipedia.org/wiki/Email_authentication
- [5] Meng Weng Wong and Jim Lyon, "Sender ID: Authenticating E-Mail," RFC 4406, April 2006.
- [6] Wayne Schlitt and Meng Weng Wong, "Sender Policy Framework (SPF) for Authorizing Use of Domains in E-Mail, Version 1," RFC 4408, April 2006.
- [7] Miles Libbey, Michael Thomas, and Mark Delany, "DomainKeys Identified Mail (DKIM) Signatures," RFC 4871, May 2007.
- [8] http://en.wikipedia.org/wiki/Forward-confirmed_reverse_DNS
- [9] <http://en.wikipedia.org/wiki/Htaccess>
- [10] <http://en.wikipedia.org/wiki/Rlogin>

- [11] Vijay Gill, John Heasley, and David Meyer, “The Generalized TTL Security Mechanism (GTSM),” RFC 3682, February 2004.
- [12] Palo Alto Networks, UserID,
<http://www.paloaltonetworks.com/products/technologies/user-id.html>
- [13] Cisco ASA firmware 8.4 Identify Firewall,
http://www.cisco.com/en/US/docs/security/asa/asa84/configuration/guide/access_idfw.html
- [14] [http://en.wikipedia.org/wiki/Tor_\(anonymity_network\)](http://en.wikipedia.org/wiki/Tor_(anonymity_network))
- [15] http://en.wikipedia.org/wiki/IPv4_address_depletion
- [16] <http://en.wikipedia.org/wiki/IPv6>
- [17] Scott Hogg and Owen DeLong, “IPv6 NAT - You can get it, but you may not need or want it,” Infoblox Blog, October 2, 2012.
<http://www.infoblox.com/community/blog/ipv6-nat-you-can-get-it-you-may-not-need-or-want-it>
- [18] <http://en.wikipedia.org/wiki/Ipsec>
- [19] Scott Hogg, “Dual-Stack Will Increase Operating Expenses,” *Network World*, July 31, 2012,
<http://www.networkworld.com/community/blog/dual-stack-will-increase-operating-expenses>
- [20] http://en.wikipedia.org/wiki/Kerckhoffs%27s_Principle
- [21] <http://en.wikipedia.org/wiki/LifeLock>
- [22] Common Access Card (CAC), <http://www.cac.mil/>
- [23] Mat Honan, “Kill the P@55W0rD,” *WIRED Magazine*, December 2012,
<http://www.wired.com/gadgetlab/2012/11/ff-mat-honan-password-hacker/all/>
- [24] Geoff Huston, “Anatomy: A Look inside Network Address Translators,” *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [25] Several articles on IPv4 Exhaustion and IPv6 Transition in *The Internet Protocol Journal*, Volume 14, No. 1, March 2011.

SCOTT HOGG is the Director of Technology Solutions at GTRI in Denver Colorado. He holds a B.S. in Computer Science from Colorado State University, a M.S. in Telecommunications from the University of Colorado, and CCIE® #5133 and CISSP #4610 certifications. Hogg is active in the IPv6 community, Chair Emeritus of the RMv6TF, author of *IPv6 Security* (Cisco Press), a member of the Infoblox IPv6 Center of Excellence, a frequent presenter, and a *Network World* blogger. He can be reached at scott@hoggnetwork.com or followed on twitter [@scotthogg](https://twitter.com/scotthogg)

Summary Report of the ITU-T World Conference on International Telecommunications

by Robert Pepper and Chip Sharp, Cisco Systems

From 3–14 December, 2012, 151 Member States of the *International Telecommunication Union* (ITU) met in Dubai^[0] at the *World Conference on International Telecommunications* (WCIT-12)^[1] to revise the *International Telecommunication Regulations* (ITRs), a treaty-level document establishing policies governing international telecommunications services. During the 2-week conference the delegates debated several proposed changes on topics such as international mobile roaming, numbering, naming, addressing, fraud, the Internet, *Quality of Service* (QoS), etc. In the end, a revised version of the treaty was finalized^[2], but only 89 of the 151 Member States attending signed it.

There have been many articles discussing different aspects of the conference and its outcomes. This article provides background on the ITRs and focuses on the potential impact of the WCIT and its revised treaty on development of the Internet.

Background

The ITRs originated from the development of international telegraphy in Europe in the late 1800s and the need for a treaty defining how the government-operated national telegraph networks would interconnect and interoperate^[3]. As telephony and radio communications were invented, new treaties were developed to regulate their international operation. Up until the 1980s most telephone and telegraph companies were government-owned monopolies with some government licensed private companies operating as a monopoly. In 1988, the separate telegraph and telephone treaties were merged into the *International Telecommunications Regulations* while the *Radio Regulations* remained a separate treaty. By 1988, though some liberalization and privatization had started in a few countries in some regions, most international telecommunications services globally were still provided by monopoly, government-owned carriers, and services were dominated by voice rather than data. International Internet connectivity and traffic were practically nonexistent in most countries. Of course, international data traffic (including Internet) was growing in importance to some countries such as the United States and some large multinational companies such as IBM (which wanted to provide international *Virtual Private Networks* [VPNs]).

One important aspect of the ITRs in 1988 was the telephony accounting rate system. Briefly, this system consisted of a calling-party-pays business model for telephony in which the originating country pays the terminating country settlements based on a bilaterally agreed-upon accounting rate. Because developed countries tended to make more calls to developing countries than conversely and the accounting rate tended to be substantially above cost in many cases, the accounting rate system effectively became a subsidy program and a source for hard currency for developing countries.

Since 1988 market liberalization, reduced regulation, increased competition, and the rise of the Internet and mobile wireless industries have drastically changed the global communications landscape. In 1997, the U.S. *Federal Communications Commission* (FCC) opted out of the accounting rate system defined in the ITRs^[4], with many countries subsequently following suit. Voice over the Internet, arbitrage, hubbing, and other factors have reduced the telephony settlements revenue for developing countries. The 1988 ITRs^[5] allowed for special arrangements between network operators outside the rules of the ITRs. These special arrangements allowed for the international physical connectivity on which growth of the international Internet depended.

As the Internet grew and the telecom market changed, there was increased pressure from some countries to revise the ITRs. Contributions submitted in the preparatory meetings for WCIT-12 reflected widely varying views on the nature and extent of possible changes to the ITRs to account for this greatly changed environment. Although some countries believed that the ITRs should set forth high-level strategic and policy principles that could adapt to further changes in the market, others proposed the inclusion of expanded regulatory provisions of a detailed and specific nature in the ITRs to address a wide range of new concerns and services, including the Internet, or even to include the intergovernmental regulation of content (for example, spam and information security).

High-Level Take-Aways

Out of 151 countries attending the conference, the treaty was signed by 89 countries, consisting of mostly emerging countries led by Russia, China, Brazil, and the Arab States; 55 countries, including the United States, Japan, Australia, Canada, United Kingdom, and most of Europe, did not sign at the time. Countries that did not sign the treaty in Dubai can accede to the treaty after the WCIT by notifying the Secretary-General of the ITU. It is quite likely that some countries that did not sign the treaty will accede to it over the next few years.

The treaty takes effect on January 1, 2015 (after the 2014 *Plenipotentiary Conference*). Each signing country has to go through its national process for approval (for example, ratification) before the treaty takes effect for that country.

Although there has been a lot of negative commentary on the WCIT in the Internet community, in the end there are some important positive results for the Internet:

- No provisions were added to treaty text explicitly concerning the Internet, Internet Governance, or information security.
- No provisions were added to the treaty text concerning naming or addressing.
- No provisions modifying the basic business models of the Internet or mandating QoS on the Internet were made.

- The updated treaty explicitly recognizes commercial arrangements in addition to the old accounting rate regime for telecommunications.
- Article 9 on *Special Arrangements* allowing for telecommunications arrangements outside the treaty was retained mostly unchanged, thus allowing such special arrangements to continue to be used even between nonsignatory and signatory countries.
- A new resolution on landlocked countries could encourage access of such countries to landing stations in other countries and ease landlocked countries' ability to acquire international connectivity.

Some results that could be of concern to the Internet follow:

- The term identifying the operators to which the treaty applies ("authorized operating agencies") was modified. The supporters of the new term claim it does not expand the scope of the treaty, but it will bear watching.
- A provision on "unsolicited bulk electronic communications," developed after a long debate on spam, could lead governments to regulate and filter e-mail in addition to having unintended consequences such as disallowing bulk electronic emergency warning systems.
- Numbering provisions and requirements to deliver *Calling Party Number* were intended by some countries to allow for restrictions on international *Voice over IP* (VoIP) and VoIP services (including VoIP over the Internet).
- A new provision on network security could encourage more multilateral discussions in an intergovernmental setting (as opposed to multistakeholder).
- A new Resolution 3 on the Internet instructs the Secretary-General to engage further in Internet Governance discussions and further supports intergovernmental Internet policy processes.
- A new Resolution 5 mentions the transition to IP-based networks. It originally was aimed at over-the-top providers, but was modified to apply to service providers of international services. The end result is rather ambiguous in many respects and will bear watching.
- A new Article was added concerning telecommunication exchange points. Although the Internet is not mentioned explicitly, the originators of this article intended for it to apply to Internet Exchange Points. This Article could be used to support development of an enabling environment for regional telecommunication connectivity, but could also be used to justify regulation of Internet Exchange Points.
- Resolution Plen/4 requires PP'14 to consider a review of the ITRs every 8 years. This provision could result in another WCIT in 2020.

Table 1 lists the Member States that signed and did not sign the treaty in Dubai^[6].

Table 1: Treaty Signatories and Nonsignatories

Signatories			Nonsignatories	
Afghanistan	Guatemala	Qatar	Albania	Latvia
Algeria	Guyana	Russia	Andorra	Lichtenstein
Angola	Haiti	Rwanda	Armenia	Lithuania
Argentina	Indonesia	Saint Lucia	Australia	Luxembourg
Azerbaijan	Iran	Saudi Arabia	Austria	Malawi
Bahrain	Iraq	Senegal	Belarus	Malta
Bangladesh	Jamaica	Sierra Leone	Belgium	Marshall Islands
Barbados	Jordan	Singapore	Bulgaria	Moldova
Belize	Kazakhstan	Somalia	Canada	Mongolia
Benin	Korea (Rep. of)	South Africa	Chile	Montenegro
Bhutan	Kuwait	South Sudan	Colombia	Netherlands
Botswana	Kyrgyzstan	Sri Lanka	Costa Rica	New Zealand
Brazil	Lebanon	Sudan	Croatia	Norway
Brunei	Lesotho	Swaziland	Cyprus	Philippines
Burkina Faso	Liberia	Tanzania	Czech Republic	Poland
Burundi	Libya	Thailand	Denmark	Peru
Cambodia	Malaysia	Togo	Estonia	Portugal
Cape Verde	Mali	Trinidad and Tobago	Finland	Serbia
Central African Rep.	Mauritius	Tunisia	France	Slovak Republic
China	Mexico	Turkey	Gambia	Slovenia
Comoros	Morocco	Uganda	Georgia	Spain
Congo	Mozambique	Ukraine	Germany	Sweden
Cote d'Ivoire	Namibia	UAE	Greece	Switzerland
Cuba	Nepal	Uruguay	Hungary	United Kingdom
Djibouti	Niger	Uzbekistan	India	United States
Dominican Rep.	Nigeria	Venezuela	Ireland	
Egypt	Oman	Vietnam	Israel	
El Salvador	Panama	Yemen	Italy	
Gabon	Papua New Guinea	Zimbabwe	Japan	
Ghana	Paraguay		Kenya	

Note: Other *United Nations* (UN) member states were not eligible to sign or did not attend the conference but might still accede to the treaty: Antigua and Barbuda, Bahamas, Bolivia, Bosnia and Herzegovina, Cameroon, Chad, Dem. People’s Republic of Korea, Dem. Rep. of the Congo, Dominica, Ecuador, Equatorial Guinea, Eritrea, Ethiopia, Fiji, Grenada, Guinea, Guinea-Bissau, Honduras, Iceland, Kiribati, Lao P.D.R., T.F.Y.R. Macedonia, Madagascar, Maldives, Mauritania, Micronesia, Monaco, Myanmar, Nauru, Nicaragua, Pakistan, Romania, Saint Kitts and Nevis, Saint Vincent and the Grenadines, Samoa, San Marino, Sao Tome and Principe, Seychelles, Solomon Islands, Suriname, Syria, Tajikistan, Timor-Leste, Tonga, Turkmenistan, Tuvalu, Vanuatu, the Vatican, and Zambia.

Proposals and Outcomes

When the conference began there were several provisions that either explicitly or implicitly applied to the Internet, including:

- A proposal to define the term “Internet” and explicitly bring the Internet into the regulatory structure of the treaty
- Proposals to bring Internet naming, addressing, and identifiers into the treaty
- A proposal to include a provision on access to Internet websites
- A proposal on “traffic exchange points” that was intended to apply to *Internet Exchange Points*
- Proposals from multiple states on spam, information security, and *cybersecurity*

Although the Secretary-General of the ITU declared that the WCIT was not about the Internet or Internet Governance^[7], by rule, the WCIT had to consider input from its Member States. Given that Member States submitted proposals on the Internet, the Internet and Internet Governance was a substantive topic of discussion.

The following sections provide a brief review of some of the more difficult discussions related to the Internet.

Security

There were several proposals^[8] going into the WCIT to include cybersecurity, including information security, in the new ITRs. These proposals generated significant discussions and negotiations during the conference. The final text (Article 5A) is a great improvement over the proposals into the conference in that it focuses on the security and robustness of networks and prevention of technical harm to networks, with no mention of information security or cybersecurity.

The new provision mentions that Member States shall “collectively endeavour,” a provision that could engender more multilateral discussions in an intergovernmental setting (for example, ITU).

Organizations to Which the Treaty Applies

The 1988 ITR treaty focused on licensed carriers and government-owned *Post, Telephone, and Telegraph* (PTT) entities. Proposals^[8] into the WCIT would have applied the treaty to a wider range of organizations and companies. In the end, the treaty developed a new term, *Authorized Operating Agencies* (AOA). The proponents of this new term argued that it does not broaden the scope of the ITRs in terms of the organizations to which it applies. This interpretation of the new term should be supported, but monitored.

Internet-Specific Proposals and Resolutions (Resolutions Plen/3 and Plen/5)

Proposals^[8] were submitted to the WCIT to define the term “Internet” and to encode into the treaty the right of countries to regulate the “national segment” of the Internet. At the end of the first week of WCIT, Algeria, Saudi Arabia, Bahrain, China, United Arab Emirates, Iraq, Sudan, and Russia announced development of a new draft set of Resolutions that contained provisions that Member States shall have the right to manage the Internet, including Internet numbering, naming, addressing, and identification resources.

Although the United States, United Kingdom, and others were successful in removing any mention of the Internet from the treaty text, Internet-related language was moved into a nonbinding resolution (*Resolution Plen/3*) proposed by Russia “to foster an enabling environment for the greater growth of the Internet.” Resolution Plen/3 instructs the ITU Secretary-General “to continue to take the necessary steps for ITU to play an active and constructive role in the development of broadband and the multistakeholder model of the Internet as expressed in § 35 of the *Tunis Agenda*.” It also invites Member States to elaborate their positions on Internet-related concerns in the relevant ITU-related fora (something they could have done anyway).

This does not look too bad until one reads Paragraph 35 of the Tunis Agenda^[9]. This paragraph lays out the roles of each type of stakeholder (private industry, civil society, *Intergovernmental Organizations* [IGOs], governments, etc.). It reserves an explicit role in “Internet-related public policy issues” for governments and intergovernmental organizations. It does not provide for any role in this area for the private sector or civil society. So although the Resolution seems to support the multistakeholder model of the Internet, it really restricts the roles of several of the main stakeholders.

Several countries pushed for inclusion of Paragraph 55 of the Tunis Agenda, recognizing that the existing arrangements have worked effectively, to balance the inclusion of Paragraph 35, but it was not included in the final Resolution.

Resolution Plen/3 may be used by some governments to reinforce the ITU’s role in Internet Governance, including at future ITU conferences in 2013 and 2014.

On the other hand, the Resolution also instructs the Secretary-General “to support the participation of Member States and *all* other stakeholders, as applicable, in the activities of ITU in this regard.” This statement supports participation of all stakeholders in the activities of the ITU, not restricted just to ITU Members, or in the case of ITU Council or some Council Working Groups just to Member States.

In signing the Final Acts, Russia added a Declaration/Reservation that it views the Internet as a new global telecommunication infrastructure and reserves the right to implement public policy, including international policy, on matters of Internet Governance. This reservation could signal that Russia plans to apply the telecommunications provisions in the ITRs to the Internet and to further regulate the Internet.

In addition to Resolution Plen/3, some of the proposals on the Internet were part of the discussion on Resolution Plen/5. This Resolution began as a basic resolution on invoicing for international telecommunication services, but ended up including numerous other provisions that did not make it into the main text of the treaty. Although the final text does not contain provisions explicitly mentioning the Internet, the introductory text of the Resolution mentions the transition of phone and data networks to IP-based networks. Also, the proposal that evolved into “resolves” originally applied to the relationship between network operators and application providers. During the discussions this proposal was modified to refer to “providers of international services” instead of application providers. Even with this modification, the application of this provision is ambiguous and could be applied to over-the-top providers.

Resolution Plen/5 is likely to reinforce work in Study Group 3 on accounting, fraud and charges for international telecommunications service traffic termination and exchange, etc.

Telecommunications Traffic Exchange Points

A proposal^[8] concerning “telecommunication traffic exchange points” was included as an Article in the ITRs. The term “telecommunication traffic exchange point” was left undefined. This article does not mention the Internet or Internet Exchange Points, but the discussion of this Article included discussion on how it related to Internet Exchange Points. At least one delegation indicated that the Article was intended to help enable development of regional Internet Exchange Points.

Although this provision raised concerns over possible regulation of Internet Exchange Points, it focuses on creating an enabling environment for creation of regional telecommunication traffic exchange points. This environment could provide support for development of trans-border telecommunications and connectivity.

Route-Related Factors

Prior to the conference, there were several proposals [8] to require transparency into the international routes used for a Member States' traffic and to allow Member States to control what routes were used between them. Note that the definition of "route" in the ITRs is different from the concept of a "route" on the Internet. In the ITRs, a route is defined as the technical facilities used for telecommunications traffic between two telecommunication terminal exchanges or offices.

Coming into the WCIT, the proposal to control routes by Member States was dropped from the proposal, so the debate centered over whether Member States should have the right to know what routes were being used. After much discussion, the final result was a provision allowing "authorized operating agencies" (not Member States) to determine which routes are to be used between them and allowing the originating operator to determine the outbound route for traffic. This provision is not much different from how network operators manage their networks today.

Quality of Service Proposals

Several proposals^[8] were made to WCIT to require QoS to be negotiated between network providers including Internet providers. Some proposals also allowed network providers to charge over-the-top providers for QoS.

The final provisions did not add any new requirements for QoS other than a nonspecific requirement related to mobile roaming. Although no new provisions were added specific to the Internet, it does not mean that countries could not try to impose the current QoS provisions to VoIP services. The debate over QoS on the Internet will continue outside the ITRs.

Naming, Numbering, and Addressing Proposals

Several countries and regions proposed^[8] to extend provisions on telephone numbering to include naming, addressing, and origin identifiers. Several proposals were made to require delivery of calling party number and to cooperate in preventing the misuse ("misuse" not defined) of numbering, naming, and addressing resources. Although the Internet was not explicitly mentioned, these proposals were intended to apply to VoIP based on comments at pre-WCIT preparatory meetings^[10].

In the end, several provisions were added related to delivery of calling party number and prevention of misuse of telecommunications numbering resources as defined in ITU-T Recommendations. Provisions to include naming, addressing, and more general "origin identifiers" were not accepted.

Even though there were no provisions specifically on the Internet, some countries could apply these provisions to VoIP services that use E.164 telephone numbers and that provide for bypass of the international telephony accounting system.

However, it is not clear that these provisions add any more authority than what these countries have today.

Content and Spam

Proposals^[8] to include spam in the treaty caused a lot of contentious discussion, in ad hoc groups, plenary, and in consultations. Some countries took a strong position that spam is a content topic that was out of scope of the ITRs. There was a concern that adding a provision on spam would legitimize content filtering by governments. Some African countries insisted on including a provision on spam, claiming that it consumed a large percentage of their international bandwidth. In the end to address concern about content, a statement was added to Article 1.1:

“These Regulations do not address the content-related aspects of telecommunications.”

To address the proposals on spam, Article 5B was added on unsolicited bulk electronic communication:

“Member States should endeavour to take necessary measures to prevent the propagation of unsolicited bulk electronic communications and minimize its impact on international telecommunication services. Member States are encouraged to cooperate in that sense.”

As written the final text, it is fairly vague and could have implications beyond spam; for example, there are no exemptions for broadcasters or for emergency alert systems (for example, tsunami alerts). It is also not clear how Article 5B can be implemented consistent with the statement on content in Article 1.1.

It was clear from the discussion that many of the delegates from countries supporting this provision do not understand spam or spam-mitigation techniques and their usage (or not) in their own countries. It is clear that many of the delegates were not aware of basic best practices from the *Messaging Anti-Abuse Working Group* (MAAWG) and other organizations. These discussions highlighted the need for capacity building for developing countries on spam-mitigation techniques.

Human Rights and Member State Access to International Telecommunications

In a plenary session on the penultimate night of the WCIT, a provision on human rights was added to the final draft of the ITRs. This discussion led to a debate concerning the right of Member States to access international telecommunication services, originating from a proposal from Sudan and Cuba creating a right of Member States to access Internet websites. This provision was targeted at U.S. and European actions taken in response to UN sanctions against Sudan due to Darfur and U.S. sanctions on Cuba.

The provision provides a right for Member States, not its citizens. Thus it did not provide any rights for citizens to access international telecommunication services. In addition, it is not clear what or whose international telecommunication service Member States have a right to. The implications of the provision were unclear, and delegations did not have time to consult their home countries before the end of the conference.

Several times during the debate the Chair of the WCIT and the Secretary-General of the ITU both tried to dissuade the proponents from pushing their proposal, to no avail. After extended debate, Iran called for a point of order and then called for a vote, the only official vote of the conference. After the text passed by majority vote, the Chair of the WCIT declared the ITRs approved. At that point the United States, followed by the United Kingdom, Sweden, and other countries, made statements that they would not sign the treaty. Supporters of the treaty read their statements in favor of the treaty. The conference was effectively over^[11].

The uncertainty caused by the addition of this text at such a late date and the way it was added created a situation in which many countries that might have signed the treaty ended up not signing. This provision more than any other disagreement in the conference caused the conference to split to the extent that it did.

Looking Forward

Much of the long-term impact of the treaty will not be felt until the signing governments ratify the treaty and start enacting provisions into either law or regulation. It is likely that some of the countries that did not sign in Dubai will accede to the treaty at a later time, including countries that did not attend the WCIT.

WCIT is only one step (though an important one) in the long-term debate over Internet Governance and the appropriate role of governments (and intergovernmental organizations) in the Internet. The debate will continue in numerous international fora going forward such as:

- World Telecommunications Policy Forum (May 2013)
- World Summit on the Information Society Action Line Forum (May 2013)
- ITU Council Working Group on Internet Public Policy (ongoing)
- ITU-T Study Group meetings (ongoing)
- ITU Plenipotentiary Conference (2014)
- WSIS+10 Review (2013–2015)

It has already been seen that many of the same topics debated at WCIT will be debated in these venues; for example, IP addressing, naming, spam, and cybersecurity. The WCIT Resolutions (especially Res. Plen/3) will likely be used to promote a larger role of the ITU in the Internet Governance debate.

The ITU's Plenipotentiary Conference in 2014 will be the next important treaty conference where the ITU's Constitution and Convention (both treaty instruments) can be revised. In the hierarchy of treaties at ITU, the ITU Constitution takes precedence over the ITRs, and many of the terms used in the ITRs are defined in the Constitution. Therefore, changes to the ITU Constitution could affect the meaning of the ITRs. The ITU Plenipotentiary will provide an opportunity for the ITU Member States to come together and heal some of the differences coming out of the WCIT, but it is also an opportunity to widen the rift.

The WSIS+10 Review will be an important process because it is likely to set the agenda for the discussion of Internet Governance for the 5–10 years after 2015, much as the Tunis Agenda from 2005 set the agenda for the last 8 years. An important aspect of the WSIS+10 Review is that it involves other UN agencies (for example, UNESCO) in addition to the ITU. Many of the events involve stakeholders whose voices are not normally heard at ITU conferences.

Some of the disagreements exhibited at WCIT brought to light opportunities for the Internet community to engage with governments and other stakeholders by providing technical and thought leadership. Capacity building with many of the developing country governments will be an important part of the preparation leading up to the major international conferences such as the ITU Plenipotentiary and WSIS+10.

Much of the growth of the Internet going forward is likely to come in the countries that signed the ITRs. Many of these countries have started to develop multistakeholder consultations and processes when dealing with Internet topics. The fact that a government signed the ITRs does not mean that the country is somehow against the Internet. On the contrary, many of these countries are looking for ways to accelerate the Internet's development within their borders and to accelerate their international connectivity to the Internet. As the Internet grows and develops in these countries, the Internet communities in these countries will likely look to play a larger role in a consultative process regarding government positions on issues related to Internet Governance. Future growth of the Internet across ITR boundaries (signatories and non-signatories) will depend on cooperation amongst all stakeholders.

References

- [0] Geoff Huston, “December in Dubai: Number Misuse, WCIT, and ITRs,” *The Internet Protocol Journal*, Volume 15, No. 2, June 2012.
- [1] World Conference on International Telecommunications (WCIT-12),
<http://www.itu.int/en/wcit-12/Pages/default.aspx>
- [2] International Telecommunication Regulations,
<http://www.itu.int/en/wcit-12/Pages/itrs.aspx>
- [3] “Discover ITU’s History,” <http://www.itu.int/en/history/Pages/DiscoverITUsHistory.aspx>
- [4] “International Settlements Policy and U.S.–International Accounting Rates,”
<http://www.fcc.gov/encyclopedia/international-settlements-policy-and-us-international-accounting-rates>
- [5] “WATTC-88 World Administrative Telegraph and Telephone Conference (Melbourne, 1988),”
<http://www.itu.int/en/history/Pages/TelegraphAndTelephoneConferences.aspx?conf=33&dms=S0201000021>
- [6] “Signatories of the Final Acts: 89 (in green),”
<http://www.itu.int/osg/wcit-12/highlights/signatories.html>
- [7] “Dr. Hamadoun I. Touré, ITU Secretary-General First Plenary of World Conference on International Telecommunications (WCIT-12),”
<http://www.itu.int/en/wcit-12/Pages/speech-toure2.aspx>
- [8] “Proposals Received from ITU Member States for the Work of the Conference,”
http://www.itu.int/md/dologin_md.asp?lang=en&id=S12-WCIT12-121203-TD-0001!!MSW-E
- [9] “Tunis Agenda for the Information Society,” *World Summit on the Information Society*, 2005. <http://www.itu.int/wsis/docs2/tunis/off/6rev1.html>
- [10] Council Working Group to Prepare for the 2012 WCIT,
<http://www.itu.int/council/groups/cwg-wcit12/index.html>
- [11] Webcast and Captioning of the WCIT,
<http://www.itu.int/en/wcit-12/Pages/webcast.aspx>

ROBERT PEPPER leads Cisco's Global Technology Policy team working with governments across the world in areas such as broadband, IP enabled services, wireless and spectrum policy, security, privacy, Internet governance and ICT development. He joined Cisco in July 2005 from the FCC where he served as Chief of the Office of Plans and Policy and Chief of Policy Development beginning in 1989 where he led teams developing policies promoting the development of the Internet, implementing telecommunications legislation, planning for the transition to digital television, and designing and implementing the first U.S. spectrum auctions. He serves on the board of the U.S. Telecommunications Training Institute (USTTI) and advisory boards for Columbia University and Michigan State University, and is a Communications Program Fellow at the Aspen Institute. He is a member of the U.S. Department of Commerce's Spectrum Management Advisory Committee, the UK's Ofcom Spectrum Advisory Board and the U.S. Department of State's Advisory Committee on International Communications and Information Policy. Pepper received his BA. and Ph.D. from the University of Wisconsin-Madison. E-mail: [**rmpepper@cisco.com**](mailto:rmpepper@cisco.com)

CHIP SHARP has 30 years in the communications industry and currently is a Director in the Research and Advanced Development Department at Cisco Systems, Inc. His current role is in Technology Policy focusing on Internet Governance issues. He participated in the US preparatory process for WCIT from 2010 including as Private Sector Advisor to the US Delegation to the ITU's Council Working Group to Prepare for the 2012 WCIT (CWG-WCIT12), mainly analyzing the impact of proposals on the Internet and Internet Governance. He helped develop many of the talking points and position papers related to the Internet for the US Delegation. He served in the same capacity on the US Delegation to WCIT. He continues to be active in many follow-on activities preparing for the World Telecommunication Policy Forum (WTPF), World Summit on the Information Society 10 year review (WSIS+10), ITU Plenipotentiary Conference 2014 and Internet Governance Forum. He also currently service on the FCC's Open Internet Advisory Committee. Prior to this role, he led a multinational, multidisciplinary team at Cisco helping drive various technologies such as LISP, DNSSEC, BGPSEC, ENUM, Lawful Intercept etc. He has also supported capacity building and development programs for developing countries, for example, deployment of Internet Exchange Points (IXPs). He started at Cisco in 1996 helping design dialup Internet access products to interface with legacy telco signaling systems. Prior to Cisco, Chip worked at Teleos Communications, AT&T Consumer Product Labs and NASA's Communications Division. E-mail: [**chsharp@cisco.com**](mailto:chsharp@cisco.com)

Letters to the Editor

Dear Ole,

I am sorry that there is some delay (more than 1 second) between the arrival of *The Internet Protocol Journal* at my desk and this e-mail. In the December 2012 issue (Volume 15, No. 4), Geoff Houston discusses the extra second on the last minute of the 31st of June. There is no 31st of June in the calendar, at least not in old Europe, but maybe in the United States. It is funny to discuss the problem of a second at the end of a nonexistent day, isn't it?

Nevertheless I could take some new knowledge from this article.

Best regards,

—Richard Schuerger
`richard.schuerger@gmx.de`

Hi Geoff (and Ole)!

I am sitting comfortably in a chair on the terrace in a Tenerife house, reading the December 2012 issue of IPJ, which I received by mail today. Since I have been working many years with the *Network Time Protocol* (NTP), I started reading your article on the subject with great interest. Having read only a few sentences I jumped in my chair:

“Back at the end of June 2012 there was a brief IT hiccup as the world adjusted the *Coordinated Universal Time* (UTC) standard by adding an extra second to the last minute of the 31st [!!] of June.”

Of course you may have received numerous notices of this hiccup [ha, ha], but still I couldn't resist writing to you. Thank you for an [otherwise] well-written and clarifying article (as always).

—Truls Hjelle
`truls@sund-hjelle.org`

PS: Thanks to Ole for this anachronism on paper still available to us oldies who prefer sitting with a paper magazine in the sun instead of gazing at a poorly lit screen and struggling with the tiny letters.


The author responds:

Back in 45 BC, Julius Caesar made same revolutionary changes to the Roman calendar, and the changes included adding one extra day to June (well not quite, as the letter “J” was not around until the 16th Century, and the letter “u” was also yet to make its debut, so it is probably less of an anachronism to record that Gaius Iulius Caesar added an extra day to the month of Iunius). Either way, this change brought the total number of days in the month of June to 30, which is where it has remained for 2058 years.

It is often said that Australia operates on a calendar all of its own, but while our isolation on a largish rock at the southern end of the Pacific Ocean has led to a number of revolutionary innovations that are easily on a par with fire and the wheel, including the world-renowned stump-jump plough and the sheep-shearing machine, we Australians have not yet turned our collective national genius to the calendar. Despite a pretty sensible suggestion from the latest meeting of the Grong Grong Shire Council for a year to be made up of 10 months of 30 days followed by a decent 65-day session at the pub, we have yet to get the blokes back from the pub after their last 65-day bender, so that plan needs some more work back at the shed before it gets another airing! Thus it looks like Australia uses the same calendar as everyone else, making the reference to the 31st of June one of those pesky brain-fade errors! Oops. Yes, it was meant to say 30th of June. Well spotted!

—*Geoff Huston*

gih@apnic.net

 The Internet Protocol Journal, Cisco Systems 170 West Tasman Drive San Jose, CA 95134-1706 USA ADDRESS SERVICE REQUESTED	<div>PSRST STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA</div> <div> FOSTER BUNNY LAGOMORPH INC 1234 MAIN STREET SAN FRANCISCO, CA 94104-1234</div>
SUBSCRIPTION ID: FBUNN006188 V15 N3 EXPIRATION DATE: 15-JAN-2013	

Don't forget to renew and update your subscription. For details see the IPJ Subscription FAQ in our previous issue (Volume 15, No. 4).

Book Review

On Internet Freedom

On Internet Freedom, by Marvin Ammori, Elkat Books, January 2013, sold by: Amazon Digital Services, Inc., ASIN: B00B1MQZNW.

Marvin Ammori has written an important book about the threats to free speech and expression that we are not only privileged to conduct on the Internet today but have come to treat as basic human rights.

On Internet Freedom looks at the past, present, and future of the Internet as a speech technology. Ammori examines how the coordinated and determined efforts by Big Content to protect content and increasing efforts by governments to censor content threaten Internet use as we embrace it today. Ammori also explains how these acts were in fact anticipated by Clark, Sollins, Wroclawski, and Braden in a paper entitled “Tussle in Cyberspace: Defining Tomorrow’s Internet,”^[1] where the authors assert:

“User empowerment, to many, is a basic Internet principle, but for this paper, it is the manifestation of the right to choose—to drive competition, and thus drive change.”

Ammori cites only the first clause of this sentence—as a technologist, I believe the second is extremely important as well—but he makes clear that the end-to-end design of the Internet establishes a fundamental thesis:

“If user choice is our design principle, then users should have the final say.”

Unfortunately, Ammori explains that users do not have the final say but are increasingly challenged by lawyers, bureaucrats, commissioners, and others who are motivated to constrain their freedoms and who want to do so by altering the fundamental design of the Internet. Ammori’s response, admittedly U.S.-centric, is simple: the Internet is a speech technology, and:

“... the ultimate design principle for any speech technology, at least in the United States: the First Amendment, which protects freedom of speech. The *First Amendment* is not generally thought of as a design principle, but, by definition, it limits what Congress or any other government actor may or may not adopt in shaping the Internet’s future.”

This statement sets the context for the remainder of the book. In Part II, Ammori looks at events leading to the 18 January 2012 Internet Blackout in protest of the *Stop Online Piracy Act* (SOPA) and *PROTECT IP Act* (PIPA) and how these and possibly future legislation threaten “...the speech tools of the many while reshaping our speech environment for the benefit of the few.”

Conveniently, Part II is largely about how the few benefit. Before judging whether you believe this theory is even-handed or not, remember that the litmus test throughout this book is the First Amendment of the U.S. Constitution. This part ought to make every Internet user or free speech advocate pause, or shiver. One of the most worrisome speculations Ammori offers is the extent to which legislation could stilt adoption of emerging technologies such as *three-dimensional* (3D) printing or stifle future innovations of this kind.

Part III looks at how the Internet as speech technology influences governments, how governments have attempted to exert influence, and how Internet users and dominant Internet forces (Google, Amazon, Facebook, and Twitter) respond. This part will probably be illuminating for most readers, because it explains situations where a *private conversation* between a government official and an *Internet Service Provider* (ISP) or hosting company can circumvent the First Amendment, and why *Terms of Service* are often more speech-restricting than the First Amendment as well.

Part IV focuses on net neutrality concerns. Ammori draws the lines of conflict: ISPs seek to differentiate, rate-control, block, or charge users differently for content that is transmitted on their networks. However, content includes speech, and if the Internet is speech technology, then ISPs should not be able to decide what you say or see, or they do so in violation of your First Amendment rights. Ammori also explains that net neutrality is not only a First Amendment concern but also an economic one: net neutrality violations can influence investments in or creation of new technology.

I began by saying that Marvin Ammori has written an important book. It is also an extremely readable book. Ammori does a commendable job explaining constitutional law and technology in easy to understand terms. I highly recommend the book not only for people who are interested in law or technology but for anyone who advocates freedom of expression.

On Internet Freedom is currently available as a Kindle download.

- [1] David D. Clark, John Wroclawski, Karen R. Sollins, and Robert Braden, "Tussle in Cyberspace: Defining Tomorrow's Internet," *IEEE/ACM Transactions on Networking*, Volume 13, Issue 3, June 2005. Available from:
<http://groups.csail.mit.edu/ana/Publications/PubPDFs/Tussle2002.pdf>

—Dave Piscitello, dave@corecom.com

Reprinted with permission from *The Security Skeptic* blog:
<http://securityskeptic.typepad.com/the-security-skeptic/>

Nominations Sought for 2013 Jonathan B. Postel Service Award

The Internet Society is soliciting nominations of qualified candidates for the 2013 *Jonathan B. Postel Service Award* by May 31, 2013. This annual award is presented to an individual or organization that has made outstanding contributions in service to the data communications community. The award is scheduled to be presented during the 87th IETF meeting in Berlin, Germany, July 28–August 2.

The award was established by the Internet Society to honor a person who has made outstanding contributions in service to the data communications community. The award is focused on sustained and substantial technical contributions, service to the community, and leadership. With respect to leadership, the award committee places particular emphasis on candidates who have supported and enabled others in addition to their own specific actions.

The award is named for Dr. Jonathan B. Postel to recognize and commemorate the extraordinary stewardship exercised by Jon over the course of a thirty-year career in networking. He served as the editor of the RFC series of notes from its inception in 1969 until 1998. He also served as the ARPANET “Numbers Czar” and *Internet Assigned Numbers Authority* (IANA) over the same period of time. He was a founding member of the Internet Architecture (nee Activities) Board and the first individual member of the Internet Society, which he also served as a Trustee.

For more information, see: <http://www.internetsociety.org/>

Upcoming Events

The *Internet Corporation for Assigned Names and Numbers* (ICANN) will meet in Beijing, China, April 7–11, 2013 and in Durban, South Africa, July 14–18, 2013. For more information, see: <http://icann.org/>

The *North American Network Operators’ Group* (NANOG) will meet in New Orleans, Louisiana, June 3–5, 2013 and in Phoenix, Arizona, October 7–9, 2013. For more information see: <http://nanog.org>

The *Internet Engineering Task Force* (IETF) will meet in Berlin, Germany, July 28–August 2, 2013 and in Vancouver, Canada, November 3–8, 2013. For more information see: <http://www.ietf.org/meeting/>

The *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) will meet in Bangkok, Thailand, February 18–28, 2014. For more information see: <http://www.apricot.net>

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ contains standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2013 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

June 2013

Volume 16, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Network Service Models	2
Looking Forward.....	14
Link-State Protocols in Data Center Networks	23
Letter to the Editor	30
Book Review.....	31
Fragments	35
Call for Papers.....	39

You can download IPJ
back issues and find
subscription information at:
www.cisco.com/ipj

ISSN 1944-1134

FROM THE EDITOR

Fifteen years ago we published the first edition of *The Internet Protocol Journal* (IPJ). This seems like a good time to reflect on where the Internet is today and where it might be going in the future, instead of looking back at earlier developments the way we did in the tenth anniversary issue of IPJ.

In our first article, Geoff Huston discusses network service models, comparing the Internet to the traditional *Public Switched Telephone Network* (PSTN) in both technical and business terms, and asks if the fundamental architectural differences between these networks might explain the rather slow deployment of IPv6. Although the number of IPv6-connected users has doubled in the last year (see page 35), IPv6 still represents a small percentage of total Internet traffic.

The mobile device dominates today's Internet landscape. Smartphones and tablets are starting to replace more traditional computers for Internet access. Many technical developments have made this possible, including high-resolution screens; powerful processors; and compact, long-lasting batteries. Combine such developments with numerous radio-based technologies (GPS, cellular, Wi-Fi, and Bluetooth) and you end up with a handheld device that is always connected to the network and can perform almost any task, using an appropriate "app." Improvements to communications technologies such as the deployment of *Long-Term Evolution* (LTE) cellular data networks and *Gigabit Wi-Fi* (IEEE 802.11ac) are already underway.

We asked Vint Cerf, known to many as one of the "Fathers of the Internet," to look beyond what is possible with today's Internet and today's devices and predict what the future might look like in a world where every imaginable appliance is "smart," connected to the network, and location-aware. His article takes us through some history and current trends, and then describes how the future Internet might shape many aspects of society such as business, science, and education.

According to Wikipedia, a *Data Center* is "a facility used to house computer systems and associated components, such as telecommunications and storage systems." In our final article, Alvaro Retana and Russ White discuss how developments in link-state protocols, usually associated with wide-area networks, can be applied to data center networks.

—Ole J. Jacobsen, Editor and Publisher
ole@cisco.com

Network Service Models and the Internet

by Geoff Huston, APNIC

In recent times we've covered a lot of ground in terms of the evolution of telecommunications services, riding on the back of the runaway success of the Internet. We have taken the computer and applied a series of transformational changes in computing power and size, battery technology, and added radio capabilities to create a surprising result. We've managed to put advanced computation power in a form factor that fits in the palms of our hands, and have coupled it with a communications capability that can manage data flows of tens if not hundreds of megabits per second—all in devices that have as few as two physical buttons! And we have created these devices at such scale that their manufacturing cost is now down to just tens of dollars per unit. The Internet is not just at the center of today's mass market consumer service enterprise, it is now at the heart of many aspects of our lives. It's not just the current fads of the social networking tools, but so much more. How we work; how we buy and sell, even what we buy and sell; how we are entertained; how democracies function, even how our societies are structured; and so much more—all of these activities are mediated by the Internet.

But a few clouds have strayed into this otherwise sunny story of technological wonder. Perhaps the largest of these clouds is that the underlying fabric of the Internet, the numbering plan of the network, is now fracturing. We have run out of IP addresses in the Asia Pacific region, Europe, and the Middle East. At the same time, the intended solution, namely the transition to a version of the IP protocol with a massively larger number space, IPv6, is still progressing at an uncomfortably slow pace. Although the numbers look like a typical “up and to the right” Internet data series, the vertical axis tells a somewhat different story. The overall deployment of IPv6 in today's Internet currently encompasses around 1.3 percent^[0] of the total user base of the Internet, and it is possible that the actions of the open competitive market in Internet-based service provision will not necessarily add any significant further impetus to this necessary transition.

We have gone through numerous phases of explanation for this apparently anomalous success-disaster situation for the Internet. Initially, we formed the idea that the slow adoption of IPv6 was due to a lack of widely appreciated knowledge about the imminent demise of IPv4 and the need to transition the network to IPv6. We thought that the appropriate response would be a concerted effort at information dissemination and awareness rising across the industry, and that is exactly what we did. But the response, as measured in terms of additional impetus for the uptake of IPv6 in the Internet, was not exactly overwhelming.

We then searched for a different reason as to why this IPv6 transition appeared to be stalling. There was the thought that this problem was not so much a technical one as a business or a market-based one, and there was the idea that a better understanding of the operation of markets and the interplay between markets and various forms of public sector initiatives could assist in creating a stronger impetus for IPv6 in the service market. The efforts at stimulation of the market to supply IPv6 goods and services through public sector IPv6 purchase programs have not managed to create a “tipping point” for adoption of IPv6.

Some have offered the idea that the realization of IPv4 exhaustion would focus our thinking and bring some collective urgency to our actions. But although IPv4 address exhaustion in the Asia Pacific region in 2011 has created some immediate interest in IPv4 address extension mechanisms, the overall numbers on IPv6 adoption have stubbornly remained under 1.5 percent of the 2 billion user base of the Internet.

Why has this situation occurred? How can we deliberately lead this prodigious network into the somewhat perverse outcomes that break to basic end-to-end IP architecture by attempting to continue to overload the IPv4 network with more and more connected devices? What strange perversity allows us to refuse to embrace a transition to a technology than can easily sustain the connection needs of the entire silicon industry for many decades to come and instead choose a path that represents the general imposition of additional cost and inefficiency?

Perhaps something more fundamental is going on here that reaches into the architectural foundations of the Internet and may explain, to some extent, this evident reluctance of critical parts of this industry to truly engage with this IPv6 transition and move forward.

Telephony Network Intelligence

Compared to today’s “smart” phone, a basic telephone handset was a remarkably basic instrument. The entire telephone service was constructed with a model of a generic interface device that was little more than a speaker, a microphone, a bell, and a pulse generator. The service model of the telephone, including the call-initiation function of dialing and ringing, the real-time synchronous channel provision to support bidirectional speech, all forms of digital and analogue conversion, and of course the call-accounting function, were essentially all functions of the network itself, not the handset. Although the network was constructed as a real-time switching network, essentially supporting a model of switching time slots within each of the network switching elements, the service model of the network was a “full-service” model.

The capital investment in the telecommunications service was therefore an investment in the network—in the transmission, switching, and accounting functions.

Building these networks was an expensive undertaking in terms of the magnitude of capital required. By the end of the 20th century the equipment required to support synchronous time switching included high-precision atomic time sources, a hierarchy of time-division switches to support the dynamic creation of edge-to-edge synchronous virtual circuits, and a network of transmission resources that supported synchronous digital signaling. Of course although these switching units were highly sophisticated items of technology, most of this investment capital in the telephone network was absorbed by the last mile of the network, or the so-called “local loop.”

Although the financial models to operate these networks varied from operator to operator, it could be argued that there was little in the way of direct incremental cost in supporting a “call” across such a network, but there is a significant opportunity or displacement cost. These networks have a fixed capacity, and the requirements for supporting a “call” are inelastic. When a time slot is being used by one call, this slot is unavailable for use by any other call.

Telephony Tariffs

Numerous models were used when a retail tariff structure for telephony was constructed. One model was a “subscription model,” where, for a fixed fee, a subscriber could make an unlimited number of calls. In other words the operator’s costs in constructing and operating the network were recouped equally from all the subscribers to the network, and no transaction-based charges were levied upon the subscriber. This model works exceptionally well where the capacity of the network to service calls is of the same order as the peak call demand that is placed on the network. In other words, where the capacity of the network is such that the marginal opportunity or displacement cost to support each call is negligible, there is no efficiency gain in imposing a transactional tariff on the user. In the United States’ telephone network, for example, a common tariff structure was that the monthly telephone service charge also allowed the subscriber to make an unlimited number of local calls.

Another model in widespread use in telephony was of a smaller, fixed service charge and a per-transaction charge for each call made. Here a subscriber was charged a fee for each call (or “transaction”) that the subscriber initiated. The components to determine the charge for an individual transaction included the duration of the call, the distance between the two end parties of the call, the time of day, and the day of the week. This model allowed a network operator to create an economically efficient model of exploitation of an underlying common resource of fixed capacity. This model of per-call accounting was widespread, used by some operators in local call zones, and more widely by telephone service operators in long distance and international calls.

This model allowed the operator to generate revenue and recoup its costs from those subscribers who used the service, and, by using the pricing function, the network operator could moderate peak demand for the resource to match available capacity.

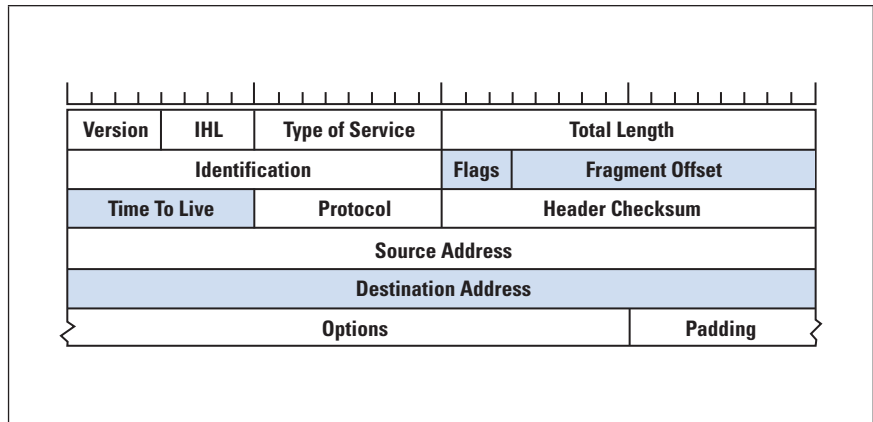
This per-transaction service model of telephony was available to the operator of the telephone service simply because the entire function of providing the telephone service was a network-based service. The network was aware of who initiated the transaction, who “terminated” the transaction, how long the transaction lasted, and what carriers were involved in supporting it. Initially this transactional service model was seen as a fair way to allocate the not inconsiderable costs of the construction and operation of the network to those who actually used it, and allocate these costs in proportion to the relative level of use. I suspect, though, that this fair cost allocation model disappeared many decades ago because these per-transaction service tariffs became less cost-based and more based on monopoly rentals.

IP Network Minimalism

The Internet is different. Indeed, the Internet is about as different from telephony as one could possibly imagine. The architecture of the Internet assumes that a network transaction is a transaction between computers. In this architecture the computers are highly capable signal processors and the network is essentially a simple packet conduit. The network is handed “datagrams,” which the network is expected to deliver most of the time. However, within this architecture the network may fail to deliver the packets, may reorder the packets, or may even corrupt the content of the packets. The network is under no constraint as to the amount of time it takes to deliver the packet. In essence, the expectations that the architecture imposes on the network are about as minimal as possible. Similarly, the information that the edge-connected computers now expose to the network is also very limited. To illustrate this concept, it is useful to look at the fields that the Internet Protocol exposes to the network.

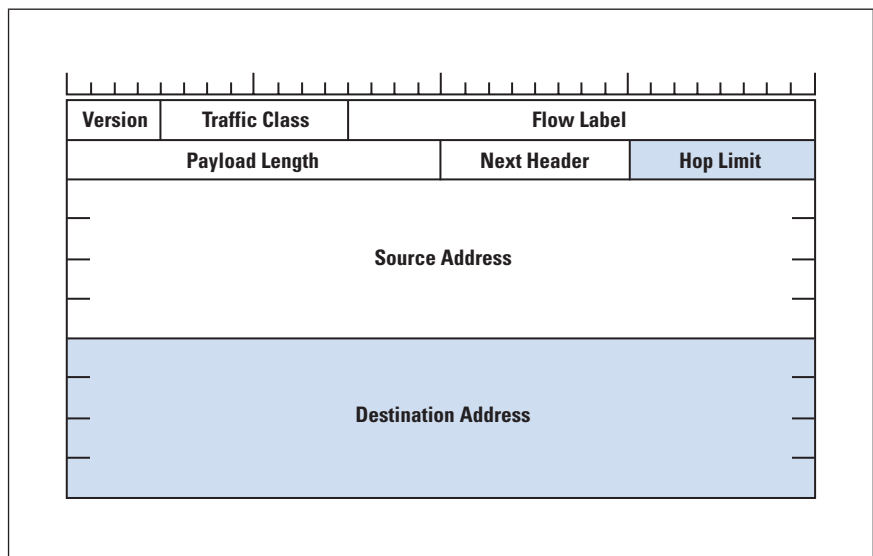
In IPv4 the fields of the Internet Protocol header are a small set, as shown in Figure 1. An IP packet header exposes the protocol *Version*, *Header Length* (IHL), *Total Length* of the IP packet, packet *Fragmentation Offset*, and *Type of Service* fields, a hop counter (*Time To Live* field), a *Header Checksum* field, and the *Source and Destination Address* fields. In practice, the Type of Service field is unused, and the Length and Checksum fields have information that is also contained in the data link frame header. What is left is the protocol Version field, packet length (Total Length field), the Fragmentation Offset field, a hop counter, and the Source and Destination Address fields. Of these fields, the Packet Length, Fragmentation Offset, hop counter, and Destination Address are the fields used by the network to forward the packet to its ultimate destination.

Figure 1: The IPv4 Packet Header



In IPv6 this minimal approach was further exercised with the removal of the Fragmentation Control fields and the Checksum fields (Figure 2). Arguably, the *Traffic Class* and *Flow Label* are unused, leaving only the *Protocol Version*, *Payload Length*, a *Hop Counter*, and the source and destination addresses exposed to the network. In IPv6 the minimal network-level information is now reduced to the packet length, the hop counter, and the destination address.

Figure 2: The IPv6 Packet Header



These fields represent the totality of the amount of information that the Internet Protocol intentionally exposes to the network. There are no transaction identifiers, no call initiation or call teardown signals, or even any reliable indication of relative priority of the packets. All the network needs to “see” in each carried packet is a hop counter, a packet length, and a destination address.

Within this model the actions of each of the network's switching elements are extremely simple, as shown in Figure 3.

Figure 3: IPv4 and IPv6 Packet Processing

```
for each received packet:
    decrement the hop counter
    if the counter value is zero then discard the packet, otherwise...
    look up the packet's destination address in a local table
    if the lookup fails then discard the packet, otherwise...
    look up the output queue from the located table entry
    if the queue is full discard the packet, otherwise...
    if the packet is too large for the outbound interface then
        fragment the packet to fit, if permitted (IPv4)
        or discard the packet (IPv6), otherwise...
    queue the packet for onward transmission
```

The Internet Service Model

What happened to “transactions” in this service model? What happened to network state? What happened to resource management within the network? What happened to all the elements of network-based communications services? The simple answer is that within the architecture of the Internet it is not necessary to expose such a detailed view of transactional state to the underlying network just to have the network deliver a packet. From a network perspective, IP has thrown all of that network level function away!

In the context of the Internet service architecture, a “transaction” is now merely an attribute of the application that is run on the end systems, and the underlying network is simply unaware of these transactions. All the network “sees” is IP packets, and each packet does not identify to the network any form of compound or multi-packet transaction.

Because a transaction is not directly visible to the IP network operator, the implication is that any effort for an IP service provider to use a transactional service tariff model becomes an exercise in frustration, given that there are no such network-visible interactions that could be used to create a transactional service model. In the absence of a network-based transactional service model, the *Internet Service Provider* (ISP) has typically used an access-based model as the basis of the IP tariff. Rather than paying a tariff per “call” the ISP typically charges a single flat fee independent of the number or nature of individual service transactions. Some basic differentiation is provided by the ability to apply price differentials to different access bandwidths or different volume caps, but this form of market segmentation is a relatively coarse one. Finer levels of transactional-based prices, such as pricing each individual video stream—or even pricing every individual webpage fetch—are not an inherent feature of such an access-based tariff structure.

The consequence for ISPs here is that within a single network access bandwidth class, this service model does not differentiate between heavy and light users, and is insensitive to the services operated across the network and to the average and peak loads imposed by these services. Like the flat-rate local telephone access model, the Internet pricing model is typically a flat-rate model that takes no account of individual network transactions. The ISP's service-delivery costs are, in effect, equally apportioned across the ISP's user base.

Interestingly, this feature has been a positive one for the Internet. With no marginal incremental costs for network usage, users are basically incented to use the Internet. In the same vein suppliers are also incented to use the Internet, because they can deliver goods and services to their customer base without imposing additional transaction costs to either themselves or their customers. For example, we have seen Microsoft and Apple move toward a software distribution model that is retiring the use of physical media, and moving to an all-digital Internet-based service model to support their user base. We have also seen other forms of service provision where the access-based tariff model has enabled services that would otherwise not be viable—here Netflix is a good example of such services that have been enabled by this flat-rate tariff structure. The attraction of cloud-based services in today's online world is another outcome of this form of incentive.

The other side effect of this shift in the architecture of the Internet is that it has placed the carriage provider—the network operator—into the role of a commodity utility. Without any ability to distinguish between various transactions, because the packets themselves give away little in terms of reliable information about the nature of the end-to-end service transaction, the carriage role is an undistinguished commodity utility function. The consequent set of competitive pressures in a market that is not strongly differentiated ultimately weans out all but the most efficient of providers from the service provider market—as long as competitive interests can be brought to bear on these market segments.

Invariably, consumers value the services that a network enables, rather than the network itself. In pushing the transaction out of the network and into the application, the architecture of the Internet also pushed value out of the network. Given that a service in the Internet model is an interaction between applications running on a content service provider's platform and on their clients' systems, it is clear that the network operator is not a direct party to the service transaction. An ISP may also provide services to users, but it is by no means an exclusive role, and others are also able to interact directly with customers and generate value through the provision of goods and services, without the involvement of the underlying network operators. It is not necessary to operate a network in order to offer a service on the Internet. Indeed, such a confusion of roles could well be a liability for such a carriage and content service provider.

The Content Business Model of the Internet

This unbundling of the service provision function from the network has had some rather unexpected outcomes. Those who made the initial forays of providing content to users believed that this function was no different from that of many retail models, where the content provider formed a set of relationships with a set of users. The direct translation of this model encountered numerous problems, not the least of which was reluctance on the part of individual users to enter into a panoply of service and content relationships. When coupled with considerations of control of secondary redistribution of the original service, this situation created some formidable barriers to the emergence of a highly valuable market for content and services on the Internet.

However, as with many forms of mass market media, the advertising market provides some strong motivation. With a traditional print newspaper, the full cost of the production of the newspaper is often borne largely by advertisers rather than by the newspaper readers. But newspaper advertising is a relatively crude exercise, in that the advertisement is visible to all readers, but it is of interest to a much smaller subset. The Internet provided the potential to customize the advertisement.

The greatest market value for advertisements is generated by those operations that gain the most information about their customers. These days it has a lot to do with knowledge of the consumer. It could be argued that Facebook's \$1B purchase of Instagram was based on the observation that the combination of an individual's pictures and updates forms an amazingly rich set of real-time information about the behavior and preferences of individual consumers. It could also be argued that Google's business model is similarly based on forming a comprehensive and accurate picture of individual users' preferences, which is then sold to advertisers at a significant premium simply because of its tailored accuracy. And the mobile services are trying to merge users' current locations with the knowledge of their preferences to gain even greater value.

These developments are heading in the direction of a multiparty service model, where the relationship between a content provider and a set of users allows the content provider to resell names of these users to third parties through advertising. This on-selling of users' profiles and preferences is now a very sophisticated and significant market. As reported in [1], some 90 percent of Google's \$37.9B income was derived from advertising revenue. The cost per click for "cheap car insurance" is reported in the same source to be \$33.97!

The Plight of the Carrier

Although the content market with its associated service plane is now an extraordinarily valuable activity, the same is not true for the network operator—whose carriage function has been reduced from complete service-delivery management to a simple packet carrier without any residual visibility into the service plane of the network.

Obviously, network carriers look at these developments with dismay. Their own traditional value-added market has been destroyed, and the former model where the telcos owned everything from the handset onward has now been replaced by a new model that relegates them to a role similar to electricity or water reticulation—with no prospect of adding unique value to the content and service market. The highly valuable service-level transactions are effectively invisible to the carriage service providers of the Internet.

There is an evident line of thought in the carriage industry that appears to say: “If we could capture the notion of a service-level transaction in IP we could recast our service profile into a per-transaction profile, and if we can do that, then we could have the opportunity to capture some proportion of the value of each transaction.”

Short of traffic interception, could the network operators working at the internet level of the network protocol stack have a means to identify these service-level transactions? The generic answer is “no,” as we have already seen, but there are some other possibilities that could expose service-level transactions to the network operator.

QoS to the Rescue?

The recent calls by the *The European Telecommunications Network Operators’ Association* (ETNO) advocating the widespread adoption of IP *Quality of Service* (QoS) appear to have some context from this perspective of restoring transaction visibility to the IP carriage provider. In the QoS model an application undertakes a QoS “reservation” with the network. The network is supposed to respond with a commitment to reserve the necessary resources for use by this transaction. The application then uses this QoS channel for its transaction, and releases the reservation when the transaction is complete.

From the network operator’s perspective, the QoS-enabled network is now being informed of individual transactions, identifying the end parties for the transaction, the nature of the transaction and its duration, as well as the resource consumption associated with the transaction. From this information comes the possibility for the QoS IP network operator to move away from a now commonplace one-sided flat access tariff structure for IP services, and instead use a transactional service model that enables the network operator to impose transaction-based service fees on both parties to a network service if it so chooses. It also interposes the network operator between the content provider and the consumer, permitting the network operator to mediate the content service and potentially convert this gateway role into a revenue stream.

Of course the major problem in this QoS model is that it is based on a critical item of Internet mythology—the myth that inter-provider QoS exists on the Internet. QoS is not part of today’s Internet, and there is no visible prospect that it will be part of tomorrow’s Internet either!

Knotting up NATs

But QoS is not the only possible approach to exposing service-level transactions to the carriage-level IP network operator. Interestingly, the twin factors of the exhaustion of IPv4 addresses and the lack of uptake of IPv6 offers the IP network operator another window into what the user is doing, and, potentially, another means of controlling the quality of the user's experience by isolating individual user-level transactions at the network level.

When there are not enough addresses to assign each customer a unique IP address, the ISP is forced to use private addresses and operate a *Network Address Translator* (NAT)^[2] within the carriage network.

However, NATs are not stateless passive devices. A NAT records every TCP and *User Datagram Protocol* (UDP) session from the user, as well as the port addresses the application uses when it creates a binding from an internal IP address and port to an external IP address and port. A new NAT binding is created for every user transaction: every conversation, every website, every streamed video, and literally everything else. If you were to look at the NAT logs that record this binding information, you would find a rich stream of real-time user data that shows precisely what each user is doing on the network. Every service transaction is now visible at the network level. How big is the temptation for the IP network operator to peek at this carrier-operated NAT log and analyze what it means?

Potentially, this transaction data could be monetized, because it forms a real-time data feed of every customer's use of the network. At the moment carriers think that they are being compelled to purchase and install this NAT function because of the IPv4 address situation. NATs offer a method for the carriage operator to obtain real-time feeds of customer behavior without actively intruding themselves into the packet stream. The NAT neatly segments the customer's traffic into distinct transactions that are directly visible to the NAT operator. I suspect that when they look at the business case for purchasing and deploying these *Carrier-Grade NAT* devices, they will notice a parallel business case that can be made to inspect the NAT logs and perhaps to either on-sell the data stream or analyze it themselves to learn about their customers' behavior.^[3] And, as noted, there is already market evidence that such detailed real-time flows of information about individual users' activities can be worth significant sums.

But it need not necessarily be limited to a passive operation of stalking the user's online behavior. If the carriage provider were adventurous enough, it could bias the NAT port-binding function to even make some content work "better" than other content, by either slowing down the binding function for certain external sites or rationing available ports to certain less-preferred external sites. In effect, NATs provide many exploitable levers of control for the carriage operator, bundled with a convenient excuse of "we had no choice but to deploy these NATs!"

Where Now?

In contrast, what does an investment in IPv6 offer the carriage provider? An admittedly very bleak response from the limited perspective of the carriage service provider sector is that what is on offer with IPv6 is more of what has happened to the telecommunications carriage sector over the past 10 years, with not even the remote possibility of ever altering this situation. IPv6 certainly looks like forever, so if the carriers head down this path then the future looks awfully bleak for those who are entirely unused to, and uncomfortable with, a commodity utility provider role.

So should we just throw up our hands at this juncture and allow the carriage providers free rein? Are NATs inevitable? Should we view the introduction of transactional service models in the Internet as a necessary part of its evolution? I would like to think that these developments are not inevitable for the Internet, and that there are other paths that could be followed here. The true value for the end consumer is not in the carriage of bits through the network, but in the access to communication and services that such bit carriage enables. What does that reality imply for the future for the carriage role? I suspect that despite some evident misgivings, the carriage role is inexorably heading to that of a commodity utility operation.

This is not the first time an industry sector has transitioned from production of a small volume of highly valuable units to production of a massively larger volume of commodity goods, each of which has a far lower unit value, but generates an aggregate total that is much larger. The computing industry's transition from mainframe computers to mass market consumer electronics is a good example of such a transformation. As many IT sector enterprises have shown, it is possible to make such transitions. IBM is perhaps a classic example of an enterprise that has managed numerous successful transformations that have enabled it to maintain relevance and value in a rapidly changing environment.

The models for electricity distribution have seen a similar form of evolution in the last century. In the 1920s in the United Kingdom, electricity was a low-volume premium product. The prices for electricity were such that to keep just 5 light bulbs running for 1 day in a household cost the equivalent of an average week's wages. The consequent years saw public intervention in the form of nationalization of power generation and distribution that transformed electricity supply into a commonly available and generally affordable commodity.

The challenge the Internet has posed for the carriage sector is not all that different from these examples. The old carriage business models of relatively low-volume, high-value, transaction-based telecommunication services of telephony and faxes find no resonance within the service model of the Internet.

In the architecture of the Internet, it is the applications that define the services, while the demands from the underlying carriage network have been reduced to a simple stateless datagram-delivery service. Necessarily, the business models of carriage have to also change to adapt to this altered role, and one of the more fundamental changes is the dropping of the transaction-based model of the provision of telecommunications services for the carriage provider. What this situation implies for the carriage sector of the Internet is perhaps as radical as the transformation of the electricity supply industry during the period of the construction of the national grid systems in the first half of the 20th century.

The necessary change implied here is from a high-value premium service provider dealing in individual transactions across the network to that of a high-volume undistinguished commodity utility operator. The architectural concepts of a minimal undistinguished network carriage role and the repositioning of service management into end-to-end applications is an intrinsic part of the architecture of the Internet itself. It is not a universally acclaimed step—and certainly not one that is particularly popular in today’s carriage industry—but if we want to see long-term benefits from the use of the Internet in terms of positive economic outcomes and efficient exploitation of this technology in delivering goods and services, then it is a necessary step in the broader long-term public interest.

References

- [0] Google’s IPv6 statistics:
<http://www.google.com/ipv6/statistics.html>
- [1] Connor Livingston, “A breakdown of Google’s top advertisers,”
<http://www.techi.com/2012/03/a-breakdown-of-googles-top-advertisers/>
- [2] Geoff Huston, “Anatomy: A Look inside Network Address Translators,” *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [3] Geoff Huston, “All Your Packets Belong to Us,” July 2012,
<http://www.potaroo.net/ispcol/2012-07/allyourpackets.html>

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001.

E-mail: gih@apnic.net

The Internet: Looking Forward

by Vint Cerf, Google

As I write, it is 2013 and 40 years have passed since the first drafts of the Internet design were written. The first published paper appeared in 1974^[1] and the first implementations began in 1975. Much has happened since that time, but this essay is not focused on the past but, rather, on the future. Although the past is plainly prologue, our ability to see ahead is hampered by the unpredictable and the unknown unknowns that cloud and bedevil our vision. The exercise is nonetheless worth the effort, if only to imagine what might be possible.

Current trends reveal some directions. Mobile devices are accelerating access and applications. The economics of mobile devices have increased the footprint of affordable access to the Internet and the World Wide Web. Mobile infrastructure continues to expand on all inhabited continents. Speeds and functions are increasing as faster processors, more memory, and improved display technologies enhance the functions of these platforms. Cameras, microphones, speakers, sensors, multiple radios, touch-sensitive displays, and location and motion detection continue to evolve and open up new application possibilities. Standards and open source software facilitate widespread interoperability and adoption of applications. What is perhaps most significant is that these smart devices derive much of their power from access to and use of the extraordinary computing and memory capacity of the Internet. The Internet, cloud computing, and mobile devices have become hypergolic in their capacity to ignite new businesses and create new economic opportunities.

In the near term, the Internet is evolving. The *Domain Name System* (DNS) is expanding dramatically at the top level. Domain names can be written in non-Latin characters. The Internet address space is being expanded through the introduction of the IPv6 packet format, although the implementation rate among *Internet Service Providers* (ISPs) continues to be unsatisfactorily slow. This latter phenomenon may change as the so-called *Internet of Things*^[2] emerges from its long incubation. Sensor networks, Internet-enabled appliances, and increasing application of artificial intelligence will transform the Internet landscape in ways that seem impossible to imagine. The introduction of IPv6 and the exhaustion of the older IPv4 address space have generated demand for application of the so-called *Network Address Translation* (NAT)^[3] system. Geoff Huston has written and lectured extensively on this topic^[4] and the potential futures involving their use. In some ways, these systems simultaneously interfere with the motivation to implement IPv6 and act as a bridge to allow both network address formats to be used concurrently.

Ironically, although most edge devices on the Internet today are probably IPv6-capable, as are the routers, firewalls, DNS servers, and other application servers, this advanced version of the Internet Protocol may not have been “turned on” by the ISP community. This situation is changing, but more slowly than many of us would like.

As the applications on the Internet continue to make demands on its capacity to transport data and to deliver low-latency services, conventional Internet technologies are challenged and new ideas are finding purchase in the infrastructure. The *OpenFlow*^[5, 6] concept has emerged as a fresh look at packet switching in which control flow is segregated from data flow and routing is not confined to the use of address bits in packet headers for the formation and use of forwarding tables. Originally implemented with a central routing scheme to improve efficient use of network resources, the system has the flexibility to be made more distributed. It remains to be seen whether OpenFlow networks can be interconnected by using an extended form of the *Border Gateway Protocol* (BGP) so as to achieve end-to-end performance comparable to what has already been achieved in single networks.

Business models for Internet service play an important role here because end-to-end differential classes of service have not been realized, generally, for the current Internet implementations. Inter-ISP or edge-to-core commercial models also have not generally been perfected to achieve multiple classes of service. These aspirations remain for the Internet of the present day. Although it might be argued that increasing capacity in the core and at the edge of the Internet eliminates the need for differential service, it is fair to say that some applications definitely need lower delay, others need high capacity, and some need both (for example, for interactive video). Whether these requirements can be met simply through higher speeds or whether differential services must be realized at the edges and the core of the network is the source of substantial debate in the community. Vigorous experimentation and research continue to explore these topics.

Ubiquitous Computing

Mark Weiser^[7] coined the term and concept of *Ubiquitous Computing*. He meant several things by this term, but among them was the notion that computers would eventually fade into the environment, becoming ever-present, performing useful functions, and operating for our convenience. Many devices would host computing capacity but would not be viewed as “computers” or even “computing platforms.” Entertainment devices; cooking appliances; automobiles; medical, environmental, and security monitoring systems; our clothing; and our homes and offices would house many computing engines of various sizes and capacities. Many, if not all, would be interconnected in communication webs, responding to requirements and policies set by users or by their authorized representatives.

To this idyllic characterization, he implied there would be challenges: configurations of hundreds of thousands of appliances and platforms, privacy, safety, access control, information confidentiality, stability, resilience, and a host of other properties.

Even modest thought produces an awareness of the need for strong authentication to assure that only the appropriate devices and authorized parties are interacting, issuing instructions, taking data, etc. It is clear that multifactor authentication and some form of public key cryptography could play an important role in assuring limitations on the use and operation of these systems. Privacy of the information generated by these systems can be understood to be necessary to protect users from potential harm.

The scale of such systems can easily reach tens to hundreds of billions of devices. Managing complex interactions at such magnitudes will require powerful hierarchical and abstracting mechanisms. When it is also understood that our mobile society will lead to a constant background churn of combinations of devices forming subsets in homes, offices, automobiles, and on our persons, the challenge becomes all the more daunting. (By this I do not mean the use of mobile smartphones but rather a society that is geographically mobile and that moves some but not all its possessions from place to place, mixing them with new ones.) Self-organizing mechanisms, hierarchically structured systems, and systems that allow remote management and reporting will play a role in managing the rapidly proliferating network we call the Internet.

For further insight into this evolution, we should consider the position location capability of the *Global Positioning System* (GPS)^[8]. Even small, low-powered devices (for example, mobile devices) have the ability to locate themselves if they have access to the proper satellite transmissions. Adding to this capability is geo-location using mobile cell towers and even known public Wi-Fi locations. In addition, we are starting to see appliances such as *Google Glass*^[9] enter the environment. These appliances are portable, wearable computers that hear what we hear and see what we see and can respond to spoken commands and gestures. The Google self-driving cars^[10] offer yet another glimpse into the future of computing, communication, and artificial intelligence in which computers become our partners in a common sensory environment—one that is not limited to the normal human senses. All of these systems have the potential to draw upon networked information and computing power that rivals anything available in history. The systems are potentially self-learning and thus capable of improvement over time. Moreover, because these devices may be able to communicate among themselves, they may be able to cooperate on a scale never before possible.

Even now we can see the outlines of a potential future in which virtually all knowledge can be found for the asking; in which the applications of the Internet continue to evolve; in which devices and appliances of all kinds respond and adapt to our needs, communicate with each other, learn from each other, and become part of an integrated and global environment.

Indeed, our day-to-day environment is very likely to be filled with information and data gathered from many sources and subject to deep analysis benefitting individuals, businesses, families, and governments at all levels. Public health and safety are sure to be influenced and affected by these trends.

Education

It is often noted that a teacher from the mid-19th century would not feel out of place in the classroom of the 21st, except, perhaps, for subject matter. There is every indication that this situation may be about to change. In 2012, two of my colleagues from Google, Peter Norvig and Sebastian Thrun, decided to use the Internet to teach an online class in artificial intelligence under the auspices of Stanford University. They expected about 500 students, but 160,000 people signed up for the course! There ensued a scramble to write or revise software to cope with the unexpectedly large scale of the online class. This phenomenon has been a long time in coming. Today we call such classes “MOOCs” (*Massive, Open, OnLine Classes*). Of the 160,000 who signed up, something like 23,000 actually completed the class. How many professors of computer science can say they have successfully taught 23,000 students?

The economics of this form of classroom are also very intriguing. Imagine a class of 100,000 students, each paying \$10 per class. Even one class would produce \$1,000,000 in revenue. I cannot think of any university that regularly has million dollar classes! There are costs, but they are borne in part by students (Internet access, equipment with which to reach the Internet, etc., for example) and in part by the university (Internet access, multicast or similar capability, and salaries of professors and teaching assistants). In some cases, the professors prepare online lectures that students can watch as many times as they want to—whenever they want to because the lectures can be streamed. The professors then hold classroom hours that are devoted to solving problems, in an inversion of the more typical classroom usage. Obviously this idea could expand to include nonlocal teaching assistants. Indeed, earlier experiments with video-taped lectures and remote teaching assistants were carried out with some success at Stanford University when I served on the faculty in the early 1970s.

What is potentially different about MOOCs is *scale*. Interaction and examinations are feasible in this online environment, although the form of exams is somewhat limited by the capabilities of the online platform used. Start-ups are experimenting with and pursuing these ideas (refer to www.udacity.com and www.coursera.org).

People who are currently employed also can take these courses to improve their skills, learn new ones, and position themselves for new careers or career paths. From young students to retired workers, such courses offer opportunities for personal expansion, and they provide a much larger customer base than is usually associated with a 2- or 4-year university or college program. These classes can be seen as re-invention of the university, the short course, the certificate program, and other forms of educational practice. It is my sense that this state of affairs has the potential to change the face of education at all levels and provide new options for those who want or need to learn new things.

The Information Universe

It is becoming common to speak of “big data” and “cloud computing” as indicators of a paradigm shift in our view of information. This view is not unwarranted. We have the ability to absorb, process, and analyze quantities of data beyond anything remotely possible in the past. The functional possibilities are almost impossible to fully fathom. For example, our ability to translate text and spoken language is unprecedented. With combinations of statistical methods, hierarchical hidden Markov models, formal grammars, and Bayesian techniques, the fidelity of translation between some language pairs approaches native language speaker quality. It is readily predictable that during the next decade, real-time, spoken language translation will be a reality.

One of my favorite scenarios: A blind German speaker and a deaf *American Sign Language* (ASL) signer meet, each wearing Google Glass. The deaf signer’s microphone picks up the German speaker’s words, translates them into English, and displays them as captions for the deaf participant. The blind man’s Glass video camera sees the deaf signer’s signs, translates the signs from ASL to English and then to German, and then speaks them through the bone conduction speaker of the Google Glass. We can do all of this now except for the correct interpretation of ASL. This challenge is not a trivial one, but it might be possible in the next 10 to 15 years.

The World Wide Web continues to grow in size and diversity. In addition, large databases of information are being accumulated, especially from scientific disciplines such as physics, astronomy, and biology. Telescopes (ground and space-based), particle colliders such as the Large Hadron Collider^[11], and DNA sequencers are producing petabytes and more—in some cases on a daily basis!

We seem to be entering a time when much of the information produced by human endeavor will be accessible to everyone on the planet. Google’s motto: “To organize the world’s information and make it universally accessible and useful,” might be nearly fulfilled in the decades ahead. Some tough problems lie ahead, however. One I call “bit rot.”

By using this term, I do not mean the degradation of digital recordings on various media, although this is a very real problem. The more typical problem is that the readers of the media fall into disuse and disrepair. One has only to think about 8-inch Wang disks for the early Wang word processor, or 3.5-inch floppy disks or their 5 ¼-inch predecessors. Now we have CDs, DVDs, and Blu-Ray disks, but some computer makers—Apple for example—have ceased to build in readers for these media.

Another, more tricky problem is that much of the digital information produced requires software to correctly interpret the digital bits. If the software is not available to interpret the bits, the bits might as well be rotten or unreadable. Software applications run over operating systems that, themselves, run on computer hardware. If the applications do not work on new versions of the operating systems, or the applications are upgraded but are not backward-compatible with earlier file and storage formats, or the maker of the application software goes out of business and the source code is lost, then the ability to interpret the files created by this software may be lost. Even when open source software is used, it is not clear it will be maintained in operating condition for thousands of years. We already see backward-compatibility failures in proprietary software emerging after only years or decades.

Getting access to source code for preservation may involve revising notions of copyright or patent to allow archivists to save and make usable older application software. We can imagine that “cloud computing” might allow us to emulate hardware, run older operating systems, and thus support older applications, but there is also the problem of basic input/output and the ability to emulate earlier media, even if the physical media or their readers are no longer available. This challenge is a huge but important one.

Archiving of important physical data has to be accompanied by archiving of metadata describing the conditions of collections, calibration of instruments, formatting of the data, and other hints at how to interpret it. All of this work is extra, but necessary to make information longevity a reality.

The Dark Side

To the generally optimistic and positive picture of Internet service must be added a realistic view of its darker side. The online environment and the devices we use to exercise it are filled with software. It is an unfortunate fact that programmers have not succeeded in discovering how to write software of any complexity that is free of mistakes and vulnerabilities.

Despite the truly remarkable and positive benefits already delivered to us through the Internet, we must cope with the fact that the Internet is not always a safe place.

The software upon which we rely in our access devices, in the application servers, and in the devices that realize the Internet itself (routers, firewalls, gateways, switches, etc.) is a major vulnerability, given the apparently inescapable presence of bugs.

Not everyone with access to the Internet has other users' best interests at heart. Some see the increasing dependence of our societies on the Internet as an opportunity for exploitation and harm. Some are motivated by a desire to benefit themselves at the expense of others, some by a desire to hurt others, some by nationalistic sentiments, some by international politics. That Shakespeare's plays are still popular after 500 years suggests that human frailties have not changed in the past half millennium! The weaknesses and vulnerabilities of the Internet software environment are exploited regularly. What might the future hold in terms of making the Internet a safer and more secure place in which to operate?

It is clear that simple usernames and passwords are inadequate to the task of protecting against unauthorized access and that multi-factor and perhaps also biometric means are going to be needed to accomplish the desired effect. We may anticipate that such features might become a part of reaching adulthood or perhaps a rite of passage at an earlier age. Purely software attempts to cope with confidentiality, privacy, access control, and the like will give way to hardware-reinforced security. Digitally signed *Basic Input/Output System* (BIOS), for example, is already a feature of some new chipsets. Some form of trusted computing platform will be needed as the future unfolds and as online and offline hazards proliferate.

Governments are formed that are, in principle, kinds of social contracts. Citizens give up some freedoms in exchange for safety from harm. Not all regimes have their citizens' best interests at heart, of course. There are authoritarian regimes whose primary interest is staying in power. Setting these examples aside, however, it is becoming clear that the hazards of using computers and being online have come to the attention of democratic as well as authoritarian regimes. There is tension between law enforcement (and even determination of what the law should be) and the desire of citizens for privacy and freedom of action. Balancing these tensions is a nontrivial exercise. The private sector is pressed into becoming an enforcer of the law when this role is not necessarily an appropriate one. The private sector is also coerced into breaching privacy in the name of the law.

"Internet Governance" is a broad term that is frequently interpreted in various ways depending on the interest of the party desiring to define it for particular purposes. In a general sense, Internet Governance has to do with the policies, procedures, and conventions adopted domestically and internationally for the use of the Internet. It has not only to do with the technical ways in which the Internet is operated, implemented, and evolved but also with the ways in which it is used or abused.

In some cases it has to do with the content of the Internet and the applications to which the Internet is put. It is evident that abuse is undertaken through the Internet. Fraud, stalking, misinformation, incitement, theft, operational interference, and a host of other abuses have been identified. Efforts to defend against them are often stymied by lack of jurisdiction, particularly in cases where international borders are involved. Ultimately, we will have to reach some conclusions domestically and internationally as to which behaviors will be tolerated and which will not, and what the consequences of abusive behavior will be. We will continue to debate these problems well into the future.

Our societies have evolved various mechanisms for protecting citizens. One of these mechanisms is the Fire Department. Sometimes volunteer, this institution is intended to put out building or forest fires to minimize risks to the population. We do not have a similar institution for dealing with various forms of “cyberfires” in which our machines are under attack or are otherwise malfunctioning, risking others by propagation of viruses, worms, and Trojan horses or participation in botnet denial-of-service or other forms of attacks. Although some of these matters may deserve national-level responses, many are really local problems that would benefit from a “Cyber Fire Department” that individuals and businesses could call upon for assistance. When the cyber fire is put out, the question of cause and origin could be investigated as is done with real fires. If deliberately set, the problem would become one of law enforcement.

Intellectual property is a concept that has evolved over time but is often protected by copyright or patent practices that may be internationally adopted and accepted. These notions, especially copyright, had origins in the physical reproduction of content in the form of books, films, photographs, CDs, and other physical things containing content. As the digital and online environment penetrates more deeply into all societies, these concepts become more and more difficult to enforce. Reproduction and distribution of digital content gets easier and less expensive every day. It may be that new models of compensation and access control will be needed in decades ahead.

Conclusion

If there can be any conclusion to these ramblings, it must be that the world that lies ahead will be immersed in information that admits of extremely deep analysis and management. Artificial intelligence methods will permeate the environment, aiding us with smart digital assistants that empower our thought and our ability to absorb, understand, and gain insight from massive amounts of information.

It will be a world that is also at risk for lack of security, safety, and privacy—a world in which demands will be made of us to think more deeply about what we see, hear, and learn. While we have new tools with which to think, it will be demanded of us that we use them to distinguish sound information from unsound, propaganda from truth, and wisdom from folly.

References

- [1] Vinton G. Cerf and Robert E. Kahn, “A Protocol for Packet Network Intercommunication,” *IEEE Transactions on Communications*, Vol. Com-22, No. 5, May 1974.
- [2] David Lake, Ammar Rayes, and Monique Morrow, “The Internet of Things,” *The Internet Protocol Journal*, Volume 15, No. 3, September 2012.
- [3] Geoff Huston, “Anatomy: A Look inside Network Address Translators,” *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [4] Geoff Huston and Mark Koster, “The Role of Carrier Grade NATs in the Near-Term Internet,” TIP 2013 Conference, <http://events.internet2.edu/2013/tip/agenda.cfm?go=session&id=10002780>
- [5] <http://www.openflow.org/>
- [6] William Stallings, “Software-Defined Networks and OpenFlow,” *The Internet Protocol Journal*, Volume 16, No. 1, March 2013.
- [7] http://en.wikipedia.org/wiki/Mark_Weiser
- [8] http://en.wikipedia.org/wiki/Global_Positioning_System
- [9] <http://www.google.com/glass/start/>
- [10] http://en.wikipedia.org/wiki/Google_driverless_car
- [11] home.web.cern.ch
- [12] Cerf, V., “Looking Toward the Future,” *The Internet Protocol Journal*, Volume 10, No. 4, December 2007.
- [13] Vint Cerf, “A Decade of Internet Evolution,” *The Internet Protocol Journal*, Volume 11, No. 2, June 2008.
- [14] Geoff Huston, “A Decade in the Life of the Internet,” *The Internet Protocol Journal*, Volume 11, No. 2, June 2008.

VINTON G. CERF is vice president and chief Internet evangelist for Google. Cerf has held positions at MCI, the Corporation for National Research Initiatives, Stanford University, UCLA, and IBM. He served as chairman of the board of the Internet Corporation for Assigned Names and Numbers (ICANN) and was founding president of the Internet Society. Cerf was appointed to the U.S. National Science Board in 2013. Widely known as one of the “Fathers of the Internet,” he received the U.S. National Medal of Technology in 1997, the Marconi Fellowship in 1998, and the ACM Alan M. Turing Award in 2004. In November 2005, he was awarded the Presidential Medal of Freedom, in April 2008 the Japan Prize, and in March 2013 the Queen Elizabeth II Prize for Engineering. He is a Fellow of the IEEE, ACM, and AAAS, the American Academy of Arts and Sciences, the American Philosophical Society, the Computer History Museum, and the National Academy of Engineering. Cerf holds a Bachelor of Science degree in Mathematics from Stanford University and Master of Science and Ph.D. degrees in Computer Science from UCLA, and he holds 21 honorary degrees from universities around the world.
E-mail: vint@google.com

Optimizing Link-State Protocols for Data Center Networks

by Alvaro Retana, Cisco Systems, and Russ White, Verisign

With the advent of cloud computing^[6, 7], the pendulum has swung from focusing on wide-area or global network design toward a focus on *Data Center* network design. Many of the lessons we have learned in the global design space will be relearned in the data center space before the pendulum returns and wide-area design comes back to the fore.

This article examines three extensions to the *Open Shortest Path First* (OSPF) protocol that did not originate in the data center field but have direct applicability to efficient and scalable network operation in highly meshed environments. Specifically, the application extensions to OSPF to reduce flooding in *Mobile Ad Hoc Networks* (MANET)^[1], demand circuits designed to support on-demand links in wide-area networks^[2], and OSPF stub router advertisements designed to support large-scale *hub and spoke* networks^[3] are considered in a typical data center network design to show how these sorts of protocol improvements could affect the scaling of data center environments.

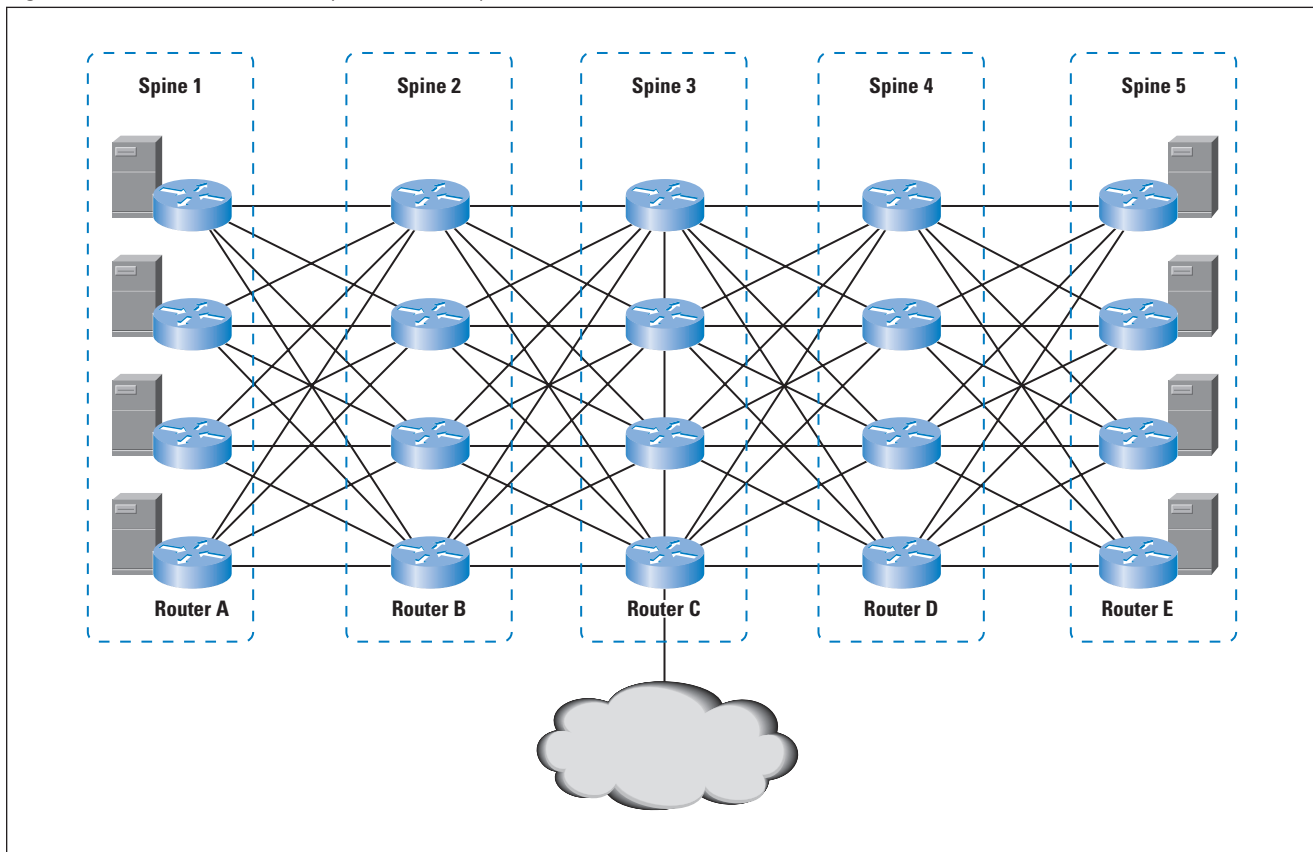
Each of the improvements examined has the advantage of being available in shipping code from at least one major vendor. All of them have been deployed and tested in real-world networks, and have proven effective for solving the problems they were originally designed to address. Note, as well, that OSPF is used throughout this article, but each of these improvements is also applicable to *Intermediate System-to-Intermediate System* (IS-IS), or any other link-state protocol.

Defining the Problem

Figure 1 illustrates a small Clos^[0] fabric, what might be a piece of a much larger network design. Although full-mesh fabrics have fallen out of favor with data center designers, Clos and other styles of fabrics are in widespread use. A Clos fabric configured with edge-to-edge Layer 3 routing has three easily identifiable problems.

The flooding rate is the first problem a link-state protocol used in this configuration must deal with. Router B (and the other routers in spine 2), for instance, will receive four type 1 *Link State Advertisements* (LSAs) from the four routers in spine 1. Each of the routers in spine 2 will reflood each of these type 1 LSAs into spine 3, so the other routers in spines 3, 4, and 5 will each receive four copies of each type 1 LSA originated by routers in spine 1, a total of 16 type 1 LSAs in all.

Figure 1: A Clos Fabric with Layer 3 to the Top of Rack



To make matters worse, OSPF is designed to time out every LSA originated in the network once every 20 to 30 minutes. This feature was originally put in OSPF to provide for recovery from bit and other transmission errors in older transport mechanisms with little or no error correction. So a router in spine 5 will receive 16 copies of each type 1 LSA generated by routers in spine 1 every 20 minutes. A single link failure and recovery can also cause massive reflooding. The process of bringing the OSPF adjacency back into full operation requires a complete exchange of local link-state databases. If the link between router A and router B fails and then is recovered, the entire database must be transferred between the two routers, even though router B clearly has a complete copy of the database from other sources.

Finally, the design of this network produces some challenges for the *Shortest Path First* (SPF) algorithm, which link-state protocols use to determine the best path to each reachable destination in the network. Every router in spine 1 appears to be a transit path to every other destination in the network. This outcome might not be the intent of the network designer, but SPF calculations deal with available paths, not intent.

This set of problems has typically swayed network designers away from using link-state protocols in such large-scale environments. Some large cloud service providers use the *Border Gateway Protocol* (BGP) (see [4]), with each spine being a separate Autonomous System, so they can provide scalable Layer 3 connectivity edge-to-edge in large Clos network topologies. Others have opted for simple controls, such as removing all control-plane protocols and relying on reverse-path-forwarding filters to prevent loops.

The modifications to OSPF discussed in this article, however, make it possible for a link-state protocol to not only scale in this type of environment, but also to be a better choice.

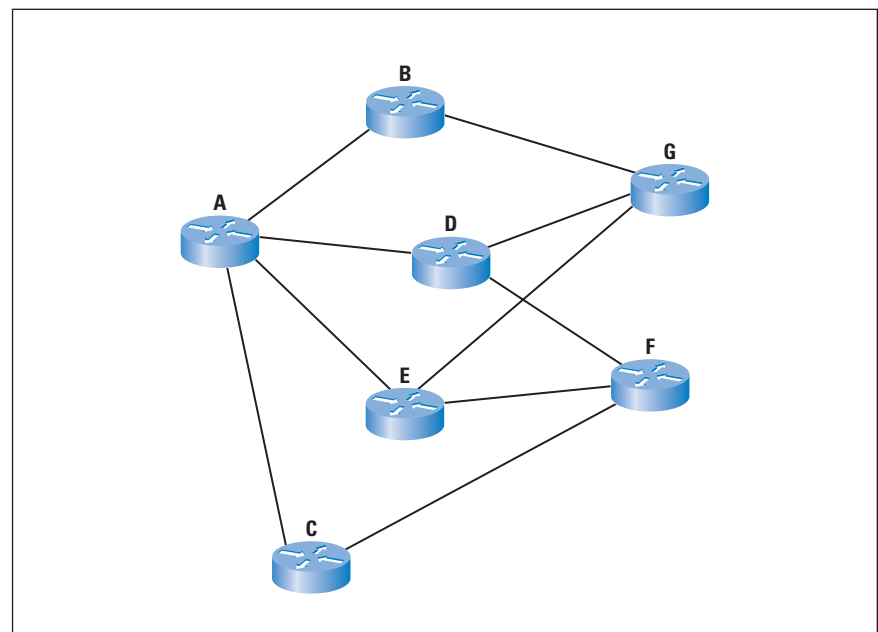
Reducing Flooding Through MANET Extensions

MANET networks are designed to be “throw and forget;” a collection of devices is deployed into a quickly fluid situation on the ground, where they connect over short- and long-haul wireless links, and “just work.” One of the primary scaling (and operational) factors in these environments is an absolute reduction of link usage wherever possible, including for the control plane.

The “Extensions to OSPF to Support Mobile Ad Hoc Networking,”^[1] were developed to reduce flooding in single-area OSPF networks to the minimal necessary, while providing fast recovery and guaranteed delivery of control-plane information. The idea revolves around the concept of an overlapping relay, which reduces flooding by accounting for the network topology, specifically groups of overlapping nodes.

Let’s examine the process from the perspective of router A shown in Figure 2.

Figure 2: Ad Hoc Extensions to OSPF



Router A begins the process by not only discovering that it is connected to routers B, C, D, and E, but also that its *two-hop neighborhood* contains routers F and G. By examining the list of two-hop neighbors, and the directly connected neighbors that can reach each of those two-hop neighbors, router A can determine that if router D refloods any LSAs router A floods, every router in the network will receive the changes. Given this information, router A notifies routers B, C, and E to delay the reflooding of any LSAs received from router A itself.

When router A floods an LSA, router D will reflood the LSA to routers F and G, which will then acknowledge receiving the LSA to routers B, C, D, and E. On receiving this acknowledgement, routers B, C, and E will remove the changed LSA from their reflood lists.

Routers F and G, then, will receive only one copy of the changed LSA, rather than four.

Applying this process to the Clos design in Figure 1 and using this extension would dramatically reduce the number of LSAs flooded through the network in the case of a topology change. If router A, for instance, flooded a new type 1 LSA, the routers in spine 2 would each receive one copy. The routers in spines 3, 4, and 5 would also receive only one copy each, rather than 4 or 16.

Reducing Flooding Through Demand Circuits

Network engineers have long had to consider links that are connected only when traffic is flowing in their network and protocol designs. Dial-up links, for instance, or dynamically configured *IP Security* (IPsec) tunnels, have always been a part of the networking landscape. Part of the problem with such links is that the network needs to draw traffic to destinations reachable through the link even though the link is not currently operational.

With protocols that rely on neighbor adjacencies to maintain database freshness, such as OSPF, links that can be disconnected in the control plane and yet still remain valid in the data plane pose a unique set of difficulties. The link must appear to be available in the network topology even when it is, in fact, not available.

To overcome this challenge, the OSPF working group in the IETF extended the protocol to support demand links. Rather than attacking the problem at the adjacency level, OSPF attacks the problem at the database level. Any LSA learned over a link configured as a demand link is marked with the *Do Not Age* (DNA) bit; such LSAs are exempt from the normal aging process, causing LSAs to be removed from the link-state database periodically.

How does this situation relate to scaling OSPF in data center network design?

Every 20 minutes or so, an OSPF implementation will time out all the locally generated LSAs, replacing them with newly generated (and identical) LSAs. These newly generated LSAs will be flooded throughout the network, replacing the timed-out copy of the LSA throughout the network. In a data center network, these refloods are simply redundant; there is no reason to refresh the entire link-state database periodically.

To reduce flooding, then, data center network designers can configure all the links in the data center as demand circuits. Although these links are, in reality, always available, configuring them as demand circuits causes the DNA bit to be set on all the LSAs generated in the network. This process, in turn, disables periodic reflooding of this information, reducing control-plane overhead.

Reducing Control-Plane Overhead by Incremental Database Synchronization

When a link fails and then recovers, the OSPF protocol specifies a lengthy procedure through which the two newly adjacent OSPF processes must pass to ensure their databases are exactly synchronized. In the case of data center networks, however, there is little likelihood that a single link failure (or even multiple link failures) will cause two adjacent OSPF processes to have desynchronized databases.

For instance, in Figure 1, if the link between routers A and B fails, routers A and B will still receive any and all link-state database updates from some other neighbor they are still fully adjacent with. When the link between routers A and B is restored, there is little reason for routers A and B to exchange their entire databases again.

This situation is addressed through another extension suggested through the MANET extensions to OSPF called *Unsynchronized Adjacencies*. Rather than sending an entire copy of the database on restart and waiting until this exchange is complete to begin forwarding traffic on link recovery, this extension states that OSPF processes do not need to synchronize their databases if they are already synchronized with other nodes in the network. If needed, the adjacency can be synchronized out of band at a later time.

The application of the MANET OSPF extensions^[1] to a data center network means links can be pressed into service very quickly on recovery, and it provides a reduction in the amount of control-plane traffic required for OSPF to recover.

Reducing Processing Overhead Through Stub Routers

The SPF calculation that link-state protocols use to determine the best path to any given destination in the network treats all nodes and all edges on the graph as equal. Returning to Figure 2, router B will calculate a path through router A to routers D, E, and C, even if router A is not designed to be a transit node in the network. This failure to differentiate between transit and nontransit nodes in the network graph increases the number of paths SPF must explore when calculating the shortest-path tree to all reachable destinations.

Although modern implementations of SPF do not suffer from problems with calculation overhead or processor usage, in large-scale environments, such as a data center network with tens of thousands of nodes in the shortest-path tree and virtualization requirements that cause a single node to run SPF hundreds or thousands of times, small savings in processing power can add up.

The “OSPF Stub Router Advertisement”^[3] mechanism allows network administrators to mark an OSPF router as nontransit in the shortest-path tree. This mechanism would, for instance, prevent router A in Figure 1 from being considered a transit path between router B and some other router in spine 2. You would normally want to consider this option only for any actual edge routers in the network, such as the top-of-rack routers shown here. Preventing these routers from being used for transit can reduce the amount of redundancy available in the network, and, if used anywhere other than a true edge, prevent the network from fully forming a shortest-path tree.

Advantages and Disadvantages of Link-State Protocols in the Data Center

Beyond the obvious concerns of convergence speed and simplicity, there is one other advantage to using a link-state protocol in data center designs: equal-cost load sharing. OSPF and IS-IS both load share across all available equal-cost links automatically (subject to the limitations of the forwarding table in any given implementation). No complex extensions (such as [5]), are required to enable load sharing across multiple paths.

One potential downside to using a link-state protocol in a data center environment must be mentioned, however—although BGP allows route filtering at any point in the network (because it is a path vector-based protocol)—link-state protocols can filter or aggregate reachability information only at flooding domain boundaries. This limitation makes it more difficult to manage traffic flows through a data center network using OSPF or IS-IS to advertise routing information. This problem has possible solutions, but this area is one of future, rather than current, work.

Conclusion

Many improvements have been made to link-state protocols over the years to improve their performance in specific situations, such as MANETs, and when interacting with dynamically created links or circuits. Many of these improvements are already deployed and tested in real network environments, so using them in a data center environment is a matter of application rather than new work. All of these improvements are applicable to link-state control planes used for Layer 2 forwarding, as well as Layer 3 forwarding, and they are applicable to OSPF and IS-IS.

These improvements, when properly applied, can make link-state protocols a viable choice for use in large-scale, strongly meshed data center networks.

References

- [0] http://en.wikipedia.org/wiki/Clos_network
- [1] Roy, A., “Extensions to OSPF to Support Mobile Ad Hoc Networking,” RFC 5820, March 2010.
- [2] Abhay Roy and Sira Panduranga Rao, “Detecting Inactive Neighbors over OSPF Demand Circuits (DC),” RFC 3883, October 2004.
- [3] Alvaro Retana, Danny McPherson, Russ White, Alex D. Zinin, and Liem Nguyen, “OSPF Stub Router Advertisement,” RFC 3137, June 2001.
- [4] Petr Lapukhov and Ariff Premji, “Using BGP for routing in large-scale data centers,” Internet Draft, work in progress, April 2013, **draft-lapukhov-bgp-routing-large-dc-04**
- [5] Daniel Walton, John Scudder, Enke Chen, and Alvaro Retana, “Advertisement of Multiple Paths in BGP,” Internet Draft, work in progress, December 2012, **draft-ietf-idr-add-paths-08**
- [6] T. Sridhar, “Cloud Computing—A Primer, Part 1: Models and Technologies,” *The Internet Protocol Journal*, Volume 12, No. 3, September 2009.
- [7] T. Sridhar, “Cloud Computing—A Primer, Part 2: Infrastructure and Implementation Topics,” *The Internet Protocol Journal*, Volume 12, No. 4, December 2009.

RUSS WHITE is a Principle Research Engineer at Verisign, where he works on the intersection of naming and routing. In the more than 20 years since he first began working in computer networking, he has co-authored 8 technical books and more than 30 patents; he has participated in the writing, editing, and guiding of numerous Internet Standards, and he has written a fiction novel. He is currently working on *The Art of Network Architecture*, to be published by Cisco Press in 2013. Russ splits his time between the Raleigh, N.C., area and Oak Island, N.C.; he teaches in a local homeschool coop and attends Shepherds Theological Seminary. He is a regular blogger and guest on the *Packet Pushers* podcast.

E-mail: **riwhite@verisign.com**

ALVARO RETANA is a Distinguished Engineer in Cisco Technical Services, where he works on strategic customer enablement. Alvaro is widely recognized for his expertise in routing protocols and network design and architecture; he has CCIE® and CCDE® certifications, and he is one of a handful of people who have achieved the CCAR® certification. Alvaro is an active participant in the IETF, where he co-chairs the Routing Area Working Group (rtgwg), is a member of the Routing Area Directorate, and has authored several RFCs on routing technology. Alvaro has published 4 technical books and has been awarded more than 35 patents by the U.S. Patent and Trademark Office. His current interests include software-defined networking, energy efficiency, infrastructure security, routing protocols, and other related topics. E-mail: **aretana@cisco.com**

Letter to the Editor

Dear Editor,

I enjoyed reading the article on “Address Authentication” in the March 2013 edition of *The Internet Protocol Journal* (Volume 16, No. 1), but I couldn’t help thinking to myself how the widespread adoption of the use of *IPv6 Privacy Addresses* (RFC 4941) would affect some of the assertions in the article about the relative merits of using IPv6 addresses for authentication. With both Microsoft and Apple operating systems now implementing IPv6 Privacy Addresses, it is now effectively impossible for any user authentication service to assume that a presented IPv6 address is going to remain constant over time. It is probably safer to assume that such IPv6 addresses are in fact not constant at all, and not to use them in any context of authentication. Given that the widespread use of NATs in IPv4 leads one to the same basic conclusion about using IPv4 addresses for authentication, isn’t the best advice these days to avoid “Address Authentication” as it is applied to Internet end users?

Regards,

—Geoff Huston
gih@apnic.net

The author responds:

I agree with Geoff’s comments. My article explores the idea that IPv6 may be more “trustworthy,” but it concludes by recommending against using any IP address as a form of authentication.

IPv4 addresses will be far less “trustworthy” with the introduction of *Carrier-Grade NATs* or *Large-Scale NATs*. We will not be able to trust IPv6 addresses if the interface identifier changes frequently. My expectation is that most enterprises would prefer *Dynamic Host Configuration Protocol for IPv6* (DHCPv6) with randomized interface identifiers, but most broadband Internet access subscribers will use a *Customer Premises Equipment* (CPE) that uses *Stateless Address Autoconfiguration* (SLAAC) and Stateless DHCPv6. IPv6 offers the ability to perform traceback to the /64 subnet level. That feature is only slightly better than IPv4 traceback.

—Scott Hogg
scott@hoggnetwork.com

Book Review

Network Geeks

Network Geeks: How They Built the Internet, by Brian E. Carpenter, Copernicus Books, ISBN 978-1-4471-5024-4, 2013.

The movie opens on a familiar scene, toward the end of a congenial dinner party at the plush home of an august personage. Conversation has been casual and wide-ranging. The group retires to the library for brandy, cigars, and more conversation. Because you are new to your profession and the august personage was involved in its early years, you ask him what it was like. As he begins his recitation, the scene fades to an earlier time... “My great-grandfather, John Winnard, was born in Wigan...”

Such is the style of Brian Carpenter’s book, *Network Geeks: How They Built the Internet*. Although indeed many other people are cited, the book really is Brian’s personal memoir, complete with his own photographs. It explores his background and work, providing a fascinating travelogue of one person’s arc through recent history. Given the breadth and scale of the 50-year process of invention and development of the global Internet, we need perhaps a thousand more such reminiscences to provide sufficiently rich detail about the many actors and acts that contributed to its success.

Brian’s experiences within that global history are certainly worthy of note. His writing paints pictures of places and topics such as the forces and attractions that drew him to computer networking; in those days, it was an outlier technical topic and people often happened into it, rather than setting out with a plan. Indeed, Brian’s doctoral work was in computer speech understanding—not networking. However, he has played a key role in many significant Internet activities. His frequent employer, the Swiss CERN^[10], was a focal point for much of the early European networking activity—as well as being the birthplace of the World Wide Web—and Brian’s various leadership roles in the *Internet Engineering Task Force* (IETF) came at pivotal times. Other popular references to Internet history tend to emphasize its American basis, making Brian’s primarily European perspective refreshing and helpful.

The book is short, just 150 pages. Although Brian makes some terse references early in the book, he does not get fully into gear talking about the Internet until a third of the way through it. He started in physics, coming fully to computer science only in graduate school. Over the course of the memoir, we hear quite a bit about his physics work at CERN and elsewhere, as well as his activities with the early European deployment of Internet services, his eventual work with Internet standards, and the like.

The IETF

Brian's reference to his great-grandfather does appear, but not until page 10 in a chapter that extensively details his family history and his own upbringing—how many other books on Internet history are likely to include an inset distinguishing the English Baptist church from the American Southern Baptist? Rather, the book begins with a description of a prototypical IETF plenary session at the thrice-annual standards meeting, and he paints the picture well enough to have prompted a guessing game about the person he was describing. IETF meetings, including the plenaries, have a great deal of audience participation, because these meetings are working meetings, not conferences. I particularly enjoyed Brian's turn of phrase when describing one participant, "...who had given several articulate but incomprehensible arguments at the microphone." Later in the book he also equitably describes a colleague as "a wise leader, decisive or even pig-headed, but willing to listen..."

After its opening sequences, the book follows Brian's life chronology, including extended periods in England, Switzerland, the United States, and New Zealand, most recently landing at the last. His employment has variously been university, research, and corporate, including roles as researcher, manager, chair, and teacher.

This book is a memoir, so Brian casually and regularly moves between discussion of personal and professional developments. From one paragraph to the next, he might describe structural aspects of an Internet organization, insulation of housing in New Zealand, the next effort at particle physics, optimizing travel when flying out of southeast England, the nature of a computer networking technology, or the personal style of a co-worker.

In particular, this work is not a tutorial on Internet technology or on its invention. Although Brian does discuss many aspects of the technologies, the pedagogy suits an after-dinner evening's reminiscences, not a classroom lecture. Some concepts are explained in great detail, while others are merely cited. For example, his early discussion of computer networking references the fact that it enables mesh topologies, in contrast to then-common star configurations, but he doesn't give much sense of what "mesh" means in technical terms. Also, the core technology of networking is *packet-switching* and although his discussion on the page after the mesh reference cites queuing theory, he never introduces the motivating design construct of "store and forward."

His discussion of addressing suggests his hardware background, and misses the essence that a name at one level of architecture is often an address at the next level up. So although `www.example.com` is the "name" of a host system attached to the Internet, it has the role of "address" in a URL, because it specifies where to go to resolve the remainder of the URL.

That said, quibbling with such an issue in a tutorial might be reasonable, but it is entirely inappropriate for a memoir. These are Brian's recollections. If they prompt the reader to explore things later, so much the better; but arguing his view will not do. Perhaps reflexively, it is convenient that the Internet makes such exploration quite easy...

NATs

Except that I remain sorry to see that Brian still has such a strikingly purist view about *Network Address Translation* (NAT)^[1, 2] devices, which map between internal (private) IP addresses and public ones. The purist view is that they are an abomination that breaks the elegance of the “end-to-end” design principle of the Internet. The principle is powerful, because it tends to greatly simplify the communications infrastructure and greatly enable innovation at the endpoints. The problem is that the real world imposes organizational and operational models that are more complex than easily supported by the basic end-to-end construct, at the least needing to include enterprise-level policies. NATs do cause problems, by replacing one IP address for another, and some mechanisms do cease to work because of these replacements. However, the operational world views NATs as being useful against multiple problems. One is address space constraints, which is the formal justification for creating the mechanism: an enterprise uses far fewer public IP addresses—a reality that is now essential as IPv4 addresses have grown scarce. Another justification is the misguided view that they improve enterprise security, and the other is the legitimate view that they simplify enterprise network administration. After more than 20 years of extensive deployment, these devices might be expected to have become tolerable to a pragmatist, possibly even forcing consideration of a more elaborate architectural model for the Internet. Yet Brian suffers no such weakness; NATs are evil.

One of the technical points that intrigued me was Brian's repeated discussion of the *Remote Procedure Call* (RPC). This mechanism makes network interaction for an application look like little more than a subroutine invocation. It was hoped that it would greatly simplify network-oriented programming and make it accessible to any software developer, rather than requiring the developer to have a deep understanding of networking interfaces and dynamics. Brian cites the mechanism as having been “invented by the ARPANET community in the mid-1970s...” and used at CERN in a programming language shortly after that. But my own recollection is of hearing a Xerox *Palo Alto Research Center* (PARC) manager in 1980 proudly announce that one of his summer interns had just developed the idea. Indeed, Wikipedia credits the late Bruce Jay Nelson, a Carnegie Mellon University graduate student who was working at PARC.^[3]

And that is the essence of a memoir. It is the remembrances of the speaker, not the formal work of a historian or journalist. It is not the diligent unfolding of a researched history, such as in *Where Wizards Stay up Late*^[4], nor the tourist approach of *Exploring the Internet: A Technical Travelogue*^[5] that seeks to name every possible person active at the time—although Brian does sometimes invoke that latter template. Instead it shares one person's sense of what happened—what he remembers doing and seeing.

Railing against architectural biases or historical nuances is essential when evaluating formal professional writing, and we do need such judicious efforts to capture the history of the Internet. But had Brian sought to produce such a tome, it would not have been as rich or as personal.

References

- [0] European Organization for Nuclear Research, Geneva,
<http://home.web.cern.ch/>
- [1] Network Address Translation,
https://en.wikipedia.org/wiki/Network_address_translation
- [2] Geoff Huston, "Anatomy: A Look inside Network Address Translators," *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [3] Remote Procedure Call,
http://en.wikipedia.org/wiki/Remote_procedure_call
- [4] Katie Hafner and Matthew Lyon, *Where Wizards Stay Up Late: The Origins of the Internet*, Simon & Schuster, ISBN 0-684-81201-0, 1996.
- [5] Carl Malamud, *Exploring the Internet: A Technical Travelogue*, Prentice-Hall, Inc., ISBN 0-13-296898-3 1992/1997,
<http://museum.media.org/eti/>

—Dave Crocker,
dcrocker@bbiw.net

Number of IPv6-Connected Internet Users Doubles

The *Internet Society* (ISOC) recently reported that the number of IPv6-connected users has doubled since *World IPv6 Launch* began on June 6, 2012, when thousands of *Internet Service Providers* (ISPs), home networking equipment manufacturers, and Web companies around the world came together to permanently enable the next generation of *Internet Protocol Version 6* (IPv6) for their products and services. This marks the third straight year IPv6 use on the global Internet has doubled. If current trends continue, more than half of Internet users around the world will be IPv6-connected in less than 6 years.

“The year since World IPv6 Launch began has cemented what we know will be an increasing reality on the Internet: IPv6 is ready for business,” said Leslie Daigle, the Internet Society’s Chief Internet Technology Officer. “Forward-looking network operators are successfully using IPv6 to reduce their dependency on expensive, complex network address translation systems (*Carrier Grade Network Address Translators*) to deal with a shortage of IPv4 addresses. Leaders of organizations that aspire to reach all Internet users must accelerate their IPv6 deployment plans now, or lose an important competitive edge.”

As IPv6 adoption continues to grow, members of the worldwide Internet community are contributing to its deployment. Statistics reported by World IPv6 Launch participants underscore the increasing deployment of IPv6 worldwide:

- Google reports the number of visitors to its sites using IPv6 has more than doubled in the past year.
- The number of networks that have deployed IPv6 continues to grow, with more than 100 worldwide reporting significant IPv6 traffic.
- Australian ISP Internode reports that 10 percent of its customers now use IPv6 to access the Internet.
- Akamai reports that it is currently delivering approximately 10 billion requests per day over IPv6, which represents a 250 percent growth rate since June of last year.
- KDDI measurement shows that the number of IPv6 users of KDDI has doubled and that IPv6 traffic has increased approximately three times from last year.

World IPv6 Launch participants have worked together to help drive adoption, leading to the creation of *World IPv6 Day* in 2011, in which hundreds of websites joined together for a successful global 24-hour test flight of IPv6.

This was followed by World IPv6 Launch in 2012, in which more than a thousand participants permanently enabled IPv6 for their products and services, including four of the most visited websites: Google, Facebook, YouTube, and Yahoo!.

As a platform for innovation and economic development, the Internet plays a critical role in the daily lives of billions. This momentum has not slowed—IPv6 adoption continues to skyrocket, fast establishing itself as the “new normal” and a must-have for any business with an eye towards the future.

For more information about companies that have deployed IPv6, as well as links to useful information for users and how other companies can participate in the continued deployment of IPv6, please visit: <http://www.worldipv6launch.org>

IPv4 has approximately four billion IP addresses (the sequence of numbers assigned to each Internet-connected device). The explosion in the number of people, devices, and web services on the Internet means that IPv4 is running out of space. IPv6, the next-generation Internet protocol which provides more than 340 trillion, trillion, trillion addresses, will connect the billions of people not connected today and will help ensure the Internet can continue its current growth rate indefinitely.

The Internet Society is the trusted independent source for Internet information and thought leadership from around the world. With its principled vision and substantial technological foundation, the Internet Society promotes open dialogue on Internet policy, technology, and future development among users, companies, governments, and other organizations. Working with its members and Chapters around the world, the Internet Society enables the continued evolution and growth of the Internet for everyone. For more information, visit: <http://www.internetsociety.org>

RIPE NCC Report on ITU WTPF-13

The RIPE NCC has published a report on the recent *ITU World Telecommunications/ICT Policy Forum* (WTPF-13). The report is available from the following URL:

<https://www.ripe.net/internet-coordination/news/ripe-ncc-report-on-the-itu-wtpf-13>

Any comments or questions are welcome on the RIPE Cooperation Working Group mailing list:

<https://www.ripe.net/ripe/mail/wg-lists/cooperation>

Google.org Awards Grant to ISOC to Advance IXPs in Emerging Markets

The *Internet Society* (ISOC) recently announced that it has been awarded a grant by Google.org to extend its *Internet Exchange Point* (IXP) activities in emerging markets. The grant will build on the Internet Society's previous efforts and will establish a methodology to assess IXPs, provide training for people to operate the IXPs, and build a more robust local Internet infrastructure in emerging markets.

IXPs play an important role in Internet infrastructure that allows *Internet Service Providers* (ISPs) and other network operators to exchange traffic locally and more cost effectively, which can help lower end-user costs, speed-up transmissions, increase Internet performance, and decrease international Internet connectivity costs. The Internet Society and Internet technical experts have been working for several years to bring IXPs to emerging markets. These efforts have resulted in locally trained experts and facilitated the development of local and regional technical infrastructures. An additional benefit of IXP development is the expansion of community governance models as well as building local Internet expertise.

Google.org, a team within Google focused on social impact, develops and supports technology solutions that can address global challenges, such as expanding Internet access to more of the world's seven billion people.

"The Internet Society has proved to be one of the most effective institutions in the Internet community," said Vint Cerf, vice president and Chief Internet Evangelist at Google. "I am confident that they will apply their grant wisely to extend their work to increase Internet access for everyone, including those in emerging markets."

Lynn St. Amour, President and CEO of the Internet Society, stated, "We are very excited to receive this grant from Google.org. With support to extend our IXP development and improvement projects, we can more quickly bring core Internet infrastructure to underserved countries and assist in building key human and governance capabilities. We will also be able to extend the Internet Society's mission to ensure the open development, evolution, and use of the Internet for the benefit of people everywhere. We look forward to working with Google.org, and we are committed to collaborating with Internet community partners around the world on this important project."

What is my “Subscription ID” for The Internet Protocol Journal (IPJ) and where do I find it?

IPJ Subscription FAQ

Your Subscription ID is a unique combination of letters and numbers used to locate your subscription in our database. It is printed on the back of your IPJ issue or on the envelope. You will also find information about your subscription expiration date near your Subscription ID. Here is an example:



How do I renew or update my subscription?

From the IPJ homepage (www.cisco.com/ipj) click “Subscriber Service” and then enter your Subscription ID and your e-mail address in the boxes. After you click “Login” the system will send you an e-mail message with a unique URL that allows access to your subscription record. You can then update your postal and e-mail details, change delivery options, and of course *renew* your subscription.

What will you use my e-mail address and postal address for?

This information is used *only* to communicate with you regarding your subscription. You will receive renewal reminders as well as other information about your subscription. We will never use your address for any form of marketing or unsolicited e-mail.

I didn’t receive the special URL that allows me to renew or update my Subscription. Why?

This is likely due to some form of spam filtering. Just send an e-mail message to ipj@cisco.com with your Subscription ID and any necessary changes and we will make the changes for you.

Do I need my Subscription ID to read IPJ online? What is my username and password?

Your Subscription ID is used *only* for access to your subscription record. No username or password is required to read IPJ. All back issues are available for online browsing or for download at www.cisco.com/ipj

I can’t find my Subscription ID and I have since changed e-mail address anyway; what do I do now?

Just send a message to ipj@cisco.com and we will take care of it for you.

Call for Papers

The Internet Protocol Journal (IPJ) is published quarterly by Cisco Systems. The journal is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. The journal carries tutorial articles (“What is...?”), as well as implementation/operation articles (“How to...”). It provides readers with technology and standardization updates for all levels of the protocol stack and serves as a forum for discussion of all aspects of internetworking.

Topics include, but are not limited to:

- Access and infrastructure technologies such as: ISDN, Gigabit Ethernet, SONET, ATM, xDSL, cable, fiber optics, satellite, wireless, and dial systems
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, network computing, and Quality of Service
- Application and end-user issues such as: e-mail, Web authoring, server technologies and systems, electronic commerce, and application management
- Legal, policy, and regulatory topics such as: copyright, content control, content liability, settlement charges, “modem tax,” and trademark disputes in the context of internetworking

In addition to feature-length articles, IPJ contains standardization updates, overviews of leading and bleeding-edge technologies, book reviews, announcements, opinion columns, and letters to the Editor.

Cisco will pay a stipend of US\$1000 for published, feature-length articles. Author guidelines are available from Ole Jacobsen, the Editor and Publisher of IPJ, reachable via e-mail at ole@cisco.com

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.



The Internet Protocol Journal, Cisco Systems
170 West Tasman Drive
San Jose, CA 95134-1706
USA

ADDRESS SERVICE REQUESTED

PRSRT STD
U.S. Postage
PAID
PERMIT No. 5187
SAN JOSE, CA

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

David Farber
Distinguished Career Professor of Computer Science and Public Policy
Carnegie Mellon University, USA

Peter Löthberg, Network Architect
Stupi AB, Sweden

Dr. Jun Murai, General Chair Person, WIDE Project
Vice-President, Keio University
Professor, Faculty of Environmental Information
Keio University, Japan

Dr. Deepinder Sidhu, Professor, Computer Science &
Electrical Engineering, University of Maryland, Baltimore County
Director, Maryland Center for Telecommunications Research, USA

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

*The Internet Protocol Journal is
published quarterly by the
Chief Technology Office,
Cisco Systems, Inc.
www.cisco.com
Tel: +1 408 526-4000
E-mail: ipj@cisco.com*

*Copyright © 2013 Cisco Systems, Inc.
All rights reserved. Cisco, the Cisco
logo, and Cisco Systems are
trademarks or registered trademarks
of Cisco Systems, Inc. and/or its
affiliates in the United States and
certain other countries. All other
trademarks mentioned in this document
or Website are the property of their
respective owners.*

Printed in the USA on recycled paper.



The Internet Protocol Journal

September 2014

Volume 17, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Gigabit Wi-Fi.....	2
A Question of DNS Protocols.....	11
Fragments	24
Call for Papers.....	26
Supporters and Sponsors	27

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

FROM THE EDITOR

It is with great pleasure and gratitude that I announce the re-launch of *The Internet Protocol Journal* (IPJ) after a hiatus of just over one year. To subscribe, send e-mail to ipj@protocoljournal.org and we will send you further information. Please bear with us as we deploy a new subscription system and web page.

The re-launch of IPJ is made possible by the support of The Internet Society, which provides administrative functions, and Cisco Systems, which allows us to use the subscriber list and journal name under license. Sponsorship for IPJ is provided by: Afiliás; Asia Pacific Internet Association (APIA); Asia Pacific Network Information Centre (APNIC); Cisco Systems; Comcast; Dyn; Google, Inc.; The Internet Corporation for Assigned Names and Numbers (ICANN); Juniper Networks; Limelight Networks; Netnod; Network Startup Resource Center (NSRC); Réseaux IP Européens Network Coordination Centre (RIPE NCC); Stichting Internet Domeinregistratie Nederland (SIDN); Team Cymru; Verisign Labs; Widely Integrated Distributed Environment (WIDE); 21Vianet Group; Dave Crocker; Jay Etchings; Dennis Jennings; Jim Johnston; Bill Manning; George Sadowsky; Helge Skrivervik; Rob Thomas; and Tom Vest. If you are interested in sponsoring IPJ, please contact us at ipj@protocoljournal.org

We have augmented our Editorial Advisory Board and welcome Fred Baker, Cisco Fellow, Cisco Systems; Steve Crocker, Chairman, ICANN; Geoff Huston, Chief Scientist, APNIC; and Olaf Kolkman, Chief Internet Technology Officer, The Internet Society. Our thanks go to our outgoing members David Farber, Deepinder Sidhu, and Peter Löthberg. The complete Editorial Advisory Board is listed on the back cover.

The increasing demand for more bandwidth in all aspects of networking has led to technology changes in wide-area, local-area, and mobile networks. Newer, faster technologies have been developed and deployed. In our first article, William Stallings gives an overview of *Gigabit Wi-Fi*, also known as IEEE 802.11ac. This relatively new version of Wi-Fi promises transmission speeds of up to 3.2 Gbps.

In our second article, Geoff Huston describes several Denial-of-Service attacks that were launched using the *Domain Name System* (DNS) as the attack vector and discusses details of DNS protocol interactions.

This journal will continue to cover all aspects of internetworking just as it always has. We welcome your feedback and suggestions!

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

Gigabit Wi-Fi

by William Stallings

Just as businesses and home users have generated a need to extend the Ethernet standard to speeds in the gigabit-per-second (Gbps) range, the same requirement exists for the wireless network technology known as *Wi-Fi*. Accordingly, IEEE 802.11, the committee responsible for wireless LAN standards, has recently introduced two new standards^[1], 802.11ac^[2] and 802.11ad^[3,4], which provide for Wi-Fi networks that operate at well in excess of 1 Gbps. These two new standards build on previous work by the IEEE 802.11 committee, which has introduced numerous versions of the wireless LAN standard over the years (Table 1).

Table 1: IEEE 802.11 Physical Layer Standards

Standard	802.11a	802.11b	802.11g	802.11n	802.11ac	802.11ad
Year introduced	1999	1999	2003	2000	2012	2014
Maximum data-transfer speed	54 Mbps	11 Mbps	54 Mbps	65 to 600 Mbps	78 Mbps to 3.2 Gbps	6.76 Gbps
Frequency band	5 GHz	2.4 GHz	2.4 GHz	2.4 or 5 GHz	5 GHz	60 GHz
Channel bandwidth	20 MHz	20 MHz	20 MHz	20, 40 MHz	40, 80, 160 MHz	2160 MHz
Antenna configuration	1 × 1 SISO	1 × 1 SISO	1 × 1 SISO	Up to 4 × 4 MIMO	Up to 8 × 8 MIMO, MU-MIMO	1 × 1 SISO

The evolution of Wi-Fi from the Mbps range to the Gbps range has required the use of three key technologies to enable the higher data rate: *Multiple-Input, Multiple-Output* (MIMO) antennas, *Orthogonal Frequency-Division Multiplexing* (OFDM), and *Quadrature Amplitude Modulation* (QAM). In this article, we first introduce each of these technologies, with a brief mention of their evolution from simpler technologies, and then look at the two new Gigabit Wi-Fi standards.

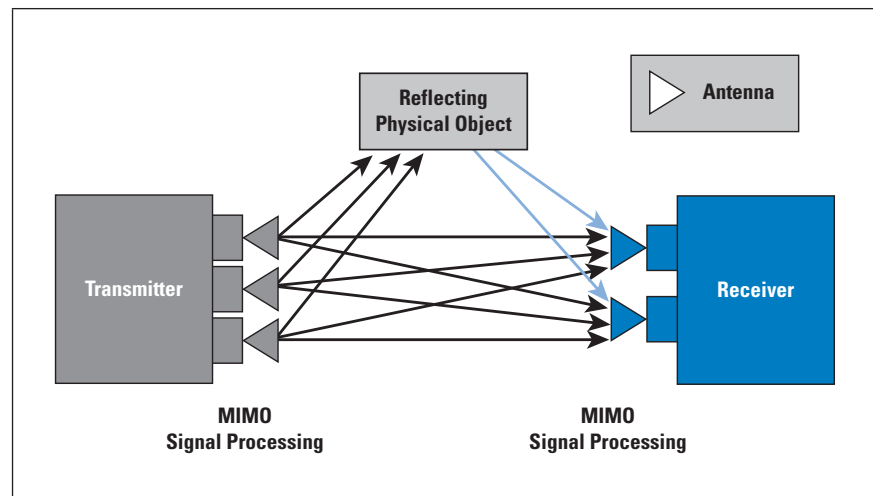
MIMO Antennas

In traditional two-way communication between two wireless stations, each station employs a single antenna for transmission and reception, referred to as *Single-Input, Single-Output* (SISO). In any wireless communication system, there are numerous forms of transmission impairments to deal with, and these impairments become increasingly significant at higher data rates. Of particular concern are noise and *multipath* effects. The latter term refers to the fact that a transmitted signal may reach a destination antenna by not just a direct path but by one or more paths that involve a reflection between source and destination.

These multiple arriving paths interfere with each other and make recovery of the data from the signal more challenging. One effective approach is to use multiple antennas, either at the transmitting end or the receiving end, or both.

In a MIMO scheme, the transmitter and receiver employ multiple antennas^[5]. The source data stream is divided into n substreams, one for each of the n transmitting antennas. The individual substreams are the input to the transmitting antennas (multiple inputs). At the receiving end, m antennas receive the transmissions from the n source antennas via a combination of line-of-sight transmission and multipath caused by reflection (Figure 1). The output signals from all of the m receiving antennas (multiple outputs) are combined. With a lot of complex math, the result is a much better received signal than can be achieved with either a single antenna or multiple frequency channels. Note that the terms *input* and *output* refer to the input to the transmission channel and the output from the transmission channel, respectively.

Figure 1: MIMO Scheme



MIMO systems are characterized by the number of antennas at each end of the wireless channel. Thus an 8×4 MIMO system has 8 antennas at one end of the channel and 4 at the other end. In configurations with a base station, such as a cellular network or a Wi-Fi hotspot, the first number typically refers to the number of antennas at the base station. There are two types of MIMO transmission schemes:

- *Spatial diversity*: The same data is coded and transmitted through multiple antennas, effectively increasing the power in the channel proportional to the number of transmitting antennas. This process improves the *Signal-to-Noise Ratio* (SNR) for cell edge performance. Further, diverse multipath fading offers multiple “views” of the transmitted data at the receiver, thus increasing robustness. In a multipath scenario where each receiving antenna would experience a different interference environment, there is a high probability that if one antenna is suffering a high level of fading, another antenna has sufficient signal level.

- *Spatial multiplexing*: A source data stream is divided among the transmitting antennas. The gain in channel capacity is proportional to the available number of antennas at the transmitter or receiver, whichever is less. Spatial multiplexing can be used when transmitting conditions are favorable and for relatively short distances compared to spatial diversity. The receiver must do considerable signal processing to sort out the incoming substreams, all of which are transmitting in the same frequency channel, and to recover the individual data streams.

Multiple-user MIMO (MU-MIMO) extends the basic MIMO concept to multiple endpoints, each with multiple antennas. The advantage of MU-MIMO compared to single-user MIMO is that the available capacity can be shared to meet time-varying demands. MU-MIMO techniques are used in both Wi-Fi and *Fourth-Generation* (4G) cellular networks.

MU-MIMO has two applications:

- *Uplink—Multiple Access Channel (MAC)*: Multiple end users transmit simultaneously to a single base station.
- *Downlink—Broadcast Channel (BC)*: The base station transmits separate data streams to multiple independent users.

MIMO-MAC is used on the uplink channel to provide multiple access to subscriber stations. In general, MIMO-MAC systems outperform point-to-point MIMO, particularly if the number of receiver antennas is greater than the number of transmit antennas at each user. A variety of multiuser detection techniques are used to separate the signals transmitted by the users.

MIMO-BC is used on the downlink channel to enable the base station to transmit different data streams to multiple users over the same frequency band. MIMO-BC is more challenging to implement. The techniques employed involve processing of the data symbols at the transmitter to minimize interuser interference.

OFDM, OFDMA, and SC-FDMA

The technologies discussed in this section all derive from one of the oldest techniques used in communications: *Frequency-Division Multiplexing* (FDM). FDM simply means the division of a transmission facility into multiple channels by splitting the frequency band transmitted by the facility into narrower bands, each of which is used to constitute a distinct channel. Common examples of FDM are cable TV, broadcast radio, and broadcast television.

A common application of FDM is *Frequency-Division Multiple Access* (FDMA), which is a technique used to share the spectrum among multiple stations. In a typical configuration, a base station communicates with numerous subscriber stations.

Such a configuration is found in satellite networks, cellular networks, Wi-Fi, and WiMAX. Typically, the base station assigns bandwidths to stations within the overall bandwidth available. Key features of FDMA include the following:

- Each channel is dedicated to a single station; it is not shared.
- If a channel is not in use, it is idle and the capacity is wasted.
- Individual channels must be separated by guard bands to minimize interference.

Thus, this scheme divides the available bandwidth into multiple nonoverlapping bands, or channels, as with FDM. The channels are allocated across multiple stations, thus allowing multiple access to the available bandwidth.

Orthogonal Frequency-Division Multiplexing (OFDM), also called *multicarrier modulation*, is a form of FDM in which a single data stream transmits over the available bandwidth, sending some of the bits on each channel. Thus, with OFDM, all of the channels are dedicated to a single data source.

Suppose we have a data stream operating at R bps and an available bandwidth of NB , centered at f . The entire bandwidth could be used to send the data stream, in which case each bit duration would be $1/R$. The alternative is to split the data stream into N substreams, using a serial-to-parallel converter. Each substream has a data rate of R/N bps and is transmitted on a separate subcarrier, with spacing between adjacent subcarriers of B . Now the bit duration is N/R .

The OFDM scheme uses advanced digital-signal-processing techniques to distribute the data over multiple carriers at precise frequencies. The relationship among the subcarriers is referred to as *orthogonality*. The result is that the peaks of the power spectral density of each subcarrier occur at a point at which the power of other subcarriers is zero. With OFDM, the subcarriers can be packed tightly together because there is minimal interference between adjacent subcarriers.

OFDM has several advantages. First, frequency selective fading affects only some subcarriers and not the whole signal. If the data stream is protected by a forward error-correcting code, this type of fading is easily handled. More important, OFDM overcomes *Intersymbol Interference* (ISI) in a multipath environment. ISI has a greater impact at higher bit rates, because the distance between bits, or symbols, is smaller. With OFDM, the data rate is reduced by a factor of N , increasing the symbol time by a factor of N . Thus, if the symbol period is T for the source stream, the period for the OFDM signals is NT . This modulation scheme dramatically reduces the effect of ISI. As a design criterion, N is chosen so that NT is significantly greater than the root-mean-square delay spread of the channel.

A variant of OFDM is *Orthogonal Frequency-Division Multiple Access* (OFDMA). Like OFDM, OFDMA employs multiple closely spaced subcarriers, but the subcarriers are divided into groups of subcarriers. Each group is named a subchannel. The subcarriers that form a subchannel need not be adjacent. In the downlink, different subchannels may be intended for different receivers. In the uplink, a transmitter may be assigned one or more subchannels.

Subchannelization defines subchannels that can be allocated to *Subscriber Stations* (SSs) depending on their channel conditions and data requirements. Using subchannelization, within the same time slot a *Base Station* (BS) can allocate more transmit power to user devices (SSs) with lower SNR, and less power to user devices with higher SNR. Subchannelization also enables the BS to allocate higher power to subchannels assigned to indoor SSs, resulting in better in-building coverage. Subchannels are further grouped into *bursts*, which can be allocated to wireless users. Each burst allocation can be changed from frame to frame as well as within the modulation order, allowing the base station to dynamically adjust the bandwidth usage according to the current system requirements.

Subchannelization in the uplink can save user-device transmit power because it can concentrate power only on certain subchannel(s) allocated to it. This power-saving feature is particularly useful for battery-powered user devices.

Another variant of OFDM is *Single-Carrier FDMA* (SC-FDMA), which is a relatively recently developed multiple access technique with similar structure and performance to OFDMA. One prominent advantage of SC-FDMA over OFDMA is the lower *Peak-to-Average Power Ratio* (PAPR) of the transmit waveform, which benefits the mobile user in terms of battery life and power efficiency. OFDMA signals have a higher PAPR because, in the time domain, a multicarrier signal is the sum of many narrowband signals. At some time instances, this sum is large and at other times small, meaning that the peak value of the signal is substantially larger than the average value.

Thus, SC-FDMA is superior to OFDMA. However, it is restricted to uplink use because the increased time-domain processing of SC-FDMA would entail considerable burden on the base station. SC-FDMA performs a complex digital-signal-processing operation, which spreads the data symbols over all the subcarriers carrying information and produces a virtual single-carrier structure. This structure then is passed through the OFDM processing modules to split the signal into subcarriers. Now, however, every data symbol is carried by every subcarrier.

For OFDM, a source data stream is divided into N separate data streams and these streams are modulated and transmitted in parallel on N separate subcarriers, each with bandwidth B . The source data stream has a data rate of R bps, and the data rate on each subcarrier is R/N bps. For SC-FDMA, it appears that the source data stream is modulated on a single carrier (hence the SC prefix to the name) of bandwidth $N \times B$ and transmitted at a data rate of R bps. The data is transmitted at a higher rate, but over a wider bandwidth compared to the data rate on a single subcarrier of OFDM. However, because of the complex signal processing of SC-FDMA, the preceding description is not accurate. In effect, the source data stream is replicated N times, and each copy of the data stream is independently modulated and transmitted on a subcarrier, with a data rate on each subcarrier of R bps. Compared with OFDM, we are transmitting at a much higher data rate on each subcarrier, but because we are sending the same data stream on each subcarrier, it is still possible to reliably recover the original data stream at the receiver.

A final observation concerns the term *multiple access*. With OFDMA, it is possible to simultaneously transmit either from or to different users by allocating the subcarriers during any one time interval to multiple users. This transmission is not possible with SC-FDMA: At any given point in time, all of the subcarriers are carrying the identical data stream and hence must be dedicated to one user. But over time, it is possible to provide multiple access by allocating the bandwidth to different users at different times.

Quadrature Amplitude Modulation

To transmit digital data over an analog signal, such as a Wi-Fi radio signal, it is necessary to encode the data onto the signal by some form of modulation. The simplest approach is to provide two different signals to be transmitted during a bit time, with one signal element representing binary one and one representing binary zero. Thus, *Amplitude Shift Keying* (ASK) involves transmitting a constant-frequency signal but varying the signal amplitude between two values. With *Phase Shift Keying* (PSK), two different phase shifts of the same carrier frequency are used to represent the two binary digits. As the data rate increases, the length of each signal element representing a single bit shortens. That is, the signal element is shorter both in duration and in physical length while being transmitted. Thus, a short noise burst or a short transmission impairment of any sort affects more bits as the data rate increases. One standard way of coping with this problem is to encode more than a single bit in each signal element. For example, if four amplitudes are used instead of two, then each signal element can encode two bits. One of the most effective techniques for encoding multiple bits per signal element is *Quadrature Amplitude Modulation* (QAM).

QAM uses two basic principles for encoding digital data onto an analog signal: ASK and PSK. QAM takes advantage of the fact that it is possible to send two different signals simultaneously on the same carrier frequency by using two copies of the carrier frequency, one shifted by 90° with respect to the other. For QAM, each carrier is ASK modulated. The two independent signals are simultaneously transmitted over the same medium. At the receiver, the two signals are demodulated and the results are combined to produce the original binary input.

If two-level ASK is used, then each of the two streams can be in one of two states and the combined stream can be in one of $4 = 2 \times 2$ states. If four-level ASK is used (that is, four different amplitude levels), then the combined stream can be in one of $16 = 4 \times 4$ states. This modulation is known as 16-QAM. Systems using 64 (64-QAM) and even 256 states have been implemented. The greater the number of states, the higher the data rate that is possible within a given bandwidth. However, the greater the number of states, the higher the potential error rate due to noise and attenuation.

IEEE 802.11ac

IEEE 802.11ac operates in the 5-GHz band, as do the older and slower standards 802.11a and 802.11n. It is designed to provide a smooth evolution from 802.11n. This new standard uses advanced technologies in antenna design and signal processing to achieve much greater data rates, at lower battery consumption, all within the same frequency band as the older versions of Wi-Fi. The new standard achieves much higher data rates than 802.11n by means of enhancements in three areas:

- *Bandwidth:* The maximum bandwidth of 802.11n is 40 MHz; the maximum bandwidth of 802.11ac is 160 MHz.
- *Signal encoding:* The 802.11n standard uses 64 QAM with OFDM, and 802.11ac uses 256 QAM with OFDM. Thus, more bits are encoded per symbol. Both schemes use forward error correction with a code rate of $5/6$ (ratio of data bits to total bits).
- *MIMO:* With 802.11n, the maximum number of antennas is 4 channel input and 4 channel output antennas. The 802.11ac standard increases this maximum to 8×8 .

Two other changes going from 802.11n to 802.11ac are noteworthy. The 802.11ac standard includes the option of MU-MIMO, meaning that on the downlink, the transmitter can use its antenna resources to transmit multiple frames to different stations, all at the same time and over the same frequency spectrum. Thus, each antenna of a MU-MIMO access point can simultaneously communicate with a different single-antenna device, such as a smartphone or tablet, thereby enabling the access point to deliver significantly more data in many environments.

IEEE 802.11ad

IEEE 802.11ad is a version of 802.11 operating in the 60-GHz frequency band. This band offers the potential for much wider channel bandwidth than the 5-GHz band, enabling high data rates with relatively simple signal encoding and antenna characteristics. Few devices operate in the 60-GHz band, meaning that communication experiences less interference than in the other bands used for Wi-Fi. However, at 60 GHz, 802.11ad operates in the millimeter range, resulting in some undesirable propagation characteristics:

- Free space loss increases with the square of the frequency, so losses are much higher in this range than in the ranges used for traditional microwave systems.
- Multipath losses can be quite high. *Reflection* occurs when an electromagnetic signal encounters a surface that is large relative to the wavelength of the signal; *scattering* occurs if the size of an obstacle is on the order of the wavelength of the signal or less; and *diffraction* occurs when the wavefront encounters the edge of an obstacle that is large compared to the wavelength.
- Millimeter-wave signals generally don't penetrate solid objects.

For these reasons, 802.11ad is likely to be useful only within a single room. Because it can support high data rates and, for example, could easily transmit uncompressed high-definition video, it is suitable for applications such as replacing wires in a home entertainment system, or streaming high-definition movies from your cell phone to your television.

Prospects for Gigabit Wi-Fi

Gigabit Wi-Fi holds attractions for both office and residential environments, and commercial products are beginning to roll out. In the office environment, the demand for ever greater data rates has led to Ethernet offerings at 10 Gbps, 40 Gbps, and most recently 100 Gbps. These stupendous capacities are needed to support blade servers, heavy reliance on video and multimedia, and multiple offsite broadband connections. At the same time, the use of wireless LANs has grown dramatically in the office setting to meet needs for mobility and flexibility. With the gigabit-range data rates available on the fixed portion of the office LAN, gigabit Wi-Fi is needed to enable mobile users to use the office resources effectively. IEEE 802.11ac is likely to be the preferred gigabit Wi-Fi option for this environment.

In the consumer and residential market, IEEE 802.11ad is likely to be popular as a low-power, short-distance wireless LAN capability with little likelihood of interfering with other devices. IEEE 802.11ad is also an attractive option in professional media production environments in which massive amounts of data need to be moved short distances.

References

- [1] Garber, L., “Wi-Fi Races into a Faster Future,” *Computer*, March 2012.
- [2] Alsabbagh, E.; Yu, H.; and Gallagher, K., “802.11ac Design Consideration for Mobile Devices,” *Microwave Journal*, February 2013.
- [3] Cordeiro, C.; Akhmetov, D.; and Park, M., “IEEE 802.11ad: Introduction and Performance Evaluation of the First Multi-Gbps WiFi Technology,” *Proceedings of the 2010 ACM International Workshop on mmWave Communications: From Circuits to Networks*, 2010.
- [4] Perahia, E., et al., “IEEE 802.11ad: Defining the Next Generation Multi-Gbps Wi-Fi,” *Proceedings, 7th IEEE Consumer Communications and Networking Conference*, 2010.
- [5] Halperin, D., et al., “802.11 with Multiple Antennas for Dummies,” *Computer Communication Review*, January 2010.
- [6] Danielyan, E., “IEEE 802.11,” *The Internet Protocol Journal*, Volume 5, No. 1, March 2002.
- [7] Sridhar, T., “Wi-Fi, Bluetooth and WiMAX—Technology and Implementation,” *The Internet Protocol Journal*, Volume 11, No.4, December 2008.

WILLIAM STALLINGS is an independent consultant and author of numerous books about security, computer networking, and computer architecture. His latest book is *Data and Computer Communications* (Pearson, 2014). He maintains a computer science resource site for computer science students and professionals at ComputerScienceStudent.com. He has a Ph.D. in computer science from M.I.T. He can be reached at ws@shore.net

A Question of DNS Protocols

by Geoff Huston, APNIC

One of the most prominent *Denial of Service* attacks in recent years was one that occurred in March 2013, launched against Spamhaus and Cloudflare.

One description of this attack is at [1]. I'm not sure about the claim that this attack "almost broke the Internet," but with a peak volume of attack traffic of some 120 Gbps, it was nevertheless a very significant attack.

How did the attackers generate such massive volumes of attack traffic? The answer lies in the *Domain Name System* (DNS). The attackers asked about domain names, and the DNS system answered. Something we all do all the time on the Internet. So how can a conventional activity of translating a domain name into an IP address be turned into a massive attack?

A few aspects of the DNS make it a readily coercible means of generating an attack:

- The DNS uses very simple *User Datagram Protocol* (UDP) transactions, where the clients send queries to resolvers, and resolvers send back responses.
- Within the DNS it is possible to send a relatively small query packet and get the resolver to reply with a much larger response. The DNS resolver becomes, in effect, a traffic "amplifier."
- There are many so-called "open resolvers," who are willing to respond to queries from any clients on the Internet. In other words, these resolvers will accept DNS queries from anyone and send DNS responses to anyone. The *Open Resolver Project*^[2] claims that there are some 28 million open resolvers on the Internet at present that represent "a significant threat." That's a disturbingly high number if you are worried about ways to subvert the DNS to launch a platform for such attacks.
- Finally, it appears that way too few networks implement source address egress filtering, as described in *Best Current Practice* (BCP) 38^[3]. This mechanism is a packet filtering mechanism, which if universally implemented, would make it far more challenging to mount attacks that relied on the ability to lie about the source address in an IP packet. If a network implemented the measures described in BCP 38, the network could emit only packets whose source address is reachable through the same network.

The combination of these four factors produces a comprehensive vulnerability for the Internet. In performing experiments about the behavior of the DNS, we see a background level of DNS “probes” that contain a query for “ANY,” often querying the domain `isc.org`. In this case the original UDP request of 64 bytes generates a UDP response of 3,475 bytes, or an amplification factor of 54. If an attacker can send this DNS UDP query to 100 of these open resolvers each second, using a common source address of the intended victim, then the attacker will be generating a query traffic volume of some 51 Kbps, and the victim will receive an incoming traffic load of 2.78 Mbps. If you then enlist a bot army to replicate this simple attack one thousand-fold, then the attack traffic volume has exceeded a gigabit per second. If a query for a domain name that is signed using the *Domain Name System Security Extensions* (DNSSEC) is used, where the query explicitly requests the DNSSEC signature information in the response, the response sizes are far larger, and the unwitting accomplices in such attacks can expand to include the authoritative name servers of DNSSEC-signed domains.

The problem with this attack vector is that, in flight, the traffic looks like all other traffic. These responses are simple DNS responses, and the network carries DNS responses as one of the more common packet types. So simple filtering of all DNS traffic is simply not an option, and we need to look deeper to see how we could mitigate this rather worrisome vulnerability.

This development is not a recent one, and the attack vector has been used for many years. For example, a presentation on this topic was given in May 2006 at an IETF meeting^[4]. The major difference between then and now is that the estimated number of open resolvers that can be used to assist in the attack has jumped from 500,000 to in excess of 25 million!

This topic has appeared in numerous threads of discussion.

One thread of conversation goes along the lines that if everyone implemented the measures described in BCP 38, endpoints would be prevented from emitting DNS query packets with a false source address, thereby preventing these reflection attacks to be mounted using the DNS. Of course this conversation is long-standing, in fact older than BCP 38 itself. BCP 38 is now 13 years old as a document, and the somewhat worrisome observation is that nothing much appears to have happened in terms of improving the situation about source address spoofing over the past 13 years, so there is not a lot of optimism that anything will change in the coming months, if not years.

Another conversation thread says that resolvers should implement *Response Rate Limiting* (RRL), and silently discard repetitive queries that exceed some locally configured threshold. Although this solution is relatively effective in the case of authoritative name servers, it is less effective in the face of a massive pool of open recursive resolvers, because in this latter case the query load can be spread across the entire resolver pool so that each resolver may not experience a detectable level of repeated queries. It is also possible to use a very large pool of authoritative name servers in the same manner. However, this consideration does not weaken the advice that authoritative name servers should implement RRL in any case. It just lowers the level of optimism that this measure alone would eliminate this vulnerability. (Randy Bush gave an informative presentation at the 2013 *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) conference^[5], illustrating the before and after states of the application of RRL for an authoritative name server.)

Still another thread of conversation is exploring ways to shut down the pool of open recursive resolvers. The *Open Resolver Project*^[2] is working on a “name and shame” approach to the problem, and is hopeful that if individual resolver operators are allowed to check their own status, these resolvers will be closed down. Like the universal adoption of BCP 38, it’s hard to be overly optimistic about this approach. Part of the problem is that a large volume of unmanaged or semi-managed systems are connected to the Internet, and the vulnerabilities they create are not necessarily known to the parties who deployed these systems in the first place. For example, one way to create more robust “plug and play” systems is to reduce the amount of environmental configuration that needs to be loaded into such systems in the first place. Equipping such standalone units with their own DNS resolver appears to be a way to remove an additional configuration element from the unit. The downside is that if these units are configured as open recursive resolvers, then they become part of the overall problem of the massive population of open resolvers on today’s Internet.

The behavior of the DNS where a small-size query generates a large response is one that appears to be intrinsic to the DNS, and particularly more so with the use of DNSSEC-signed domain names. If we want some form of security in the DNS so that the client can be assured that the response it receives is authentic and current and has not been tampered with in any way, then the overheads of cryptographic signature blocks and the size these signature blocks take up appears to be an intrinsic part of the security architecture of DNS. We appear to want a lightweight, fast DNS, so we like the performance properties of a DNS resolution framework that uses UDP, coupled with a ubiquitous distribution of recursive resolvers and the associated resolver caches.

But we also want DNS responses that can be verified, so we like DNSSEC. So the responses get larger, and the outcome is that the DNS is a highly effective platform for massive traffic attacks where there is a very limited space to mitigate the associated risks.

If we want to close the door on using the DNS to mount large-scale attacks, there does not appear to be much space left to maneuver here. However, there is a conversation that has not quite petered out yet. That conversation is about the use of the *Transmission Control Protocol* (TCP) in the DNS.

The original specification of the DNS^[6] allowed for the use of both UDP and TCP as the transport service for DNS queries and responses. The relevant discussion of this specification in “Requirements for Internet Hosts—Application and Support”^[7] reads:

“... it is also clear that some new DNS record types defined in the future will contain information exceeding the 512 byte limit that applies to UDP, and hence will require TCP. Thus, resolvers and name servers should implement TCP services as a backup to UDP today, with the knowledge that they will require the TCP service in the future.”

Why 512 bytes? The 512-byte limit was based on the IPv4 “Requirements for Internet Hosts—Communication Layers,”^[8] where it is required that all IPv4 host systems must accept an IP packet that is at least 576 octets in size. Allowing for 20 bytes of IP header, room for the maximum of 40 bytes of IP options and 8 bytes of UDP header, the implication is that the maximum payload in a UDP packet that will be accepted by all IPv4 hosts is restricted to 512 bytes.

It should be noted that it is theoretically possible that there are links that do not support IP packets of 576 bytes, so even these 576-byte IP packets may possibly be fragmented in flight. The limit being referred to here is the largest (possibly reassembled) packet that a host will assuredly accept.

Now of course it is possible to generate larger packets in IPv4, to a theoretical maximum of 65,535 bytes, which, by the same reckoning, accommodates a UDP payload of 65,507 bytes, but such larger packets will probably be fragmented in flight. In such a case, when this behavior is combined with typical firewall behavior, then the trailing packet fragments may not be passed onward to a client at all, because many firewall configurations use acceptance rules based on the UDP and TCP port numbers. Because the trailing fragments of a fragmented IP datagram have no UDP or TCP packet header, the firewall is left with a quandary. Should the firewall simply accept all IP fragments? In this case the security role of the firewall may be compromised through the transmission of fragments that form part of a hostile attack.

Or should the firewall discard all IP fragments? In this case a sender of a large packet should be aware that if there is fragmentation, then the trailing packet fragments may not be passed to the receiver. So with the two considerations that hosts are not required to accept UDP datagrams that are larger than 576 bytes, and firewalls may discard trailing fragments of a fragmented IP datagram, the original DNS response to this situation was to limit all of its UDP responses to 512 bytes, and always use TCP as the backup plan if the DNS response was larger than 512 bytes.

However, clients may not know in advance that a DNS response is larger than 512 bytes. To signal to a client that it should use TCP to retrieve the complete DNS response, when the response would be greater than 512 bytes the DNS resolver will respond in UDP with a partial response that fits within the 512-byte limit, and set the “truncated” bit in the DNS flags that form part of the response.

We continued for some years with this approach. The DNS used UDP for the bulk of its transactions, which all fit within the 512-byte limit, and for the relatively infrequent case where larger DNS responses were being generated, the UDP response was truncated, and the client was expected to repeat the question over a TCP connection.

This situation was not altogether comfortable when we then considered adding security credentials to the DNS through DNSSEC. The inclusion of digital signatures in DNS responses implied that very few DNSSEC responses would fit within this 512-byte limit. But if the process of switching to TCP was to respond to the UDP query with a UDP response that essentially said “Please use TCP,” then this response adds a considerable delay to the DNS function. Each query would now involve a minimum of three round-trip time interactions with the server, rather than just the single round-trip time interval for UDP (a TCP transaction still includes one round-trip time transaction for the UDP query and truncated UDP response, one round-trip interval for the TCP connection establishment handshake, and one for the TCP query and response). The next refinement of the DNS was to include a way to signal that a client was able to handle larger DNS responses in UDP, and thereby bypass the fall-back to TCP in the case where the client is able to handle larger IP packets and the client is confident that there is no IP fragment filtering in intervening middleware.

As pointed out in RFC 5966^[9]:

“Since the original core specifications for DNS were written, the *Extension Mechanisms for DNS* (EDNS0)^[10] have been introduced. These extensions can be used to indicate that the client is prepared to receive UDP responses larger than 512 bytes. An EDNS0-compatible server receiving a request from an EDNS0-compatible client may send UDP packets up to that client’s announced buffer size without truncation.”

EDNS0 allows for the UDP-based DNS response to grow to far higher sizes. As a result, it appears that the DNS largely uses UDP and EDNS0, and EDNS0 is used to allow for the larger-size responses, most commonly set at 4096 bytes. TCP is now often regarded as a somewhat esoteric option, being used only for DNS zone transfer operations, and if the zone does not want to support zone transfers as a matter of local policy, then it is commonly thought that the role of TCP is no longer essential for the DNS.

Section 6.1.3.2 of “Requirements for Internet Hosts—Application and Support”^[7] states:

“DNS resolvers and recursive servers **MUST** support UDP, and **SHOULD** support TCP, for sending (non-zone-transfer) queries.”

In the world of standards specifications and so-called normative language, that “SHOULD” in the quoted text is different from a “MUST.” It’s a little stronger than saying “well, you can do that if you want to,” but it’s a little weaker than saying “You really have to do this. There is no option not to.” Little wonder that some implementors of DNS resolvers and some folks who configure firewalls came to the conclusion that DNS over TCP was an optional part of the DNS specification that does not necessarily need to be supported.

But DNS over TCP is not just a tool to allow for large DNS responses. If we review the preconditions for the use of the DNS in large-scale reflector attacks, namely the widespread support of UDP for large packet responses, and the relatively sparse use of BCP 38, then we’ve effectively allowed attackers to mount reflection attacks by co-opting a large set of open resolvers to send their responses to the target system, by using UDP queries whose IP source address is the IP address of the intended victim.

TCP does not have the same vulnerability. If an attacker were to attempt to open a TCP session using an IP source address of the intended victim, the victim would receive a short IP packet (IP and TCP header only, which is a 40-byte packet) containing only the SYN and ACK flags set. Because the victim system has no preexisting state for this TCP connection, it will discard the packet. Depending on the local configuration, it may send a TCP RESET to the other end to indicate that it has no state, or the discard of this unsolicited packet may be completely silent. This behavior removes one of the essential preconditions for a reflector amplification attack. If the attack traffic with the spoofed source address of the intended victim uses a 40-byte SYN TCP packet, then the victim will receive a 40-byte SYN/ACK TCP packet. The DNS attack amplification factor would be effectively removed.

If the DNS represents such a significant vulnerability for the Internet through these UDP-based reflection attacks, then does TCP represent a potential mitigation? Could we realistically contemplate moving away from the ubiquitous use of EDNS0 to support large DNS responses in UDP, and instead use DNS name servers that limit the maximal size of their UDP responses, and turn to TCP for larger responses? Could we contemplate a rate-limiting approach where, instead of not responding to what are possibly considered to be “excess” queries, the DNS server responds with a truncated UDP response to indicate that the client should use TCP and repeat the query?

Again, let’s turn to RFC 5966:

“The majority of DNS server operators already support TCP, and the default configuration for most software implementations is to support TCP. The primary audience for this document is those implementors whose failure to support TCP restricts interoperability and limits deployment of new DNS features.”

The question we are looking at here is: Can we quantify the extent to which DNS resolvers are capable of using TCP for DNS queries and responses? How big is this “majority of DNS server operators” being referred to in the quote?

The Experiment

We conducted an experiment using a modified DNS name server, where the maximal UDP packet size was configured to 512 bytes, and then set up an experiment where a simple query to resolve a DNS name would generate a response that could not fit within 512 bytes.

Although such a response is relatively easy trigger if it includes DNSSEC, if we want to set up a condition where all DNS responses are larger than 512 bytes for a domain, then we need to use a slightly different approach. The approach used in this iteration of the experiment is to use the DNS name alias function, the *Canonical Name* (CNAME) record.

Here is a sample zone:

```
$TTL1h
@ IN      SOA      nsx.dotnxdomain.net. research.apnic.net. (
                        2013011406      ; Serial
                        3600             ; Refresh
                        900              ; Retry
                        1                ; Expire
                        1 )              ; Minimum
IN        NS       nsz1.z.dotnxdomain.net.

z1        IN       A        199.102.79.186

*         IN       A        199.102.79.186

4a9c317f.4f1e706a.6567c55c.0be33b7b.2b51341.a35a853f.59c4df1d.3b069e4e.87ea53bc.2b4cfb
4f.987d5318.fc0f8f61.3cbe5065.8d9a9ec4.1ddfa1c2.4fee4676.1ffb7fcc.ace02a11.a3277bf4.22
52b9ed.9b15950d.db03a738.dde1f863.3b0bf729 IN CNAME
33d23a33.3b7acf35.9bd5b553.3ad4aa35.09207c36.a095a7ae.1dc33700.103ad556.3a564678.16395
067.a12ec545.6183d935.c68cebfb.41a4008e.4f291b87.479c6f9e.5ea48f86.7d1187f1.7572d59a.9
d7d4ac3.06b70413.1706f018.0754fa29.9d24b07c

33d23a33.3b7acf35.9bd5b553.3ad4aa35.09207c36.a095a7ae.1dc33700.103ad556.3a564678.16395
067.a12ec545.6183d935.c68cebfb.41a4008e.4f291b87.479c6f9e.5ea48f86.7d1187f1.7572d59a.9
d7d4ac3.06b70413.1706f018.0754fa29.9d24b07c IN A 199.102.79.187
```

The use of the combination of long label strings and a CNAME construct forces a large response, which, in turn, triggers DNS to send a truncated UDP response in response to a conventional query for the address record for the original domain name. The truncated UDP response should force the client resolver to open a TCP session with the name server, and ask the same query again.

To ensure that the authoritative name server directly processed every name query, we used unique labels for each presented experiment, and ensured that the DNSSEC signed zones were also unique, ensuring that local DNS resolver caches could not respond to these name queries.

The first question we are interested in is: How many clients were able to successfully switch to TCP following the receipt of a truncated response in UDP?

We conducted an experiment by embedding numerous test cases inside an online ad, and over 8 days at the end of July 2013 we presented these tests to 2,045,287 end clients. We used four experiments. Two experiments used a name where the query and response would fit within a 512-byte payload as long as the query did not include a request for DNSSEC. One of these domain names was unsigned, and the other was signed. The other two experiments used the CNAME approach to ensure that the response would be larger than 512 bytes.

Again, one zone was signed, and the other was unsigned (Table 1).

Table 1: DNS Resolution Using TCP, with CNAME Names

Experiment	UDP Queries	Truncated UDP Responses	TCP Responses	Truncated UDP to TCP Fail
Short, unsigned	2,029,725	2	6	0
Short, signed	2,037,563	1,699,935 (83.4%)	1,660,754 (81.5%)	39,101 (1.9%)
Long, unsigned	2,023,205	2,021,212 (99.9%)	1,968,927 (97.3%)	52,285
Long, signed	2,033,535	2,032,176 (99.9%)	1,978,396 (97.3%)	53,780 (2.6%)

This data appears to point to a level of failure to follow up from a truncated UDP response to a TCP connection of some 2.6 percent of clients.

That level of failure to switch from a truncated UDP response to rephrase the query over TCP is large enough to be significant. The first question: Is this failure due to some failure of the DNS authoritative name server or a failure of the client resolver? If the name server is experiencing a high TCP session load, it will reject new TCP sessions by responding to incoming TCP session establishment packets with a TCP RESET. We saw no evidence of this session overload behavior in the packet traces that the authoritative name server gathered. So the TCP failure looks to be occurring closer to the client resolver than to the authoritative name server.

We can also look at this example from the perspective of the set of DNS resolvers. How many resolvers will switch to use TCP when they receive a UDP response that is truncated? Before looking at the results, it needs to be noted that the only resolvers that are exposed in this experiment are those resolvers that directly query our authoritative name server (*visible* resolvers). If a resolver is configured to forward its queries onto a recursive resolver, then its behavior will not be directly exposed in this experiment. It should also be noted that even when the visible recursive resolver is forced to use TCP to query the authoritative name server, this resolver may still relay the response back to its client by UDP, using EDNS0 to ensure that the larger UDP response is acceptable to the client. Thus the scope of these measurements refers specifically to this subset of resolvers who are visible resolvers.

The 2 million clients in this experiment used a total of 80,505 visible resolvers. Some 13,483 resolvers, or 17 percent of these visible resolvers, did not generate any TCP transactions with the authoritative name server. These 13,483 UDP-only resolvers made a total of 4,446,670 queries, and of these some 4,269,495 responses were truncated, yet none of these resolvers switched to TCP.

There is a second class of filtering middleware that operates on incoming traffic. In such cases the authoritative server sees an incoming TCP SYN packet to establish the DNS connection, and the server responds with a SYN+ACK packet. Because this packet will be blocked by the filtering middleware, it will never get passed through to the resolver client, and the TCP connection will not be established. It has been observed in discussions on DNS reliability that some security middleware permits only inbound traffic on UDP port 53, and discards inbound TCP port 53 traffic. The way in which this filtering behavior would manifest itself at the authoritative name server is that the name server would see and respond to the initial TCP SYN, and not see the ensuing TCP ACK that would complete the TCP handshake. This SYN-only behavior was observed in just 337 resolvers, a count that represents 0.4 percent of the set of visible resolvers. These resolvers generated a total of 1,719,945 queries, and received 1,575,328 truncated UDP responses.

Why is the client-level TCP failure rate at 2.6 percent of clients, whereas at the resolver level the TCP failure rate is 17 percent of visible resolvers? There are at least three possible reasons for this result:

- First, in some cases we observe service providers using DNS forwarder farms, where queries are spread across many query engines. When a DNS query is rephrased using TCP, it may not use the same forwarder to make the query.
- Second, we should factor in end-client failover to another DNS resolver that can support DNS transactions over TCP. Most clients are configured with multiple resolvers, and when one resolver fails to provide a response the client asks the query of the second and subsequent resolvers in its resolver set. If any of the visible resolvers associated with the resolvers listed in the client's resolver set are capable of using TCP, then at some stage we will see a TCP transaction at the authoritative name server. In this more prevalent case of TCP failure, either the resolver itself is not capable of generating a DNS query using TCP (presumably because of local limitations in the resolver software or local configuration settings), or some network middleware is preventing the resolver from performing TCP connections to port 53.
- Finally, the distribution of end clients across the set of visible resolvers is not even, and whereas some resolvers, such as the set used by Google's Public DNS service, serve some 7 percent of all end clients, others serve a single end client. We observed that 53,000 experiments, out of a total of 2 million experiments, failed to complete a TCP-based DNS resolution, so it is also possible that these 13,483 visible resolvers that do not support TCP queries are entirely consistent in volume with this level of end-client failure to resolve the DNS label of the experiment.

There is a slightly different way to look at this question. Although we saw some 53,000 experiments that failed to complete the DNS resolution at all, how many experiments were affected by this deliberate effort to force resolvers to use TCP? How many clients were affected in terms of longer DNS resolution time through the use of DNS resolvers that failed to switch to use TCP?

Table 2 shows that slightly more than 6 percent of all clients used a DNS resolver that was unable to repeat the DNS query over TCP. Some 75,000 clients used an alternate resolver that was capable of performing the TCP query, whereas the remainder were unable to resolve the DNS name at all.

Table 2: Client Use of Resolvers That Fail to Complete a TCP Query

Experiment	UDP Queries	Truncated UDP Responses	TCP Responses	Truncated UDP to TCP Fail	Used TCP Fail Resolvers
Long, unsigned	2,023,205	2,021,212 (99.9%)	1,968,927 (97.3%)	52,285 (2.6%)	124,881 (6.1%)
Long, signed	2,033,535	2,032,176 (99.9%)	1,978,396 (97.3%)	53,780 (2.6%)	129,555 (6.4%)

After running this initial experiment, we considered our use of the CNAME construct to inflate the DNS response to more than 512 bytes, and wondered if this additional DNS indirection created some problems for some resolver clients. Another approach to coerce client resolvers to use TCP is to modify the name-server code used by the authoritative name server, and drop its UDP maximum size to 275 octets, so that the name server will truncate the UDP response for any response of 276 bytes or larger. In this way a DNS query for the short unsigned name would fit within the new UDP limit, but in all other cases the UDP response would be truncated.

The results we saw for this second experiment, which removed the CNAME entry and used an authoritative name server with a 275-byte UDP payload limit, with 3 days of collected data, are summarized in Table 3.

Table 3: DNS Resolution Using TCP, Using 275-Byte UDP Truncation

Experiment	UDP Queries	Truncated UDP Responses	TCP Responses	Truncated UDP to TCP Fail
Short, unsigned	936,007	0	3	3
Short, signed	936,116	936,116 (100.0%)	911,751 (97.4%)	24,365 (2.6%)
Long, unsigned	920,613	920,613 (100.0%)	896,953 (97.4%)	23,530 (2.6%)
Long, signed	934,446	934,446 (100.0%)	910,757 (97.5%)	25,573 (2.5%)

These results are consistent with the results of the original experiment, indicating that the use of the CNAME construct is not causing additional problems for resolver clients.

Conclusion

The original specification of the DNS called for resolvers to use UDP when the response was 512 bytes or smaller, and TCP was to be used for larger DNS transactions. DNS clients would interpret the truncated flag in a DNS UDP response to trigger a re-query using TCP.

With the introduction of EDNS0, clients can now signal their capability to accept larger UDP datagrams, with the result that the fallback to TCP for large DNS responses is used less frequently, to the extent that there is now a concern that a significant set of clients cannot resolve a DNS name if that resolution operation is required to occur using TCP.

However, DNS UDP is being used in various forms of malicious attacks, using DNS queries where the response is far larger than the query. The combination of source address spoofing and DNS over UDP is presenting us with some significant concerns. For that reason there is a renewed consideration of the viability of reverting to TCP for various forms of larger DNS responses, which effectively prevents source address spoofing in the DNS query/response interaction.

In this experiment we've looked at the impact a forced switch to DNS over TCP would have on clients. In particular, what proportion of clients would no longer be able to complete a DNS name-resolution process if the process necessarily involves the use of TCP? Our measurements of a sample of 2 million clients in early August 2013 points to a DNS resolution failure rate for 2 percent of clients.

The picture for individual DNS resolvers appears to be somewhat worse, in that 17 percent of visible resolvers do not successfully follow up with a TCP connection following the reception of a truncated UDP response.

Although that 17-percent number is surprisingly high, there are two mitigating factors here.

It appears that clients use multiple DNS resolvers in their local DNS configuration, so that failure of an initially selected resolver to respond to a query because of a lack of support for TCP may be resolved by the clients selecting the next resolver from their local resolver set. For this set of clients, which appears to encompass some 4 percent of the total client population, the penalty is increased DNS resolution time, where the resolution of a name requires the client to fail over to the other resolvers listed in their local DNS resolver set.

Secondly, the more extensively used visible DNS resolvers appear to be capable of supporting TCP-based queries, so the problems with TCP support in the DNS appear to be predominately concerned with resolvers that are used by a relatively small pool of end clients.

References

- [1] “The DDoS That Almost Broke the Internet,”
<http://blog.cloudflare.com/the-ddos-that-almost-broke-the-internet>
- [2] <http://openresolverproject.org>
- [3] Paul Ferguson, “Network Ingress Filtering: Defeating Denial of Service Attacks which employ IP Source Address Spoofing,” RFC 2827, BCP 38, May 2000.
- [4] Frank Scalzo, “Recent DNS Reflector Attacks From the Victim and the Reflector POV,”
<http://www.iepg.org/july2006/1-frank-scalzo.pdf>
- [5] Randy Bush, “DNS Rate Limiting—a Hard Lesson,”
http://conference.apnic.net/__data/assets/pdf_file/0011/58880/130226.apops-dns-rate-limit_1361839670.pdf
- [6] Paul V. Mockapetris, “Domain Names—Implementation and Specification,” RFC 1035, November 1987.
- [7] Robert Braden, “Requirements for Internet Hosts—Application and Support,” RFC 1123, October 1989.
- [8] Robert Braden, “Requirements for Internet Hosts—Communication Layers,” RFC 1122, October 1989.
- [9] Ray Bellis, “DNS Transport over TCP—Implementation Requirements,” RFC 5966, August 2010.
- [10] Paul Vixie, “Extension Mechanisms for DNS (EDNS0),” RFC 2671, August 1999.

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. E-mail: gih@apnic.net

IANA Transition

In March, 2014, the United States *National Telecommunications and Information Administration* (NTIA) announced its intent to “...transition Key Internet Domain Name Functions to the global multistakeholder community.” See [1] for the full announcement.

Quoting from the announcement: “As the first step, NTIA is asking the *Internet Corporation for Assigned Names and Numbers* (ICANN) to convene global stakeholders to develop a proposal to transition the current role played by NTIA in the coordination of the Internet’s *Domain Name System* (DNS).

NTIA’s responsibility includes the procedural role of administering changes to the authoritative root zone file—the database containing the lists of names and addresses of all top-level domains—as well as serving as the historic steward of the DNS. NTIA currently contracts with ICANN to carry out the *Internet Assigned Numbers Authority* (IANA) functions and has a Cooperative Agreement with Verisign under which it performs related root zone management functions.

ICANN is uniquely positioned, as both the current IANA functions contractor and the global coordinator for the DNS, as the appropriate party to convene the multistakeholder process to develop the transition plan. NTIA has informed ICANN that it expects that in the development of the proposal, ICANN will work collaboratively with the directly affected parties, including the *Internet Engineering Task Force* (IETF), the *Internet Architecture Board* (IAB), the *Internet Society* (ISOC), the *Regional Internet Registries* (RIRs), top level domain name operators, VeriSign, and other interested global stakeholders.

NTIA has communicated to ICANN that the transition proposal must have broad community support and address the following four principles:

- Support and enhance the multistakeholder model;
- Maintain the security, stability, and resiliency of the Internet DNS;
- Meet the needs and expectation of the global customers and partners of the IANA services; and,
- Maintain the openness of the Internet.

While stakeholders work through the ICANN-convened process to develop a transition proposal, NTIA’s current role will remain unchanged. The current IANA functions contract expires September 30, 2015.”

Since the announcement, a lot of discussion has taken place in many fora and comments have been produced by various groups. Below we include some links to various documents, groups and events that captures some of these activities. The IANA transition is likely to remain a “hot topic” for the next couple of years.

References

- [0] Internet Assigned Numbers Authority (IANA)
<http://www.iana.org/about>
- [1] “NTIA Announces Intent to Transition Key Internet Domain Name Functions,”
<http://www.ntia.doc.gov/press-release/2014/ntia-announces-intent-transition-key-internet-domain-name-functions>
- [2] “IANA Functions and Related Root Zone Management Transition Questions and Answers,”
<http://www.ntia.doc.gov/other-publication/2014/iana-functions-and-related-root-zone-management-transition-questions-and-answ>
- [3] “Comments of the Internet Architecture Board (IAB) Regarding the ‘Draft Proposal, Based on Initial Community Feedback, of the Principles and Mechanisms and the Process to Develop a Proposal to Transition NTIA’s Stewardship of the IANA Functions,’”
<http://www.iab.org/wp-content/IAB-uploads/2014/04/iab-response-to-20140408-20140428a.pdf>
- [4] “NTIA IANA Functions’ Stewardship Transition,” ICANN microsite, <https://www.icann.org/stewardship>
- [5] Jari Arkko, “ICANN and Transition of NTIA’s Stewardship,” IETF Blog,
<http://www.ietf.org/blog/2014/07/icann-and-transition-of-ntias-stewardship/>
- [6] ARIN, “IANA Globalization,” <http://teamarin.net/education/internet-governance/iana-globalization/>

The Internet Protocol Journal Needs your Help

We are delighted to be publishing IPJ once again, but we cannot do so without your help. Please consider sponsoring the journal. We have a range of sponsorship levels to suit most budgets, just send an e-mail to ipj@protocoljournal.org for more information. We also need suggestions for topics as well as actual articles, so do get in touch!

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

The Internet Protocol Journal (IPJ) is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. Publication of IPJ is made possible by:

Supporters



Diamond Sponsors



Ruby Sponsor



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



Individual Sponsors

Dave Crocker, Jay Etchings, Dennis Jennings, Jim Johnston, Bill Manning, George Sadowsky, Helge Skrivervik, Rob Thomas, Tom Vest.

For more information about sponsorship, please contact ipj@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

ADDRESS SERVICE REQUESTED

PRSRT STD U.S. Postage PAID PERMIT No. 5187 SAN JOSE, CA
--

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Fred Baker, Cisco Fellow
Cisco Systems, Inc.

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Steve Crocker, Chairman
Internet Corporation for Assigned Names and Numbers

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

December 2014

Volume 17, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
What's Special About 512?	2
ZigBee IP Protocol Stack	19
Letters to the Editor.....	39
Book Review.....	42
Fragments	44
Supporters and Sponsors	47

FROM THE EDITOR

As announced in September, *The Internet Protocol Journal* (IPJ) has been re-launched with the generous support of numerous organizations and individuals. You will find a complete list of our sponsors and supporters on page 47. As we bring you this second issue, I want to say a few words about subscriptions to this journal. If you are already a subscriber to IPJ and in the past have received a printed copy by mail, you will find your Subscription ID printed on the back page along with your delivery address. We will send you an account activation e-mail in the near future that will allow you to update and renew your subscription. However, if you have changed e-mail address you can contact us with the new information by sending a message to ipj@protocoljournal.org. If, on the other hand, you picked up a copy of IPJ at a conference or other event and you wish to create a *new* subscription, just follow the instructions on our website at www.protocoljournal.org, where you will also find all our back issues, index files, and sponsorship information. Subscriptions to IPJ are free of charge.

The growth of the Internet has been a recurring theme in this journal since its inception in 1998. We've covered various aspects of IPv4 address-space depletion and IPv6 deployment, and of course explored many challenges related to *Network Address Translation* (NAT). But address depletion isn't the only scaling issue facing the Internet. In this issue, Geoff Huston explains what happened in August 2014 when the number of reachable networks as announced by the *Border Gateway Protocol* (BGP) exceeded 512,000 for a time.

The Internet is also growing in a different arena, namely that of intelligent embedded systems that use Internet protocols for communication. This emerging technology area is referred to as *The Internet of Things*. In our second article, Douglas Comer describes the *ZigBee IP Protocol Stack*, which is under development and standardization.

As always, we would love your feedback on anything you read in this journal. With your permission we can include your comments in the form of a Letter to the Editor, or you may consider writing a Book Review. Send your message to ipj@protocoljournal.org

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

What's So Special About 512?

by Geoff Huston, APNIC

August 12, 2014 was widely reported as a day when the Internet “collapsed.” Despite the sensational media reports, the condition was not fatal for the Internet, and perhaps it could be more reasonably reported that some parts of the Internet were having a “bad hair day.” The root cause of this event was the Internet routing system, and in this article I will review the behavior of this system and the relationship between the routing system and routing hardware.

A Lightning Introduction to the Internet Routing System

The Internet is a collection of some 50,000 component networks. Some are large, spanning multiple continents, while others are the size of a small office. Most of these networks sit somewhere between these two extremes. All networks announce those IP addresses (or “network address prefixes”) that are located within their network (or “originate” from their network).

The routing protocol used to disseminate these announcements across the Internet is the *Border Gateway Protocol* (BGP)^[0]. To simplify the operation of BGP, these announcements of reachable network address prefixes are not made to each of the 50,000 other networks, but instead are made to their immediately connected adjacent network neighbors (or routing “peers”). If one (or more) of these peer networks is willing and capable of passing traffic through its network to reach the originating network (that is, act as “transit”), then these networks will further announce these network address prefixes to their peers, and so on. In this way a BGP speaker will hear, via its immediately connected peers, announcements for all reachable network address prefixes, and can assemble a complete picture of all reachable addresses on the Internet.

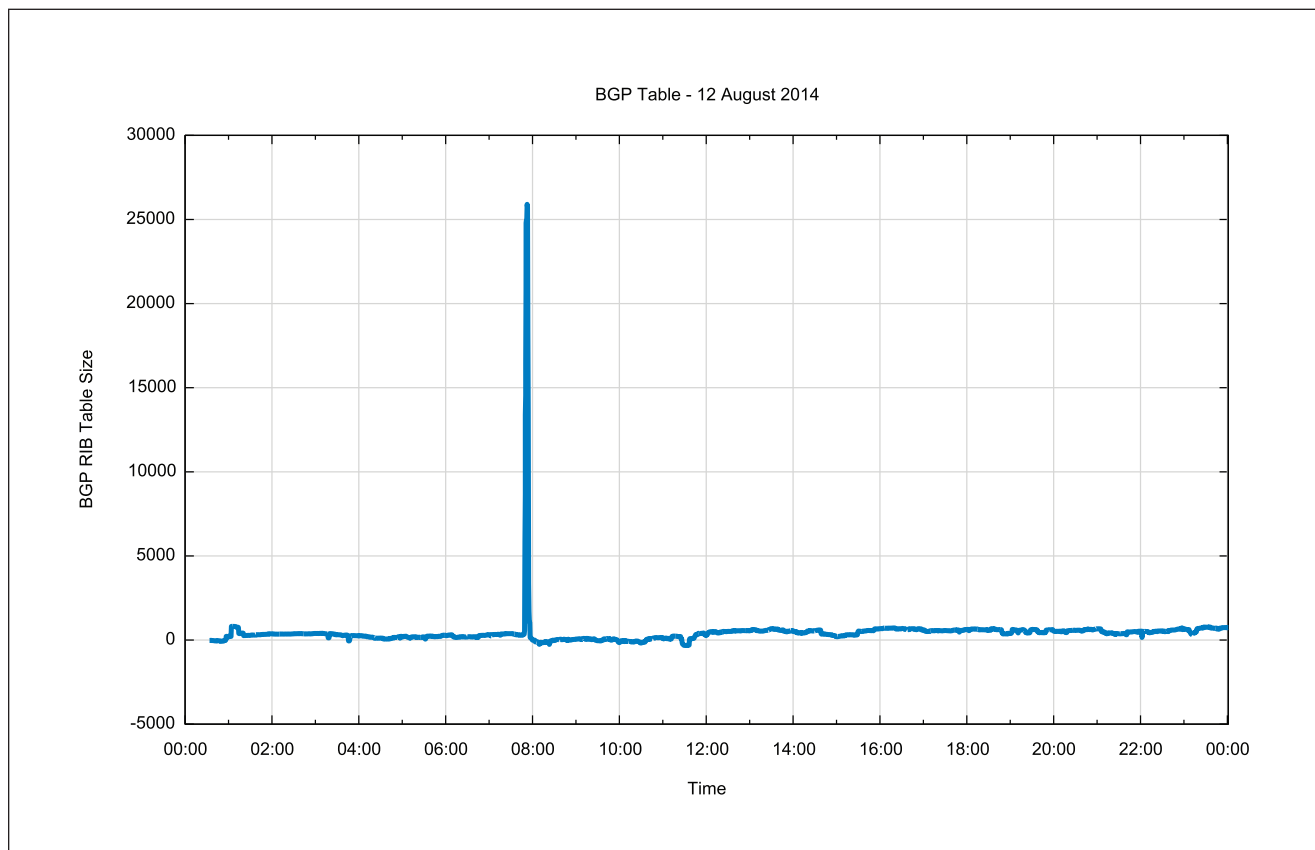
When a network can no longer reach a network address prefix, it needs to inform its peers that it is no longer a useful transit to reach these addresses. The network informs its peers by passing a *withdraw* message to them, withdrawing the network address prefix that it can no longer reach. When a BGP speaker receives a withdraw message, it checks to see if any other peers are still announcing reachability to that address. If so, it adjusts its internal record of the preferred path to reach those addresses to be via a peer that has not withdrawn its path to those addresses, and updates its peers with this new path. If no alternate path to the address exists, the BGP speaker marks those addresses as unreachable and passes a withdrawal message to its BGP peers in turn.

Given that there are some 520,000 network address prefixes and 50,000 component networks, this whole process sounds extremely chatty in terms of protocol interaction. However, BGP uses the *Transmission Control Protocol* (TCP) as its transport protocol, and because TCP provides reliable carriage services, BGP does not repeat its announcements or withdrawals. A BGP protocol session carries only the changes that occur in reachability to network address prefixes. Secondly, BGP uses internal timers to damp the frequency of updates sent to each peer, so each BGP speaker waits for its peers to stabilize before propagating further reachability changes for each network address prefix. The result is surprising stability most of the time. But, from time to time, exceptional events occur.

August 12, 2014

Analysis of BGP activity on August 12, 2014, in terms of the net change in size of the routing table from midnight of that day, shows that something certainly happened just before 08:00 *Coordinated Universal Time* (UTC) (Figure 1).

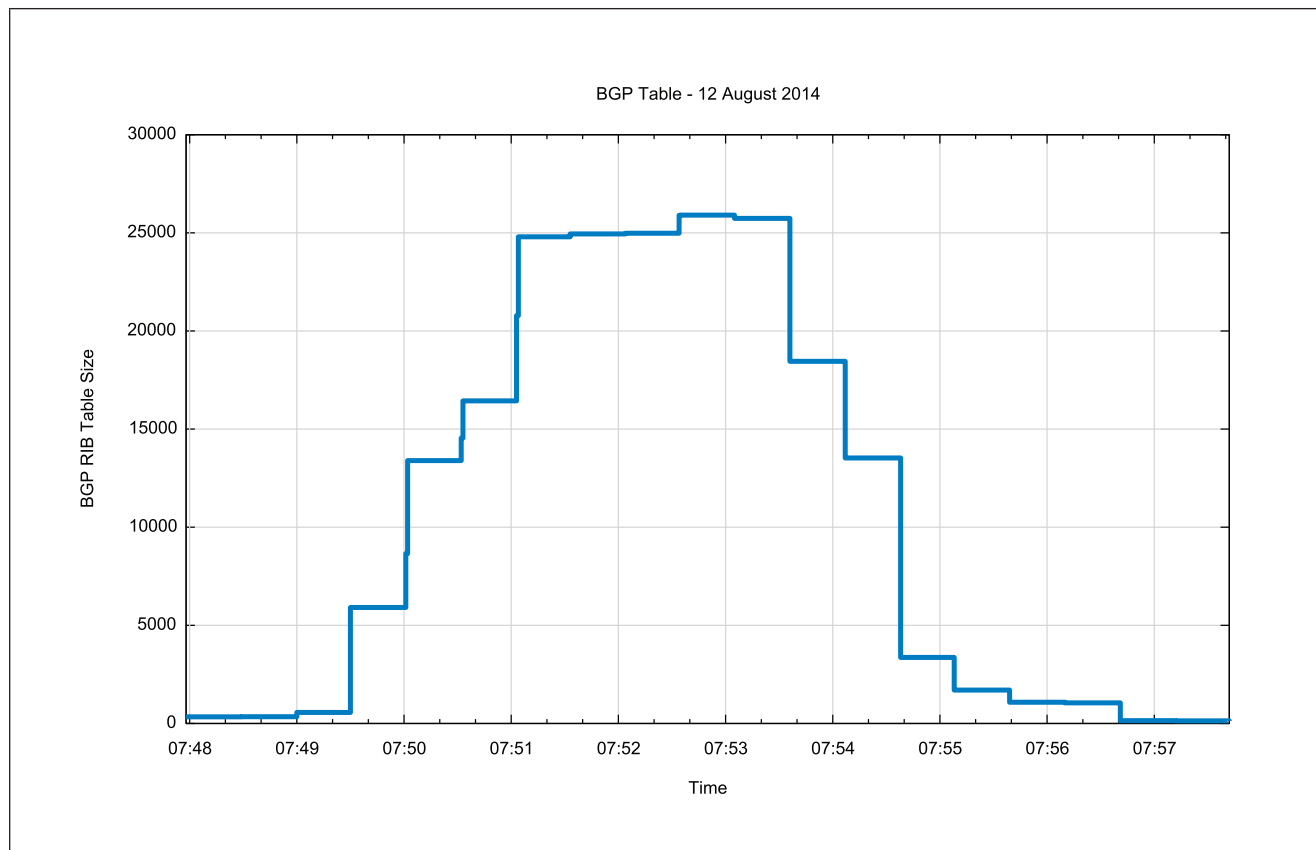
Figure 1: BGP Table Size for August 12, 2014



What dominates this picture is the spike that occurred a few minutes before 08:00 UTC on that day, when the Internet was flooded with what appears from the graph to be 26,000 new prefixes, which were withdrawn very rapidly thereafter. All these new routes shared a common origin, *Autonomous System 701* (AS 701). They did not convey any change in routing information, because they were announcements of more specific prefixes of already announced network prefixes. The announcements were short-lived, and were withdrawn soon after their announcement. The most likely explanation of this event was a *route leak*, where routing detail that was internal to this network was inadvertently leaked into the larger inter-AS routing space, either as a result of a filter reset or a BGP community tag failure, or other forms of transient failure within the route control apparatus of the network.

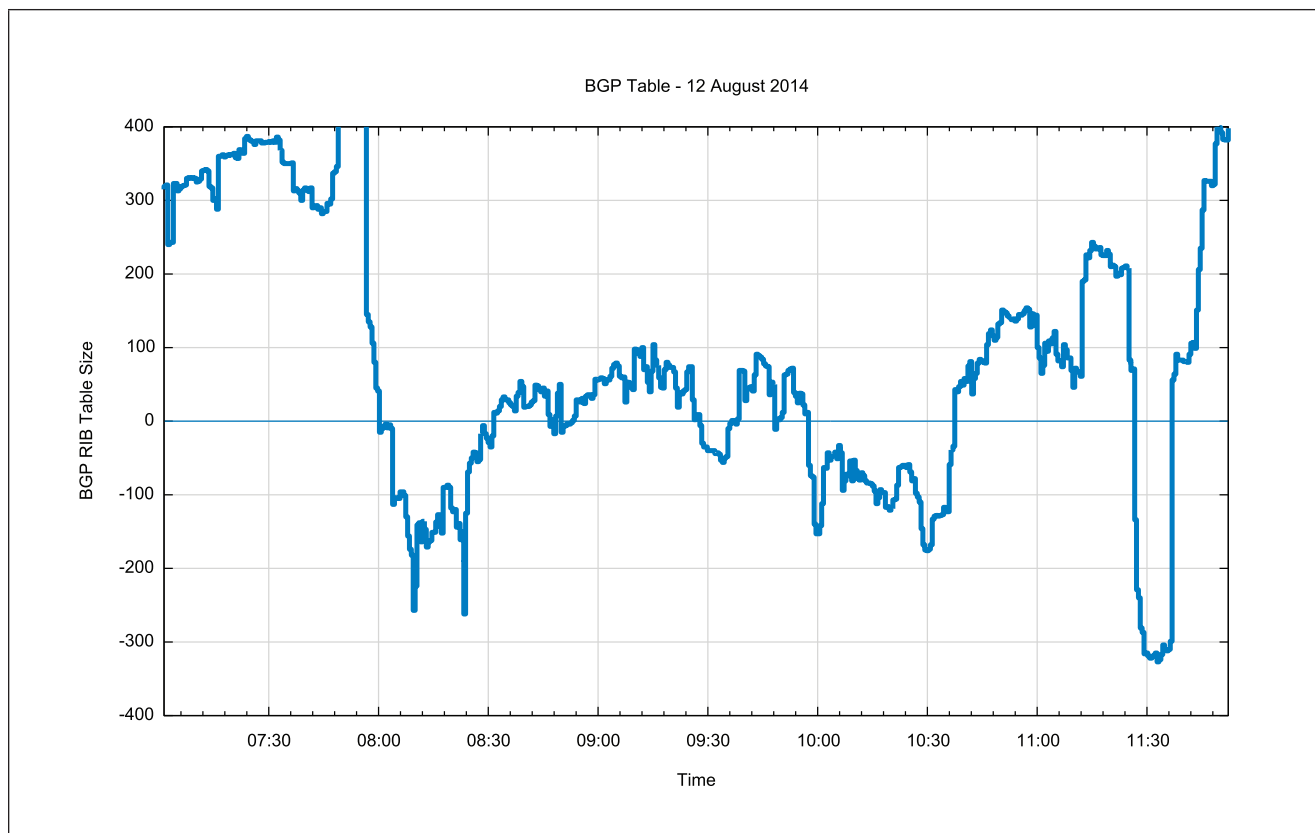
Figure 2 gives us a closer look at that route leak, as seen at AS 131072. An initial burst of 25,000 prefixes was received over a 150-second interval starting 07:49:30 UTC. The leaked routes remained in place for 200 seconds, and then were withdrawn over the ensuing 150 seconds. The routing table returned to its earlier state by 07:57 UTC.

Figure 2: BGP Routing Table: 08:00 August 12, 2014



But while the scale of the data in Figure 2 makes it look as if all returned to “normal” at the end of the route leak, that’s not actually what happened. Figure 3 shows the BGP activity profile up to midday UTC of that day, and it appears that in the immediate aftermath of the route leak a further 550 routes were missing. In the period from midnight until the second immediately prior to the route leak, we had seen a net gain of 300 routes added to the routing table. When the leaked routes were removed, the table size dropped to a net loss of 250 routes for the day. That is, some 600 routes that were present immediately prior to the route leak were missing immediately following the route leak. Many (350) of these routes were subsequently restored over the ensuing 90 minutes, and then a further period of instability involving some 500 routes occurred, until the routing table was at its pre-leak level just before midday.

Figure 3: BGP Routing Table: 07:00–12:00 August 12, 2014



Is this event uncommon? Unfortunately, it’s relatively common. If you look closely at the behavior of the interdomain routing system across any week in the Internet, you will see evidence of some route leaks. If route leaks are so common, then why was the leak on the 12th of August so special?

The first part of the answer to that question concerns the origin of this route leak. AS 701 is one of the so-called “Tier 1” service providers. AS 701 does not purchase upstream transit from any other provider, and, more critically in this context, its route advertisements are, in general, not filtered. Further down in the transit and peering hierarchy the chances of having a filter applied to your advertised routes is high, and the implication is that any form of route leak is quickly suppressed. But if AS 701 is the origin of a route leak, then no other *Internet Service Provider* (ISP) filters the leak, and the advertisements flow across the entire Internet. So the first factor of this leak was that every nondefault BGP speaker was exposed to the route leak.

Secondly, it was a large route leak, in which an additional 26,000 prefixes were added to the Internet routing table for the duration of the leak. While leaks of a few thousand routes are commonplace, they are generally locally contained and appear to be readily absorbed, but 26,000 routes is a somewhat different proposition. It's a significantly large leak.

The third factor is the Internet routing table. At this time many BGP speakers were holding routing tables of around 500,000 routes. There is no single view of the BGP routing system; every BGP speaker gathers its own view, but there is a common core of routes, and at the time of this route leak the common core of advertised routes was around 500,000 routes for most BGP speakers. The leak of 27,000 additional routes on August 12th pushed most BGP speakers to carry in excess of 512,000 routes for a short period of time.

This size of the routing table (512,000 routes) is a default limit point for some commonly deployed items of equipment. The specifications from some commonly used switching equipment have some references to the number 512,000 in the fine print as a default setting for the number of IPv4 entries that are carried in a high-speed lookup cache.

When the number of routes in the routing table exceeds this number, numerous potential scenarios are possible. Note that I am not describing the exact behavior of any particular equipment or configuration here, just the options for failure.

The worst possible option is that the unit crashes and awaits an operator intervention before rebooting; this option may be related to the additional withdrawals that are seen in Figure 3.

Another option is that the condition triggers a reset of the equipment. In the case of the route leak, the reset of the local equipment would take longer than the period of the route leak, and as the equipment came back on line it would be loaded with the “normal” load of some 500,000 routes, and it would function normally once more.

Another possibility is that new and updated routes are simply discarded by the unit in its forwarding caches. This action would result in a rather subtle condition where, for packets addressed to a relatively small number of prefixes, the equipment would silently discard the packet. However, the operating BGP process on the equipment would not necessarily be aware of this action and would report that all was normal.

Something to note about this particular event is that it is more in the way of a warning of what is to come. The continued growth of the routing table is basically inevitable, and by late November 2014 the pool of common BGP routes was passing across the same 512,000 level, and this time it didn't recede, but continued to grow. So in some ways the route leak of August was a warning of what was to follow in a few months.

But before looking at the dynamics of routing-table growth, it's useful to ask why is this particular value—512,000 routes—presents such a problem for some items of routing equipment. And to do that we need to look inside a router.

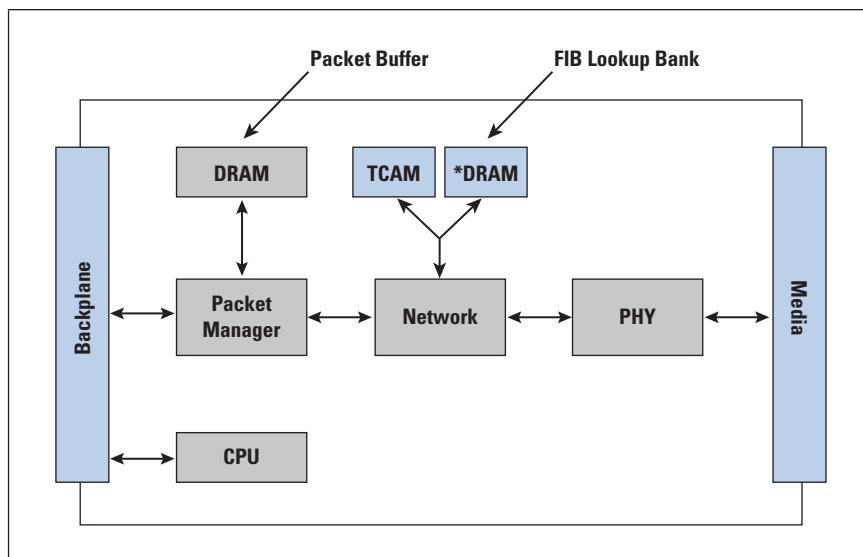
Router Internals

How do these routing-table limits, like 512,000, arise in routing equipment in the first place? Why isn't it possible to design routing equipment that is not arbitrarily limited in this manner?

The internal design of a high-speed router can be described in an analogous way as similar to the old mainframe computer architectures: as a set of specialized modules attached by a common *backplane*. These modules include a *controller*, a *switch fabric card*, and a collection of *line interface cards*.

The purpose of each line interface card is to perform as much as it can in an autonomous manner. It is designed to minimize the load that is imposed on other components of the unit, meaning that each line interface card has many roles to perform; these roles can be summarized as a *Physical Interface*, a *Network Forwarding* unit, and a *Packet Manager* (Figure 4). The Physical Interface unit includes the digital signal processing units that support the interface to the physical media. The interface that this unit presents to the remainder of the line interface card is essentially one of an assembled data packet. For incoming data packets, the network unit performs the initial part of the switching function, where for each received packet the line card looks up a local forwarding table, using the destination parameters from the packet as the lookup key.

Figure 4: Logical Structure of a Line Interface Card



The result of the forwarding-table lookup is the hardware address of the outgoing interface. If this interface is located on the same line card, then the packet is queued to the output structure associated with that local interface. If the interface is located on another card, then the packet is passed to the packet manager for transmission along the backplane to the switching unit to be passed to the outbound line interface card.

A critical aspect of the design of the line interface card is the memory structure used to support the network-level destination-address lookup. This lookup must be completed within the time defined by the maximal packet arrival rate, so for high-speed line cards the performance of this forwarding-table memory structure is critical.

An approach used in some routers is to use *Ternary Content Addressable Memory* (TCAM). TCAMs store a routing prefix in each memory “slot,” using a ternary-state representation of the bits within the prefix (“ternary” because the stored values in the table are either “1,” “0,” or “don’t care”). When presented with an IP address, the TCAM module returns the address of the router interface slot that is the longest match network prefix against the destination address. The advantage of TCAM is that it is consistent, in that every lookup takes just one TCAM cycle time. However, TCAM memory requires a significantly higher gate count per stored bit (up to 24 gates per bit), and the storage structure can be somewhat inefficient, so although TCAM offers consistent performance, it is expensive and consumes a significant level of power on the line card. TCAM is an expensive approach.

An alternate approach is a *trie* (a radix-tree lookup structure) lookup using conventional memory and an *Application-Specific Integrated Circuit* (ASIC) front end. The advantage of this approach is that the routing table can be stored in conventional high-speed *Dynamic Random-Access Memory* (DRAM), which is much cheaper than TCAM, but it does require an ASIC front end. The lookup function also requires multiple comparisons, and the number of comparisons to complete an address search is variable, so this approach does not provide an answer within a consistent time interval. In general, the larger the overall table, the slower the lookup, but the exact performance of a trie depends on the distribution of prefixes in the table, the choice of trie structure, and the specific lookup algorithm that is built into the ASIC.

The question when designing a line card is how much lookup memory should be provisioned on the card, how fast the memory should be, and whether to use a TCAM or a trie structure. The larger the memory and the faster the lookup, the higher the cost, so a trade-off is made between provisioning enough memory and fast enough memory for the expected service life of the unit and at the same time avoiding the cost of overprovisioning.

Two important questions must be answered when looking at this aspect of router design. How quickly will the routing tables grow in the coming years? And how quickly will transmission speeds grow? The answer to the first will influence the size of the forwarding tables in the line interface cards, and the answer to the second will influence the desired memory cycle time.

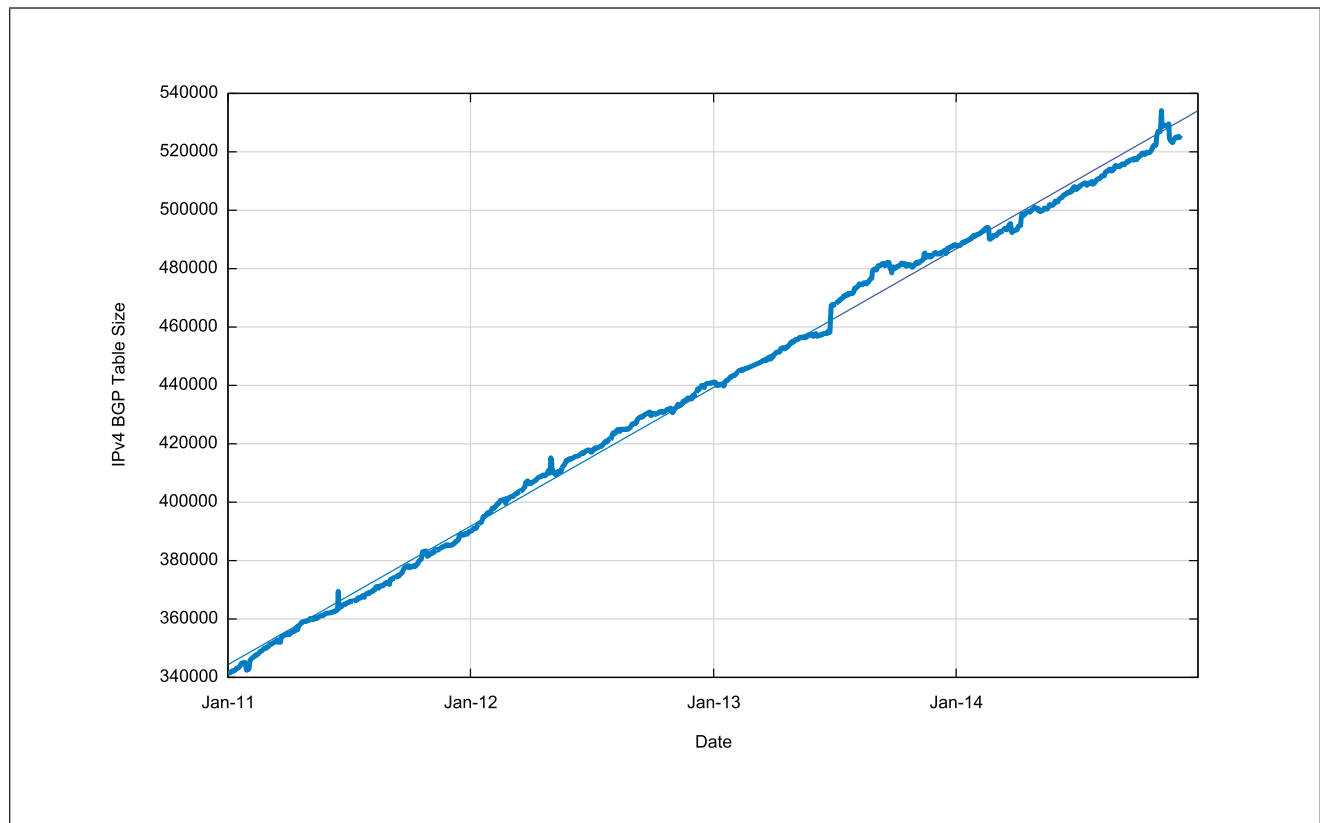
Let's look at these two questions in further detail.

Predicting Routing-Table Growth

Sometimes these table-overflow events are unpredictable, and the route leak on the morning of August 12, 2014, certainly falls into the category of an unpredicted event. But the subsequent growth of the common pool of advertised routes is a predictable event. How quickly is the routing table growing?

Since January 2011 the Internet routing table has increased from some 355,000 entries to the current (late November 2014) level of some 523,000 entries. As can be seen in Figure 5, the overall trend of growth in the past 3 years is that of constant, or linear, growth. What is perhaps anomalous here is that during this same period three of the five *Regional Internet Registries* (RIRs) exhausted their general use pools of IPv4 addresses, yet the momentum of growth in the routing table was largely unaffected by these events. We saw neither a massive change to a large number of more specific advertisements being added to the routing table nor a marked decline in the number of new prefixes appearing.

Figure 5: IPv4 BGP Table Growth 2011–2014



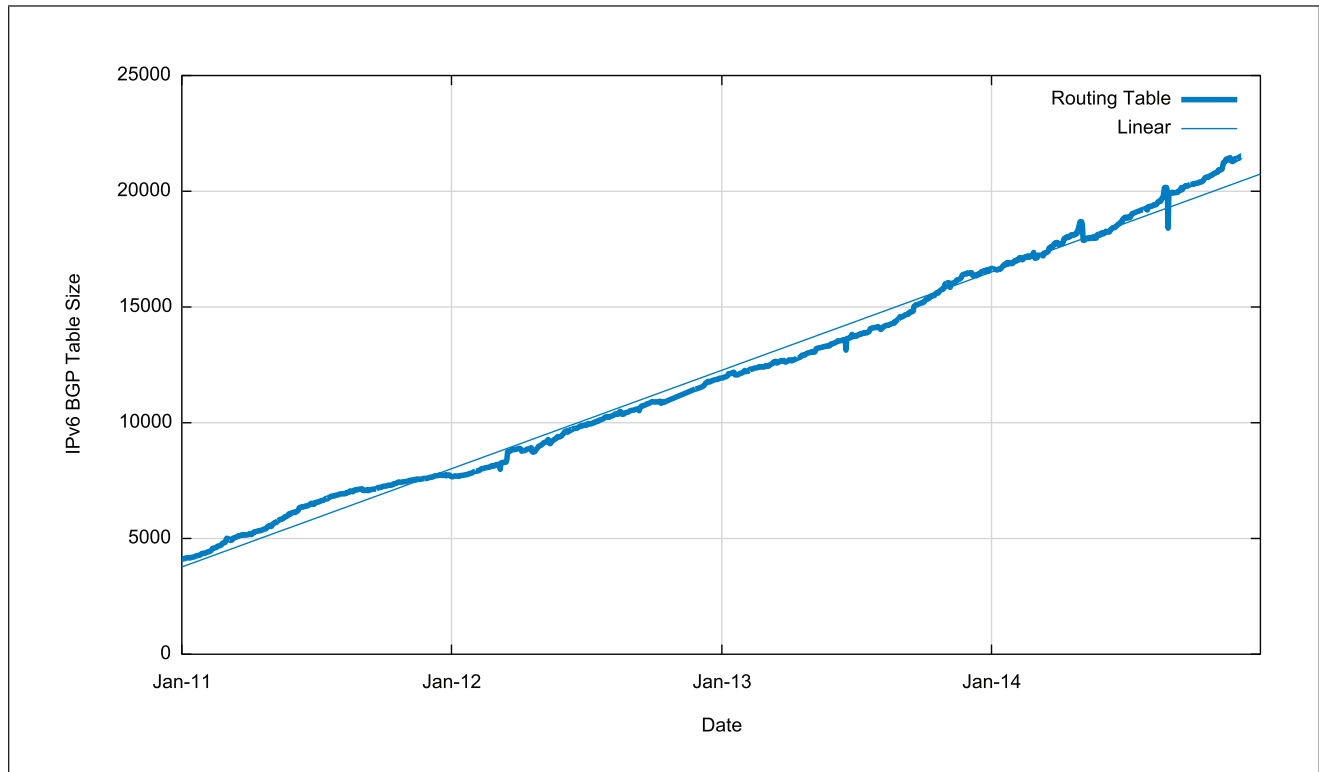
The overall metrics of Internet IPv4 routing-table growth in this period are a modest level of between 8 to 9% per year for most of the basic routing metrics (Table 1). Of course this point is the one where address exhaustion is a little more visible, and the growth in the total span of addresses has grown at a far lower rate of just 2%.

Table 1: IPv4 Routing Metrics 2013–2014

IPv4 Routing Table	Jan-13	Nov-14	Relative Growth (p.a.)
Prefix Count	440,000	523,000	9%
Root Prefixes	216,000	255,000	9%
More Specifics	224,000	268,000	9%
Address Span	156/8s	162/8s	2%
AS Count	43,000	49,000	7%
Transit ASs	6,100	7,200	9%
Stub ASs	36,900	41,800	7%

These days any consideration of the overall routing environment must also include consideration of the IPv6 network. Since the start of 2010 the IPv6 routing table has expanded fivefold, from some 4,000 entries to more than 20,000 entries at the end of 2014. However, this growth has also been predominately a linear growth since 2011, with the table size growing by approximately 4,000 entries per year over this period (Figure 6).

Figure 6: IPv6 BGP Table Growth 2011–2014



The overall metrics of growth in the IPv6 routing table since January 2013 are shown in Table 2.

Table 2: IPv6 Routing Metrics 2013–2014

IPv6 Routing Table	Jan-13	Nov-14	Relative Growth (p.a.)
Prefix Count	11,500	20,580	31%
Root Prefixes	8,451	14,030	27%
More Specifics	3,049	6,550	41%
Address Span	65,127	73,936	7%
AS Count	6,560	9,038	17%
Transit ASs	1,260	1,728	17%
Stub ASs	5,300	7,300	17%

Over this period, when the IPv4 network added a further 172,000 routing entries, the IPv6 network grew at a somewhat more modest level, at least in absolute terms. The number of routing entries grew from 11,500 to 20,500 routes, adding an additional 9,000 prefixes over this period. However, in relative terms this growth represents an annual growth rate of some 31%, which is considerably higher than the equivalent metric in IPv4.

Since 2011 the average growth of routing entries in the routing table has been relatively consistent at a long-term average of some 140 net additional entries per day. In relative terms this growth represents a steady decline in relative growth, falling from a relative growth rate of some 15% per year in 2011 to around 9% by the third quarter of 2014. This slowing down of growth in the IPv4 network could be attributed to market saturation factors in many markets in the developed world, or possibly due to the exhaustion of IPv4 addresses, which has pushed much of the growth activity behind various forms of *Network Address Translators* (NATs). What these figures indicate is that the outlook for IPv4 table growth would be best modeled on a simple linear model, looking at a medium-term growth rate of some 50,000 additional entries per year. This model implies a prediction of the IPv4 routing table reaching some 750,000 entries 5 years from now, in 2019.

However, it must be stated that this model *is* just a model, and it assumes continuity of the environment that accelerates routing-table growth, and of course continuity is simply not going to happen. In what I could describe as the most rational direction for the Internet, the momentum of IPv6 adoption should gather pace, and at some stage within this 5-year outlook, there will be a critical mass of IPv6 deployment such that an IPv6-only end client will have a seamless experience when using the Internet. At that point the momentum behind further IPv4 growth should taper off, and then we will see the IPv4 network shrink as IPv6 assumes the role of the protocol platform for further growth of the Internet. But such a rational perspective of the medium-term future has been constant over the past 5 years at least, and yet it has not eventuated so far. We have to recognize the possibility that we will continue to use IPv4 over the coming 5 years, and absorb the growth pressures through more efficient use of addresses. This paradigm would imply increasing the pressures in address sharing in NATs, looking at ways to intensify the use of public address pools across larger populations of served clients, but may also imply the emergence of fine-grained routing advertisements.

The current convention of a minimum advertised routing prefix size in the default-free zone of the Internet of /24 routes is indeed a common convention across network operators, and it is conceivable that the increasing address scarcity pressures may alter this convention. If we move to an Internet that can support the common acceptance of /25 routes, and even /32 routes, the predictions of the resultant routing-table size are of course far more uncertain.

The growth rates for the IPv6 routing table have increased from an early rate of less than 1 entry per day in 2006 to an average rate of some 17 new entries per day at present, with admittedly a high rate of variance. In relative terms, when this growth is expressed as a proportion of the routing table, the growth rate is slowing down, and the current relative growth rate is somewhere between 20 and 40% p.a. for IPv6.

Within the obvious bounds of uncertainty that accompany any such predictions, these numbers are not particularly alarming in terms of requirements for routing hardware. The routing table is stored in a memory structure, and the capacity and price of memory are related to the number of gates that can be placed into a single integrated circuit. So far *Moore's Law*, postulated some 50 years ago, continues to hold sway, and the silicon industry has been able to double the number of gates on an integrated circuit chip every 18 months or so. If the routing space were growing at a faster rate than this, then there may be some cause for concern about the future cost-effectiveness of routers, but in the IPv4 network it is simply not happening. In IPv4 the linear growth model is far lower than the exponential growth model of Moore's Law, so there is little cause for concern in that domain.

For IPv6 the numbers are a little closer to Moore's Law; if we take a model of the IPv6 routing table doubling in size every 2 years, then the IPv6 routing table is growing at a comparable pace. The mitigating factor here is that the absolute size of the IPv6 table is relatively small, and a 5-year growth outlook from 20,000 entries to some 120,000 entries is not an overly concerning prospect.

Predicting BGP Routing Update Growth

Are there other aspects of the growth of the routing system that we should be concerned about? The BGP protocol is a distance vector protocol, and a common weakness of such protocols is that the protocol reaches convergence by a process of repeated iteration of communication of updates between peer BGP speakers. Each time a BGP speaker receives information of a better path to a destination, it passes this updated information to each of its other peers.

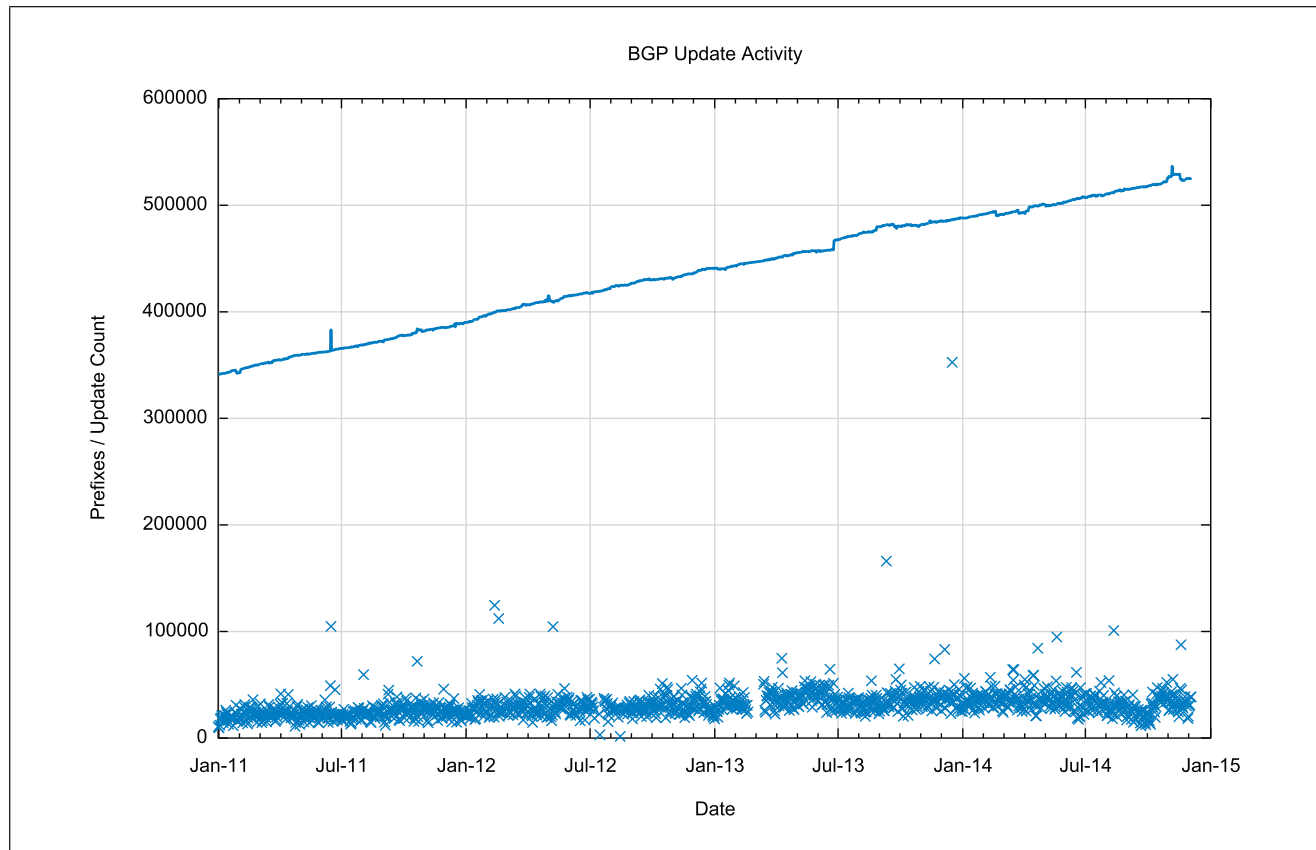
One would expect that as the number of routed entities increases, and as the number of Autonomous Systems (BGP "networks") increases, then the number of updates in BGP should increase at a comparable rate. Of course BGP has many attributes that damp this growth in updates, including the use of TCP as a transport protocol, which removes the need for periodic flooding updates between peers; the use of a *Minimum Route Advertisement Interval* (MRAI) timer, which damps the update rate between BGP speakers; and the use of the AS Path attribute, which prevents the "count to infinity" problem of conventional distance vector routing protocols.

However, these measures should not prevent any growth in the number of BGP updates. At best, they might mitigate such growth, but one would expect that, over time as the Internet grows, the amount of bandwidth and processor capacity devoted to routing should increase as the size of the Internet increases. Over time routers should need faster processors and higher bandwidth to support the operation of BGP. At the same time a larger network with fixed protocol-defined timers should take more time to converge to a stable state. So we should expect to see an increase in the update message counts of the protocol for each BGP speaker and extended convergence times as the Internet grows.

What do we see?

Nothing visible in the observed data supports these expectations. Over the past 4 years the number of entries in the IPv4 routing table has risen from 330,000 to 520,000 entries, yet the number of prefix updates in BGP has remained constant at some 40,000 prefix updates per day (Figure 7). The number of prefix withdrawals was relatively constant, averaging some 40,000 prefix updates per day. In terms of protocol performance, the average time to converge has remained relatively constant at some 70 seconds across this period.

Figure 7: BGP Daily Update Activity for IPv4



The major reason for the observed behavior varying so greatly from orthodox expectations of distance vector routing protocol behavior lies in the overall profile of the inter-AS topology of the Internet. As the number of component networks increases, the new networks all try to cluster towards the “core” of the Internet, and try to avoid attaching to the periphery. The result is an Internet that, as it grows, becomes denser rather than larger, and this increasing density assists BGP to scale.

The efforts with local peering, local exchange points, and large-scale multinational transit providers all assist in absorbing growth without increasing the “diameter” of the Internet, and these efforts offer a direct benefit in preserving the performance of the routing protocol itself.

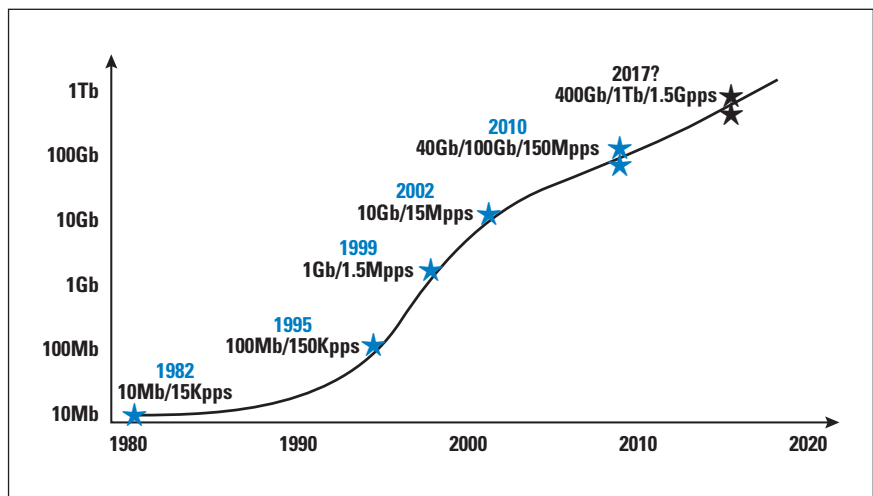
The major conclusion here is that the dynamic growth of updates is also not a cause for significant concern at this time. As long as further growth of the Internet is expressed in terms of increasing the density of the network, and as long as prefix announcements are relatively stable, the operation of the BGP routing protocol will not place extraordinary processor demands on routing equipment.

Predicting Speed

The other important parameter in terms of routing hardware is speed. The router should be capable of processing each packet, implying that in the worst case the amount of processing time available is equal to the time taken for the shortest possible packet to arrive.

In the original 10-Mbps Ethernet specification, the minimum packet size is 64 bytes, and the interpacket gap and frame preamble accounted for a further 20 bytes, implying that for the 10-Mbps Ethernet carrier, the maximum packet rate is 14,880 packets per second (pps), or one packet every 67 microseconds. Since the original 10-Mbps Ethernet specification was standardized in the early 1980s, the speed of transmission systems has increased dramatically. The pace of change in Ethernet speeds is shown in Figure 8.

Figure 8: Ethernet Speeds



Across this evolution of carrier speed, the basic unit of the minimum packet size has remained constant at 64 bytes, implying that for today's 100-Gbps systems the maximal packet rate is some 150 million packets per second, and the interpacket arrival interval is now 6.7 nanoseconds (ns). Taking this thought further with the anticipated 1-Tbps Ethernet specification, the interpacket arrival interval would be cut by a further factor of 10, to 0.67 ns.

Beyond the Tbps threshold, predictions of carriage speeds are difficult to make. Between 1995 and 2002 we saw the carriage speed rise from 10 Mbps to 10 Gbps, a thousandfold increase in 7 years. But a further 8 years elapsed until the standardization of the 100-Gbps system in 2010. To some degree we expect to see a 1-Tbps standard in 2017, but beyond that there is no clear consensus on where and how any further speed increases may be realized.

Router Limitations

What if you wanted to purchase a router today and wanted it to have a production lifetime of many years? What are the basic specifications that such a unit would need to meet in order to address the anticipated demands of, say, a 5-year service life routing the public Internet using BGP?

The processor speeds are not a major issue in terms of processing BGP routing updates. It appears that the push from network operators to maximize connectivity has a positive feedback in terms of limiting the growth of network updates, and the processing capability required to keep pace with today's BGP would not be significantly different from that required in 10 years. The processing capability required in today's routers is not going to vary significantly in the coming years.

However, it's not the same story in terms of the forwarding-table size of the line cards. At the start of 2014 a TCAM with capacity for 512,000 IPv4 entries and 25,000 IPv6 entries would have still been adequate, but now at the end of the year these numbers are inadequate. Ten years is possibly an adequate amount of time to see the transition to IPv6 through to completion in an optimistic scenario, in which case it may no longer be necessary to provide any residual IPv4 support. But this transition has so far taken longer than anyone predicted even 10 years ago, so perhaps in terms of estimating future needs for routing equipment, we should take a more conservative outlook. That conservative outlook would see further fragmentation of the IPv4 address space, and that pessimistic scenario would see the IPv4 routing table approach 1 million entries in late 2019. In addition, we need to include consideration of the IPv6 forwarding table. Assuming some form of momentum behind continued uptake of IPv6 in the coming years, we can anticipate that the IPv6 routing table will grow to some 125,000 entries by 2019.

Beyond that it's more challenging to predict. If we predicted that we would continue to use fine-grained routing control to perform traffic engineering, and use prefix blocks for network policy discrimination, then we could anticipate that the level of routing fragmentation in IPv6 would rise to the same levels we see in IPv4 today. If that's the case, then at the 10-year point we can anticipate an IPv6 routing table of some 512,000 entries.

While Moore's Law talks about the number of gates in an integrated circuit, it does not make the same prodigious predictions over the speed of the chip clock, and clock speeds certainly have not doubled every 1 or 2 years. The fastest available commodity DRAM uses a clock cycle time of between 40 and 50 ns, which is far too slow for 100 Gbps, let alone 1 Tbps. Router memory uses specialist high-speed memory, such as *Double Data Rate Type Three Synchronous Dynamic Random-Access Memory* (DDR3DRAM) and *Reduced-Latency Dynamic Random Access Memory* (RLDRAM), which have clock speeds of up to 9 and 1.9 ns, respectively. This speed is comparable to a 100-Gbps transmission system, which is the form of memory used in today's routers.

If we want this router to survive a production lifetime of 5 years, then the line speeds present a challenge. If the network sits on 100-Gbps transmission systems over this period, then current state-of-the-art high-speed memory would be adequate, but that's a rather unrealistic expectation. Within this 5-year span we will most likely see the emergence of 1-Tbps transmission systems, and if that happens then we will have to improve the clock speeds both in memory and in the line-card packet-processing engines to operate at sub-nanosecond clock speeds. I suspect that this clock speed issue may be the harder challenge and may call for the more imaginative solutions in router design in the continuing effort to meet the demands of an ever-growing Internet.

For production processors the clock rate has remained relatively constant for the past decade. The state-of-the-art in 2002 was a 3-GHz processor, and it has increased only to a 5-GHz processor today. In the computing world the quest for ever-faster computers quickly turned from a quest for faster clock speeds across a giant monolithic system into a quest for ever-larger amounts of parallelism. That way the computer industry was able to meet escalating demands for processing capability and throughput without resorting to exotic technologies in order to support extremely high clock speeds. If this story is less than reassuring, the picture with memory speeds is no better. High memory speeds are achieved through pipelining of memory access requests, rather than in a basic increase in the clock rate.

The Internet may be on the verge of a similar threshold in the design of transmission and switching systems. To date the effort has been largely one of increasing clock speeds in what is essentially a serial paradigm.

BGP is a single-best-path selection routing protocol, and efforts to introduce serialism, such as in equal-cost multipath selection or other forms of dispersed traffic across multiple paths in parallel, have not proved to be all that robust in an inter-AS routing environment. But we can't rely on turning up the clock speed indefinitely.

At some point we may need to take some of the intra-AS approaches to traffic management across parallel paths, using various forms of path pinning, segment routing, and multipath routing, and apply it to the inter-AS routing space, so that we would be looking at further speed increases through the explicit approach of parallelism.

Of course there is also "Plan B." If we really want to reduce the maximal packet rate on high-speed transmission systems, we can always contemplate lifting the minimum packet size. If the minimum packet size had kept itself in proportion to carriage speed, a 64-byte minimum packet on a 10-Mbps system would be a 64-kilobyte minimum packet on a 10-Gbps system, and a 1.2-megabyte packet on a 1-Tbps system. Lifting the minimum packet size to 1.2 megabytes on very-high-speed systems is perhaps heading too far, but when we contemplate these 1-Tbps systems, then perhaps we should reserve some time to think about speed and capability and whether it's time to revise the minimum packet sizes on these ultra-high-speed systems.

Either way, while the next 5 years of Internet growth can be predicted within some acceptable levels of uncertainty, trying to push this range of visibility out to 10 years is a tough task. The continual pressures of scale and speed don't look as if they are stopping anytime soon, so no doubt sometime in the future we will encounter more Internet "bad hair days," as deployed equipment trips over further basic limitations in their size and speed in the face of the inexorable continuing growth of the Internet.

For Further Reading

- [0] Yakov Rekhter, Susan Hares, and Tony Li, "A Border Gateway Protocol 4 (BGP-4)," RFC 4271, January 2006.
- [1] Geoff Huston, "The BGP Routing Table," *The Internet Protocol Journal*, Volume 4, No. 1, March 2001.
- [2] Kris Foster, "BGP Communities," *The Internet Protocol Journal*, Volume 6, No. 2, June 2003.
- [3] "BGP: the Border Gateway Protocol Advanced Internet Routing Resources," www.bgp4.as

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. E-mail: gih@apnic.net

The ZigBee IP Protocol Stack

by Douglas Comer, Purdue University

A major area of networking is emerging: communication among intelligent embedded systems. The idea is that computing systems can be embedded in many devices, and the systems will be able to use the Internet to communicate with other devices. Of course, some embedded devices will provide information that humans view, and humans may run applications that control embedded devices. Consequently, the new networking paradigm does include some human interaction. However, the major emphasis is on systems that can interact with their environment and with each other, rather than on conventional computers that store data and run applications. Researchers and professionals have coined the terms *Internet of Things* (IoT, or iThings) and *Machine-to-Machine* (M2M) *Applications* to capture the idea. Despite sounding awkward, the phrase “Internet of Things” appears to have become widely accepted.

This article begins by presenting examples of intelligent embedded systems and the way they communicate. It then discusses one technology in detail: a protocol stack for wireless mesh networks that is designed for smart-grid applications. The article assesses the protocol stack, the use of IPv6, and some of the consequences of the design.

Sensing, Monitoring, and Control Applications

We use the term *embedded system* to refer to a computational system that is an integral part of another mechanism or device. The chief difference between an embedded system and a conventional computer system arises from their external connections. A conventional computer system deals with information: the computer can store, access, and manipulate data. In contrast, an embedded system can *sense* and *control* the physical world around it.

As an example, consider a thermostat used to control a heating and air conditioning system. A modern thermostat (called a *smart thermostat*) constitutes an embedded system: the thermostat contains an embedded processor that runs software to perform all functions. A user can configure the thermostat to change settings according to the time of day. The thermostat has connections to a variety of sensors that might include an indoor temperature sensor, an outdoor temperature sensor, a sensor that detects airflow (that is, whether the fan is operating), and a sensor connected to push buttons that allow a user to set the desired temperature. More advanced systems have sensors for the relative humidity of the air. In addition, the thermostat has connections to controls that allow the processor to turn the heater or air conditioner on or off, regulate the speed of the fan, and control humidifier and dehumidifier functions.

Most computer users are already familiar with embedded systems. For example, a printer connected to a computer incorporates an embedded system. When a computer sends a document to the printer, the embedded system of the printer controls the motors and mechanisms in the printer, causing them to feed sheets of paper through the printer, move the ink jet mechanism, and spray drops of ink. The printer also contains sensors that can detect a paper jam or low ink supplies.

General Electric (GE), the largest industrial company in the United States, is going beyond items for the consumer market. GE produces aircraft engines, power-plant turbines, locomotives for railroads, medical imaging equipment, and heavy-duty machinery that transports people, heats homes, and powers factories. Using the phrase *Industrial Internet*, GE has launched a major initiative to incorporate communicating embedded systems in both its factories and its products.

Power Conservation and Energy Harvesting

Some embedded systems, such as the embedded control system in a printer, attach to a reliable source of continual power. However, many embedded systems rely on temporary power, and are designed to conserve energy. For example, cell phones run on batteries and environmental sensors located in remote locations (for example, a desert) may use photocells.

As a special case, some embedded systems are designed to *harvest energy* from the environment around them. For example, a sensor in the ocean might use the motion of waves to generate power, and a sensor near a hot spring might use thermal energy. Energy harvesting even includes the kinetic energy that humans generate merely by opening a door or flipping a light switch. An embedded system that uses harvested energy may need to operate periodically—the system might need to accumulate energy until a sufficient charge is available (for example, to run a radio transmitter).

A World of Intelligent Embedded Devices

To understand the vision for the Internet of Things, we have to imagine that powerful embedded systems will be everywhere: houses, office buildings, vehicles, shopping malls, and street corners. For example, consider vehicles. In addition to systems that provide entertainment and navigation, designers envision systems that allow a vehicle to calculate the distances to surrounding vehicles, sense objects in the roadway, warn of changes in pavement (for example, from construction), and sense lanes and warn the driver when a car drifts. A vehicle with intelligent embedded systems can communicate with nearby vehicles, and can coordinate braking. An intelligent embedded system can use facial recognition to identify drivers when they enter a vehicle, adjust settings to their preferences, monitor an individual driver's driving habits and watch for unusual driving, and adjust warnings to accommodate a specific driver's reaction times.

In office buildings, embedded systems already sense the presence of individuals and adjust lights and heating or cooling accordingly. A system can use sensors to change the heating and cooling systems when windows are open. More important, an intelligent embedded system will be able to use learning algorithms to accumulate patterns. For example, if given employees move in and out of their office frequently during the work day, the system can learn not to turn off the heat until they are absent for a longer time. Similarly, if an employee tends to work late, the system can learn the pattern and control the office environment accordingly. Thus, if the employee usually arrives early and goes home at 3 p.m. each day, an intelligent system can learn the pattern and anticipate the employee's arrival and departure.

The Importance of Communication

Why is emphasis shifting to intelligent embedded systems that can communicate? There are many advantages. For example, in addition to local coordination, communication allows systems to exploit *Cloud Computing*^[10, 11] to analyze data from a set of embedded devices. Communication means individual embedded systems can have smaller memories and less-powerful processors, meaning they will use less energy. In short, small embedded systems can achieve complex functions by working together with nearby embedded systems or by accessing remote information.

As an example, consider a set of sensors used to assess the stress on civil infrastructure, such as bridges. Measuring stress is important in understanding whether a bridge can tolerate the load, whether reinforcement is warranted, or when a bridge should be replaced. To measure stress, engineers place small battery-powered sensors at various points along a bridge. Without communication, each sensor must have a local store to keep measurements along with a timestamp for each. If each sensor includes a radio, the set of sensors can form a wireless network in which measurement data is passed along to a collection point, which may be located across the Internet, far from the bridge. In terms of measurement, the important difference arises from coordination and rapid assessment. Communication allows sensor nodes to run a protocol that takes readings simultaneously. Uploading data in real time makes it possible to detect dangerous situations quickly and take action before a disaster occurs.

Communication can also lower costs. For example, consider *smart meters* used by utility companies. The traditional approach used to assess usage consists of placing a meter outside each customer's location and sending a person to record the meter reading each month. A smart meter incorporates wireless communication, meaning the utility company can read the meter from a remote location. Even if a smart meter uses a wireless transmitter that reaches only to the street, the meter can be read from a passing vehicle rather than an individual on foot, lowering the cost of reading meters dramatically.

Embedded Systems in Shopping Malls

In addition to the traditional sensor systems described previously, the Internet of Things includes unexpected applications. For example, many shopping malls now have large video display panels that show ads. Stores use the displays to advertise products and services as well as discounts and special promotions. Communication is needed to download ads dynamically because the content and schedules can change at any time—management needs the ability to control which ad is displayed on a given screen and how long the ad remains visible.

Where is the control system for the video displays located, and what networking technology is needed to connect the control system to the individual displays? The answer is interesting and a bit surprising: in current implementations, each display contains a TCP/IP protocol stack and a connection to the global Internet. An Internet connection allows the controller to be located anywhere. In particular, the mall can outsource IT functions by placing the controller “in the cloud.” More important, providing each display with an Internet connection and protocol software means that the content does not need to reside on the same physical host as the controller.

Allowing content to be separate from the controller is important because it differentiates information owned by a retailer from the control system for a given mall, and allows a retailer to serve content to multiple malls. For example, a retailer such as Apple can place video content for ads on a cloud server, and then issue commands to the controllers for each mall to specify a schedule of items to display along with the URL of each item. Because they have Internet access, individual display systems can download a copy of an item they have been assigned to show (possibly through a local cache to improve performance).

Uploading Data from the Internet of Things

Displays in shopping malls illustrate another important capability of networked systems: the ability to upload data. Although it is not obvious to customers, some of the displays in shopping malls are equipped with a camera. When people approach the display, the system uses the camera to detect their presence. The system identifies human faces and applies analysis algorithms that use features, such as the distance between the eyes, to characterize the individual. With high accuracy, software in a mall display is able to tell whether the individual in front of the camera is male or female as well as the person’s approximate age group. Thus, instead of merely following a predetermined schedule of ads to display, the system can use the characteristics of the individual to choose an appropriate ad. For example, a middle-aged male might be shown an ad for a sports car instead of an ad for women’s apparel.

In addition to using video information to select ads, mall systems also gather and report data about the interaction. For example, a system uploads statistics about how many people watched a given ad, their sex and age, and how long each person or group remained in front of the display during a particular set of ads. Grocery stores are using the same approach: they are installing cameras with embedded processing capability over freezers and at other locations in the store. The cameras record whether customers stop at a given product display, how long each person looks, and how many customers finally select a product or merely move on. When it is uploaded to a server, the information from a given location can be combined with information from other locations. The key idea is that combining data from multiple sites increases the accuracy of the analysis.

Wireless Networking and IEEE 802.15.4

How should devices be connected to the Internet of Things? Wireless networking technologies are popular, even in the case of semipermanent deployments across a small area. For example, consider the electronic displays in malls. Although they are semipermanent (that is, a display usually remains stationary for weeks), wireless networking means a display can be moved without installing a new network connection.

What wireless networking technologies should be used to connect an intelligent embedded system? The answer depends on several factors, including the geographic distance between nodes of the network, the desired data rates, and the power requirements. Power is important in two ways. In the case of embedded sensor systems that run on battery power, overall power consumption must be minimized to maximize battery life. Although the power used by a radio transmitter can dominate the overall battery drain, using a smaller memory or reducing processor speed can also lower overall power requirements. In the case of embedded systems that have a continuous power source (for example, connect to a wall outlet), it may be necessary to limit radio transmissions to avoid interfering with other devices or other transmissions.

Several low-power wireless networking technologies have been standardized. This article considers a network technology defined by the IEEE standard 802.15.4^[1]. Various versions of 802.15.4 have been produced; they differ in the frequency bands and modulation techniques used as well as the *Maximum Transmission Unit* (MTU). The IEEE standard specifies the physical and *Media Access Control* (MAC) layers of the network, and other groups have defined upper-layer protocols for use on low-power wireless networks. For our purposes, it is necessary to understand only the general characteristics of 802.15.4 technology:

- The data rate is relatively low (a maximum of 250 kbps).
- The MTU is extremely small (127 octets).
- The distance is limited (a maximum of 10 meters with a conventional antenna and power from a battery).

A Mesh Network for Smart-Grid Sensors

One application of 802.15.4 wireless technology arises from an effort to add intelligent computer management to the electrical power-distribution system. Known as *Smart Grid*, the overall plan includes placing sensors in all devices that use electricity. In addition to large systems, such as those used for heating and cooling, the designers envision sensors in kitchen appliances (for example, ovens, refrigerators, dishwashers, and even toasters), computer systems, and entertainment systems (for example, televisions and stereos), and small handheld appliances. The utility companies want to charge more during peak hours, and the sensors will communicate with the utility company to determine pricing and warn users when prices are high. Alternatively, the sensors will be capable of disabling certain uses during peak hours.

The most obvious design for a system of sensors in a residence consists of placing a base station in the residence and using a wireless technology that allows the base station to communicate with each sensor. The approach is known as a *hub-and-spoke* topology. For example Wi-Fi (802.11) systems use a hub-and-spoke approach. However, such a design does not work well for all situations. In particular, metal pipes and other obstructions in a building can interfere with wireless signals and make it impossible to reach all locations from a single point, especially in the case of portable appliances that can be moved from one room to another. Therefore, the smart-grid designers envision an adaptive system in which the set of sensor nodes automatically forms a *self-organized mesh network*, simply called a *mesh*. Each node in the mesh performs two tasks: communication for the device to which it attaches and forwarding for other nodes.

A residence will contain a *border router* that connects the mesh network to the outside world. When a node boots, it joins the mesh and tries to establish a connection to the border router. If it can reach the border router directly, the node communicates directly with the border router. If it cannot reach the border router directly, the node searches for a nearby neighbor node that has a path to the border router. In essence, the neighbor agrees to act like a router and forward packets.

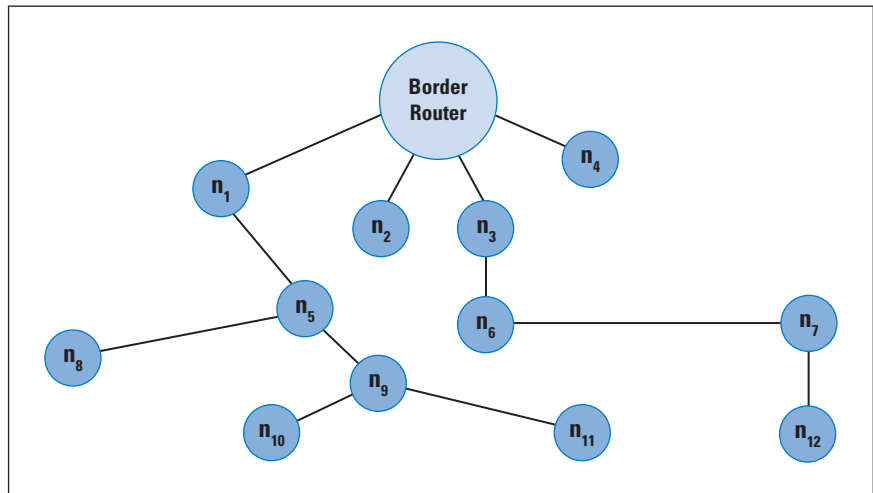
If multiple neighbors have a path to the border router, a node applies a selection algorithm to choose one. The selection algorithm can account for several metrics, including quality of the radio signal (that is, the level of interference), latency, link capacity, and capacity of intermediate forwarding nodes. When a node selects a path to the border router, the node informs its neighbors and agrees to forward their packets. Thus, if it is possible to create a network that gives each node connectivity, the nodes will form a mesh network automatically. We will return to the question of how a neighbor makes a selection later in the article.

A Forwarding Tree for a Mesh Network

Our description makes the selection of paths in a mesh seem trivial. In fact, a mesh may offer many possible paths. The distinction between routing (choosing paths among networks) and forwarding (choosing paths within a network) is blurred because we can think of each radio link as a network. More important, a path with fewest hops may not be optimal because connectivity is pairwise among nodes, and a low signal strength between a given pair may make the path unreliable. Furthermore, if mesh nodes are mobile or changes in the environment interfere with signals (for example, people moving through a room), forwarding must change dynamically. Despite the potential complexity of mesh routing, solutions have been created to handle basic cases. In particular, we will see that protocols have been created to handle forwarding between the border router and individual nodes in a semistatic mesh topology.

The key idea that simplifies mesh forwarding is an observation: if the primary goal is to enable sensors to communicate with a remote server, each sensor node needs to choose only one way to reach the border router. Traffic from a remote server to a mesh node can be directed along the path in the reverse direction. That is, the forwarding paths in the network form a tree in the graph-theoretic sense (that is, a graph with no cycles), with the border router as the root. Figure 1 illustrates a set of sensor nodes and one possible forwarding tree.

Figure 1: Example Forwarding Tree Imposed on a Set of Sensor Nodes. A Tree Results if Each Node Chooses One Path to the Border Router.



As the figure indicates, a border router is usually larger than other nodes in the network (that is, has more processing power and memory). We will see that a border router is expected to run routing and server software that provides forwarding service for the entire mesh.

Using Internet Protocols in a Mesh

In terms of Internet protocols, the question is whether to use the traditional paradigm of assigning an IP prefix to the entire mesh or to use the paradigm of treating each radio link as a separate point-to-point network.

The two approaches are known as:

- *Mesh-under*: The mesh acts like a single network, and Layer 2 protocols handle broadcast and multicast.
- *Route-over*: The mesh acts like a set of point-to-point links, and Layer 3 protocols perform all forwarding.

With the *mesh-under* approach, which is favored by IEEE, a node uses Layer 2 protocols to form a forwarding tree. The idea is similar to the bridge and spanning-tree protocols used for Ethernet. For example, a node uses a Layer 2 broadcast to discover which neighbors are within radio range. Each neighbor responds, allowing the pair to learn that they can reach each other and the signal quality. The border router also uses Layer 2 broadcast to advertise itself. Nodes in the mesh pass along information about which nodes they can reach. Eventually, each node in the mesh is aware of other nodes and how to reach them, as well as how to forward broadcast packets. Thus, given the Layer 2 address of the border router, nodes in the mesh will know how to forward packets (that is, they will have formed a forwarding tree).

With the *route-over* approach, which is favored by the *Internet Engineering Task Force* (IETF), a node uses Layer 3 protocols to identify neighbors and form a forwarding tree. Of course the underlying hardware does not understand IP addresses or the IP datagram format. Thus, Layer 3 packets are carried in Layer 2 frames. For example, to find neighbors, Layer 3 software generates an IPv4 datagram with a local broadcast address or an IPv6 datagram with a link-local multicast address. The datagram is sent via hardware broadcast.

In either approach, when a new node enters the mesh, the node must choose how to link itself into the tree. There are two steps. In the first step, a node must find the set of neighboring nodes that can be reached directly and assess the quality of the radio link to each neighbor. In the second step, the node must choose either to communicate with the border router directly or to use one of the neighbors when forwarding packets. Note that the quality of a signal is of key importance—even if a node can reach the border router directly, the node may choose an indirect path if the signal quality of the direct connection is sufficiently low.

In terms of an IP protocol stack, a major difference between mesh-under and route-over arises in IP forwarding. The mesh-under approach follows the traditional paradigm of treating the entire mesh as a single network with familiar characteristics of a single broadcast domain and the ability for an arbitrary pair of nodes to communicate.

IP assigns a single prefix to the mesh network, and IP forwarding causes any datagram destined for a node on the network to be passed to the underlying hardware interface for delivery. When it accepts an outgoing datagram from IP, the network interface uses the information that has been gathered by Layer 2 routing protocols to choose a next hop to which the datagram will be forwarded. As the packet travels across the mesh, the packet is processed only by Layer 2 on each intermediate node. When it arrives at the destination, the datagram is passed up to Layer 3. Thus, all mesh details are hidden from Layer 3.

The route-over approach makes IP aware of the mesh topology. That is, IP becomes aware that although some nodes on the network are reachable directly, others are not. In particular, using route-over breaks a standard assumption in IP protocols that if two hosts share an IP prefix, the two attach to a network that allows them to exchange packets directly. In a route-over mesh, all nodes in the mesh share a prefix, even though a given node can communicate directly only with nearby neighbors. Using IPv6 terminology, we say that a node is either *on link* or *off link*. To handle nodes that are not directly reachable, IP uses *source routing*. IP must understand the topology of the mesh and be able to specify a path through the mesh to the destination (for example, go to node 9, then to node 5, then to node 1, and finally to the border router). The next sections describe the *ZigBee* protocol stack that uses the route-over approach, and later sections assess some of the problems that arise.

The ZigBee IPv6 Protocol Stack

The ZigBee Alliance^[2] and the IETF^[3] are cooperating to define the use of IPv6 in a mesh design that follows the route-over approach. The ZigBee Alliance has defined an open standard known as *ZigBee IP*^[4]. Table 1 lists three key IETF working groups related to the ZigBee effort. The next sections describe the protocols that are being developed.

Table 1: IETF Working Groups Related to the ZigBee Effort

Name	Main Contribution
6LoWPAN	IP-over-802.15.4 Shim Layer
ROLL	RPL – A Routing Protocol for Mesh Networks
CoRE	CoAP – Constrained Application Protocol

The basic idea of the ZigBee route-over design is to use IPv6 when possible and introduce modifications as needed. A protocol has been created to compress IPv6 datagrams and send them over an 802.15.4 radio link. A modified form of IPv6 *Neighbor Discovery* is used to find the IP addresses of directly reachable neighbors, and a protocol has been invented to allow neighbors to exchange characteristics.

In addition, a new routing protocol is used to gather information about connectivity throughout the mesh and compute forwarding information. Finally, an IPv6 source route header is used to forward each datagram hop-by-hop across the mesh. The next sections describe some of the basic protocols.

IPv6 over Low-Power Wireless Networks

The *IPv6 over Low power Wireless Personal Area Networks* (6LoWPAN) effort defines the transmission of IPv6 over a 802.15.4 radio link. The primary problem is a conflict between the IPv6 requirement for an MTU of at least 1280 octets^[5] and the 802.15.4 protocol that species a maximum of 127 octets. In fact, if AES-CCM-128 encryption is used, the available payload size is reduced to 81 octets. To send an IPv6 datagram over such a link, 6LoWPAN introduces an extra shim layer that performs compression and transmission. The shim layer accepts an outgoing IPv6 datagram, compresses the header, divides the datagram into a series of pieces we will call *fraglets*, and sends each fraglet in a separate packet. On the receiving side, the 6LoWPAN shim layer accepts incoming fraglets, recombines them into a single datagram, decompresses the header, and passes the result to the IP layer. Thus, IPv6 is configured to send and receive complete datagrams without knowing that the shim layer breaks a datagram into fraglets for transmission.

There are two reasons 6LoWPAN does not use conventional IPv6 fragmentation. First, 6LoWPAN needs to operate over only a single link. Thus, the protocol is much simpler because all the fraglets of a datagram must arrive in order. Second, IPv6 fragmentation cannot handle an MTU of 127 octets.

6LoWPAN Neighbor Discovery

Traditional *IPv6 Neighbor Discovery* (IPv6-ND) provides an address-resolution mechanism that can be used for, among other things, *Duplicate Address Detection* (DAD). Unfortunately, IPv6-ND makes a fundamental assumption that an IPv6 prefix maps to a broadcast domain. Therefore, a node can use IPv6 multicast, which maps to hardware broadcast, to reach all other nodes that share the prefix. In a mesh network, however, a broadcast transmission may reach only some of the nodes in the mesh, making some of the nodes off link. As a result, conventional duplicate address detection will not work correctly.

6LoWPAN Neighbor Discovery (6LoWPAN-ND) defines several changes and optimizations of IPv6-ND that are intended specifically for lossy, low-power wireless networks that have limited range. In general, 6LoWPAN-ND avoids all mechanisms that flood packets across the mesh. Instead of requiring individual nodes to engage in duplicate address detection, 6LoWPAN-ND uses a *registration* approach in which each node in the mesh registers its address with software that runs on the border router.

As nodes register their addresses, software on the border router flags any duplicates (recall that a border router has the processing power and memory needed to handle networkwide services, such as address registration). Finally, 6LoWPAN-ND allows nodes to *sleep* (that is, go into a stasis state to conserve power). When a node reawakens, it must renew its address registration in case some other node registered a duplicate address during the sleep period.

Mesh Link Establishment

IPv6 is designed with the assumption that underlying hardware links have been configured before IP software runs. In particular, IPv6 expects exchanges to be authenticated, meaning that links must already be in place. For an 802.15.4 mesh, a node must choose how to link into the forwarding tree. Link configuration is complex because radio transmissions can be asymmetric. A node cannot merely listen for transmissions from neighbors and choose the neighbor with the strongest signal, because the question is how well a neighbor can receive transmission of a node. Consider the case of a border router with a powerful transmitter and large antenna. It may be possible for a node to receive a strong signal from the border router, even if the node transmitter is too weak to reach the border router. Thus, before IPv6 can be used in a route-over mesh, a lower-level protocol is needed that allows a node to learn the level of signals that neighbors observe when receiving the transmissions of the node.

ZigBee uses the *Mesh Link Establishment* (MLE) protocol for link configuration. MLE employs a two-way packet exchange: one node transmits a message and a receiving node sends a reply. The reply reports the quality of the signal that was observed. When a node receives an MLE reply from each neighbor, the node will know how well each neighbor can receive its transmissions. Of course, signal strength can change over time if nodes move or electrical interference is introduced (for example, a large electrical motor starts to run). Therefore, the measurements must be repeated periodically.

MLE includes facilities for more than signal-strength assessment. During packet transmission, MLE allows nodes to exchange configuration information. The two nodes exchange address information, and choose a type of security for the link. Most important, MLE allows a node to inform a neighbor that it can reach a border router. Thus, when it joins a mesh, a new node runs MLE, gathers information about which neighbors have a path to the border router, and uses signal strength to choose one of the neighbors as a parent node in the forwarding tree.

Interestingly, not all the ZigBee protocols account for asymmetric signal strength. In particular, when it builds a forwarding tree, the *Routing Protocol for Low-Power and Lossy Networks* (RPL) selects links that are bidirectional; if communication can proceed only from node *A* to node *B* and not from node *B* to node *A*, RPL does not include a link between *A* and *B*.

However, the ability to communicate does not imply that the quality is the same in both directions. In the case of a border router communicating with another node, it seems reasonable to assume that the signal sent by the border router could be stronger than the signal sent by the other node. Even in the case of two nodes that are not border routers, it may be that the signal received in one direction is much stronger than the signal received in the reverse direction. Nevertheless, once it chooses a path, RPL sends traffic in both directions over the path.

Forwarding in a ZigBee Route-Over Mesh

Conventional IP forwarding uses the IP prefix when choosing a next hop. As pointed out previously, however, some nodes of the network will be on link (that is, directly reachable) and others will be off link (that is, reachable only indirectly). Therefore, when deciding how to forward a datagram, IP must use more information than the network prefix. The problem can be handled in two ways, and ZigBee IP uses the terms *storing mode* and *non-storing mode* to characterize the two approaches.

As the term *storing* implies, each node in the network stores a significant amount of information. In addition to storing the address of a parent (that is, the next hop to the border router), each node learns and stores a next-hop address for each node that is downstream. In a graph-theoretic sense, a node in the tree stores a next hop to each node in its subtree. The worst case occurs in networks where a single node, N , connects all other nodes to the border router. For such a case, node N stores a next hop to all other nodes in the network.

When it needs to forward a datagram, the IP software on a node consults the forwarding information that has been stored locally. If the destination is downstream, the information specifies the next hop to use to reach the destination. If the destination is not downstream, the node forwards the datagram along the path toward the border router.

The memory requirements for storing mode are more extensive than the previous description implies. Once RPL has computed routes, a node needs to store only a next hop for destinations in its subtree ($N - 1$ destinations in the worst case). However, additional memory is needed during route computation because RPL uses a link-state algorithm. Therefore, a node must collect pairwise link advertisements for all links in the subtree and then run the shortest-path algorithm to compute next hops. The memory required is still proportional to the number of nodes in the subtree, but the computation can require twice the memory used to store next-hop information.

Although it handles transfer along the edges of a forwarding tree, storing mode does not provide optimal routing in all cases. For example, consider two nodes that lie in close proximity but are not in the same forwarding subtree.

When one sends to the other, the packet is forwarded up the tree toward the border router. If the packet reaches a node that is common to both forwarding subtrees, the packet is sent down the other tree to its destination. The worst case occurs when the border router is the only node in common with both subtrees: the packet must travel all the way up one subtree to the border router before being sent down the other subtree to the destination. The IETF is developing a protocol that will find and use routes that lie outside the forwarding tree.

The concept of a forwarding tree and a routing protocol, such as RPL, that computes and maintains the tree arises from three assumptions: the topology will remain relatively static, nodes will communicate frequently, and latency should be minimized when communication occurs. By precomputing a forwarding tree, the mesh nodes remain ready to forward any packet at any time. In situations where the topology changes or traffic is infrequent (for example, a sensor mesh where sensor values are collected once per week), the overhead incurred in maintaining routes may not be warranted. Instead, it may be more efficient to use an on-demand approach in which the mesh nodes wait for a packet, find a route, send the packet, and then delete the route.

Non-storing mode is designed for networks in which nodes have limited memory and CPU resources. To minimize local storage and processing, each node learns only two things: the set of directly reachable neighbors and the identity of one neighbor that leads to the border router. The border router is assumed to have substantial amounts of memory and processor capability, and can perform all the necessary path computation. The border router learns the complete topology of the mesh, and handles all source routing. When a node has a datagram to send, the node forwards the datagram to the border router. If the datagram is destined for a site on the Internet, the border router forwards the datagram. If the datagram is destined for another node in the mesh, the border router encapsulates the datagram in an outer datagram, uses its copy of the topology information to insert a source-route header in the outer datagram, and forwards the encapsulated datagram across the mesh. At each step, a node in the mesh finds the address of a neighbor in the source-route header, and uses the address to send the datagram to the specified neighbor. Perhaps the most serious consequence of using IPv6 to implement a route-over mesh arises from the requirement for source routing in non-storing mode and the size of an IPv6 source-route header.

Non-storing mode may seem to waste network resources, because a datagram sent from one node to another first goes to the border router and then to its destination. The worst case occurs when a packet is sent between a pair of nodes that are two hops apart but whose forwarding trees intersect only at the border router—instead of two hops, the packet may traverse N hops, where N is the number of nodes.

Nevertheless, the ZigBee Alliance selected non-storing mode to permit individual nodes to have extremely small memories and slow CPUs, minimizing both energy consumption and cost.

An important point about storing mode arises from limitations on scaling: in the worst case, information about mesh connectivity requires space proportional to N , the number of nodes in the network. Although most ZigBee mesh networks are expected to have fewer than 24 nodes, some mesh networks contain thousands of nodes. Consequently, for a large mesh, storing mode implies that each node must store tables that are large relative to the memories available on small devices (for example, 64 or 128 kilobytes of RAM). Using non-storing mode allows small battery-operated nodes to store only the following items in memory:

- A list of directly reachable neighbors and the MAC address of each
- The identity of the neighbor that is currently serving as the path to the border router

Note that the memory required when using non-storing mode is proportional to the number of reachable neighbors, which may be substantially less than the number of nodes in the mesh. In fact, even in a dense mesh, a node can restrict the list to the top K neighbors (that is, the K neighbors that report the highest signal strength).

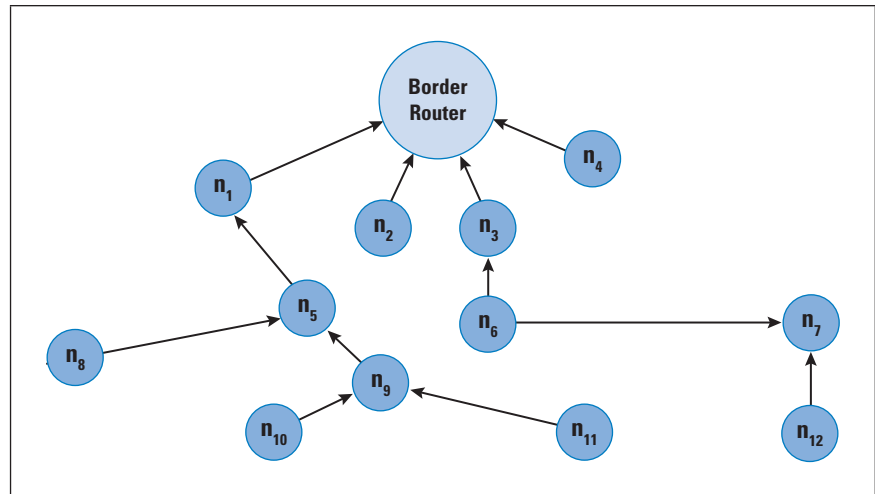
Routing Protocol for Low-Power and Lossy Networks

The IETF has defined a routing protocol that can be used with IPv6 in a route-over mesh network. Known as *Routing Protocol for Low-Power and Lossy Networks* (RPL), the protocol allows nodes to advertise direct connections and to learn about other connections in the mesh. RPL defines an IPv6 header so that datagrams can carry RPL information in addition to a payload. The ZigBee IP standard specifies the use of non-storing mode. In non-storing mode, RPL propagates connection information *upward* to the border router. The border router runs a special version of RPL software that gathers the information, meaning that it learns the topology of the entire mesh.

When it learns the topology, the border router computes a forwarding tree. Instead of imposing a tree on an undirected graph, RPL makes each link directed, with the direction toward the root (that is, toward the border router). Thus, RPL calls the graph of the mesh topology a *Destination-Oriented Directed Acyclic Graph* (DODAG). Figure 2 illustrates the DODAG form of the tree in Figure 1.

Although links in the DODAG point toward the root, the representation is merely a detail of the protocol, and does not dictate packet flow. In particular, when a border router needs to send a datagram to one of the nodes, the border router uses the DODAG in the reverse direction by composing a source-route header that lists nodes down the tree (that is, in the reverse of the arrows in the figure).

Figure 2: The DODAG RPL Defined for the Tree in Figure1.



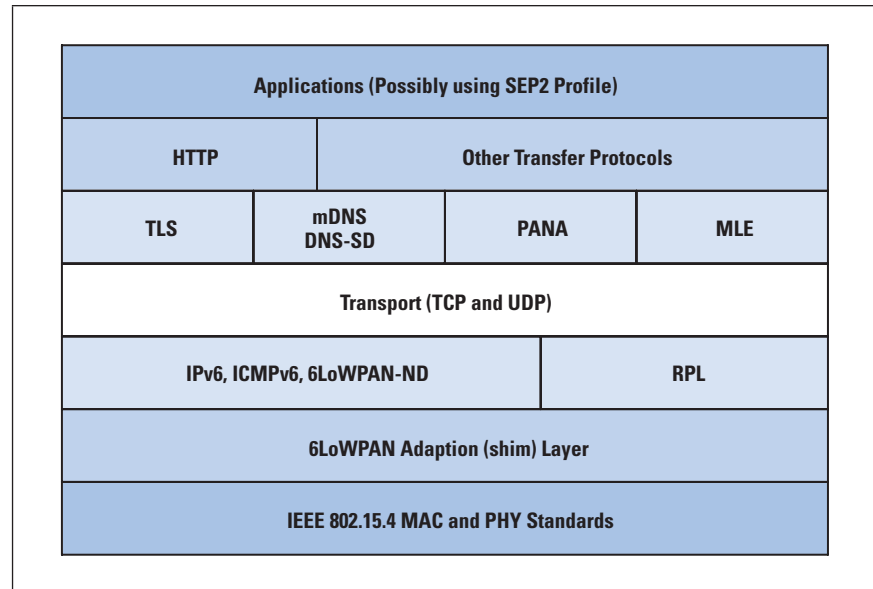
RPL separates nodes in the tree into three types: *root* (the border router acts as the root of the DODAG), *leaf* (that is, a node that has only one connection), and *intermediate* (a node that has at least two connections and forwards datagrams). Because intermediate nodes forward packets, the ZigBee IP standard uses the term *ZigBee IP router*, abbreviated *ZIP router*. Leaf nodes do not run RPL. Instead, each leaf node connects to a ZIP router, which handles all forwarding responsibilities for the leaf node. That is, a ZIP router informs the border router about each leaf node to which it connects.

RPL is much more complex than indicated here. For example, RPL can be used with storing mode; messages must specify the mode being used. In addition, RPL can distribute information down the tree. For example, it is possible to use RPL to inform nodes about the IP prefixes being used.

Other Protocols in the ZigBee IP Specification

Besides the major protocols defined previously, the ZigBee IP specification includes many other protocols. ZigBee IP uses IPv6, *Internet Control Message Protocol Version 6* (ICMPv6), TCP, *User Datagram Protocol* (UDP), *Protocol for carrying Authentication for Network Access* (PANA), *multicast DNS* (mDNS), *DNS Service Discovery* (DNS-SD), and *Transport Layer Security* (TLS), which is used in conjunction with PANA, *Extensible Authentication Protocol* (EAP), and *Extensible Authentication Protocol Transport Layer Security* (EAP-TLS) for authentication. Applications that conform to the *Smart Energy Profile 2.0*^[6] use traditional HTTP as a transfer protocol. As an alternative, the IETF is defining a new protocol known as the *Constrained Application Protocol* (CoAP)^[7] that is designed for constrained networks, including an 802.15.4 mesh. Figure 3 summarizes the major protocols.

Figure 3: The Layering of Major Protocols Used in a ZigBee Mesh Network



Assessment of Using IPv6 Route-Over for a Mesh

Many questions arise about the use of IPv6 and a route-over approach in a mesh of low-power wireless sensor nodes. Is a route-over design better than a mesh-under design? How much additional protocol overhead is needed for route-over? Will a route-over implementation require more memory than a mesh-under implementation? If so, how much more? Does it make sense to use IPv6 in a mesh network that has an extremely small MTU and extremely slow links? If RPL uses non-storing mode, what is the overhead in terms of additional packet forwarding? How much of IPv6 and the IPv6 support protocols must change to make them operational in the mesh environment?

As mentioned previously, the IPv6 standard specifies that IPv6 cannot be used over a network with an MTU smaller than 1280. Thus, 6LoWPAN adds a shim layer to divide a datagram into fraglets for transmission. Fraglet transmission has an advantage over conventional fragmentation: all fraglets must arrive sequentially and in order. That is, after the first fraglet reaches a receiver, subsequent frames from the sender must contain the rest of the fraglets, with no fraglets from other datagrams. Thus, if an incoming frame does not contain the expected fraglet, a receiver discards the entire datagram. The advantage of fraglets compared to traditional fragmentation lies in reduced memory usage: a receiver does not need to store buffers for a set of partial datagrams.

The chief disadvantage of the fraglet approach arises from a combination of three factors: large datagram size, an extremely small MTU, and higher probability of loss. The datagram size is especially large for datagrams sent from the border router to an individual mesh node because the extra level of encapsulation adds an IPv6 base header and a source-route header.

The result is that a minimum-size datagram (1280 octets) will be divided into 11 fraglets. If the probability of losing a given packet is p ($0 < p \leq 1$), the probability of losing the entire datagram is much higher than p . Although the 6LoWPAN specification recognizes lossy behavior, relying on fraglet transmission can increase retransmissions (and latency).

Another issue related to MTU arises because a border router must choose an MTU for datagrams that arrive from outside the mesh. The IPv6 standard specifies a minimum link MTU of 1280. If the border router enforces an MTU of 1280 on external links, when a datagram arrives from the outside the datagram must be further encapsulated before transmission across the mesh. Unfortunately, the additional header increases the datagram size to $1280 + \delta$, making it larger than the 6LoWPAN MTU of 1280. One solution defines the 6LoWPAN MTU to be $1280 + \delta$, but requires the border router to enforce an MTU of 1280 for external sources. Unfortunately, embedding such special restrictions in IPv6 code reduces the generality—employing nonstandard techniques to handle a mesh makes the protocol stack brittle. For example, if someone invents a new path MTU discovery mechanism, the new mechanism cannot be integrated into the border router until it has been modified to honor the special MTU rules.

Designers realized that conventional IPv6 protocols cannot be used to configure a radio link or to perform a two-way signal assessment, so they created MLE. They also had to replace IPv6 Neighbor Discovery because nodes in the mesh share a single IP prefix even though they do not share a single broadcast domain. On the surface, it might appear that MLE and IPv6-ND could be modified to work together: MLE finds neighbors and IPv6-ND uses the information from MLE to maintain a list of MAC addresses for the neighbors. However, IPv6-ND also handles other tasks. For example, IPv6-ND uses ICMPv6 messages to propagate network information, including network prefixes. Unfortunately, RPL also provides a way to propagate a network prefix downward from a border router to mesh nodes. Should IPv6-ND or RPL be used? Whatever one decides, one of the two protocols must be modified to avoid having a race in which both protocols attempt to propagate address prefixes at the same time. ZigBee IP solves the problem by specifying that only 6LoWPAN-ND is to be used for propagating network configuration information.

We said that to conserve power, nodes in a ZigBee mesh can choose to sleep. Interestingly, the address-detection mechanism can require a node to spend extra power. To see why, we must know two facts. First, an address registration is assigned a lifetime interval during which the registration remains valid. Second, to conserve power, sleep mode shuts down as much hardware as possible. Thus, on a low-power node, the clock may not continue running during a sleep cycle.

Now consider what happens when a node awakens. If the sleep cycle is sufficiently long or its clock was not running, a node cannot know whether another node arrived and registered a duplicate IP address with the border router. Thus, after awakening, the node must send a message to re-register, and must receive a reply before using its address. In a conventional network, the IPv6 use of a /64 address prefix and embedded MAC addresses makes duplicate addressing unlikely. However, 802.15.4 includes a 16-bit MAC address, meaning nodes must follow the protocol to avoid duplicate addresses.

Another unexpected complication arises from mesh routing. Routing protocols used in conventional networks choose shortest paths (unless policy dictates an alternative). Mesh routing protocols must contend with multiple free variables: the signal strength of radios along the path, the length of the path, and the probability of interference. Thus, the shortest path through the mesh may not be optimal. More important, there is no easy way to estimate the probability of interference or to quantify the tradeoff between path length and signal quality. In fact, it may even be difficult to calculate the relationship between signal quality and effective throughput. Power sources can further complicate routing. For example, we can imagine a routing system that prefers nodes that obtain power from a continuous power source over nodes that obtain power from a battery. As a consequence of multiple items to optimize and no obvious objective function, mesh routing protocols can be more complex and more difficult to tune than conventional routing protocols.

Part of the inefficiency in a ZigBee mesh arises from a fundamental IPv6 design decision: an IPv6 datagram header cannot be modified after the datagram is in transit. To transit a mesh, a datagram must include a source-route header and an RPL header. However, the extra headers are relevant only within the mesh, and must be removed if a datagram leaves the mesh. If extra headers are allowed to remain, the datagram might pass across another ZigBee mesh, where the information could be misinterpreted, causing datagrams to be misrouted. Similarly, if a datagram passes from the outside into the mesh, the appropriate headers must be added. As a consequence of the IPv6 design, headers cannot be added or removed. Therefore, the only viable option is encapsulation: each IPv6 datagram sent across the mesh must be encapsulated in an IPv6 datagram that has the appropriate headers. In a network with a large MTU, IP-in-IP encapsulation adds a small overhead. With an 802.15.4 radio, however, the hardware MTU is only 127 octets. Thus, one pair of source and destination IPv6 addresses occupies more than 25% of the MTU. By contrast, a pair of IPv4 addresses accounts for just over 6% of the MTU. Unlike IPv4, which allows options to be inserted in an existing header, IPv6 requires an encapsulating header to hold a source route option. The approach requires adding multiple IPv6 addresses, which can easily generate one or more extra fraglets. As a result, the choice of IPv6 instead of IPv4 increases the overhead significantly, and makes the resulting network less efficient.

From the observations discussed previously, we conclude:

Although it is possible to use a route-over paradigm and IPv6 for an 802.15.4 ZigBee mesh, doing so means inventing alternative protocols, creating special exceptions, and incurring significant overhead.

Summary

An emerging trend focuses on an Internet of Things, in which intelligent embedded systems that sense and control their environment use Internet technology to communicate. Examples include sensors in vehicles, residences, office buildings, shopping malls, and civil infrastructure.

A consortium of vendors known as the ZigBee Alliance is specifying standards for a networking technology that uses IEEE 802.15.4 wireless network hardware to form a mesh of Smart Grid sensors. A ZigBee network has a border router that connects to the outside; other nodes in the mesh self-organize to connect and establish forwarding.

In terms of protocols, the ZigBee Alliance is working with the IETF on a route-over approach that uses IPv6. In principle, a route-over system uses IP for all forwarding. In practice, IPv6 and standard IPv6 support protocols are insufficient for a low-power, lossy wireless mesh technology that has a small MTU and a low data rate. Consequently, work has focused on replacing many parts of IPv6, adding a shim layer to accommodate small MTU, inventing a new protocol that tests signal strength and establishes links, and building a new routing protocol that constructs a forwarding tree. Even with the changes, the design of IPv6 does not match IEEE 802.15.4 radio technology well.

Acknowledgment

Material in this article has been taken with permission from [8].

References

- [1] The IEEE 802.15.4 standard can be purchased from:
<http://webstore.ansi.org/RecordDetail.aspx?sku=IEEE%20802.15.4-2011>
- [2] Information about the ZigBee Alliance can be found at:
<http://www.zigbee.org/>
- [3] Information about the Internet Engineering Task force can be found at: <http://www.ietf.org/>
- [4] The ZigBee IP specification can be found at:
<http://zigbee.org/zigbee-for-developers/network-specifications/zigbeeip/>

- [5] Stephen E. Deering and Robert M. Hinden, “Internet Protocol, Version 6 (IPv6) Specification,” RFC 2460, December 1998.
- [6] A complete list of ZigBee specifications can be found at:
`http://zigbee.org/zigbee-for-developers/applicationstandards/`
and:
`http://zigbee.org/zigbee-for-developers/network-specifications/`
- [7] The specification for the Constrained Application Protocol (CoAP) can be found at:
`http://datatracker.ietf.org/doc/draft-ietf-core-coap/`
- [8] Douglas E. Comer, *Internetworking with TCP/IP Volume 1: Principles, Protocols, and Architecture*, sixth edition, Prentice Hall, ISBN 0-13-608530-X.
- [9] David Lake, Ammar Rayes, and Monique Morrow, “The Internet of Things,” *The Internet Protocol Journal*, Volume 15, No. 3, September 2012.
- [10] T. Sridhar, “Cloud Computing—A Primer: Part One,” *The Internet Protocol Journal*, Volume 12, No. 3, September 2009.
- [11] T. Sridhar, “Cloud Computing—A Primer: Part Two,” *The Internet Protocol Journal*, Volume 12, No. 4, December 2009.

DOUGLAS E. COMER is a Distinguished Professor of Computer Science at Purdue University. Formerly, he served as VP of Research and Research Collaboration at Cisco Systems. As a member of the original IAB, he participated in early work on the Internet, and is internationally recognized as an authority on TCP/IP protocols and Internet technologies. He has written a series of best-selling technical books, and his three-volume Internetworking series is cited as an authoritative work on Internet technologies. His books, which have been translated into 16 languages, are used in industry and academia in many countries. Comer consults for industry, and has lectured to thousands of professional engineers and students around the world. For 20 years he was editor-in-chief of the journal *Software-Practice and Experience*. He is a Fellow of the ACM and the recipient of numerous teaching awards. E-mail: `comer@cs.purdue.edu`

Letters to the Editor

Hi Geoff,

I found a copy of *The Internet Protocol Journal*, Volume 17, No. 1, on my desk this morning and really enjoyed your article “A Question of DNS Protocols.” If you have a couple of spare minutes, some questions occurred to me:

Could the number of open resolvers be due to the implementer’s lack of experience and/or “open” out-of-the-box resolver configuration? If the default configuration provided by vendors was restricted (for example, for Internet Systems Consortium’s *BIND* specifying **edns-udp-size**, **max-udp-size**, **rate-limit** and so on), do you think this restriction might slowly reduce the number of open resolvers? I’m happy to push this idea within Red Hat if you think it could be worthwhile.

To your knowledge was there any follow-up research into which clients were failing to transition to TCP queries? If it’s one or two resolvers, maybe the maintainers could be engaged directly to push TCP support. Theoretically, if people start restricting their *Extension Mechanisms for DNS* (EDNS) size to 512, these clients will break anyway, correct?

Thanks again for your time,

Cheers,

—Morgan Weetman, Red Hat Consulting
mweetman@redhat.com

The author responds:

Thanks for your questions.

The basic problem we observe with these open resolvers is that there is a large-scale use of *Customer Premises Equipment* (CPE) network interface devices that include DNS resolution. Evidently, the intent was to take DNS queries from the “inside” and pass them on to the *Dynamic Host Configuration Protocol* (DHCP)-provided “outside” DNS resolver, and cache the results that come back. The local cache provides a small increment in perceived performance on the inside, and the DNS forwarder removes the pressure of the *Network Address Translation* (NAT) function of the CPE for these *User Datagram Protocol* (UDP) transactions.

All good, but there are an annoyingly large number of units, evidently numbering in the tens of millions, where the DNS resolver function on the units has no idea what is the “inside” and what is the “outside.” It will happily treat queries coming on in the “outside” and resolve these names in the same fashion as if the query was received on the “inside,” and then send the response back to the “outside” query agent.

The critical configuration element that appears to be missing on these units is a filter to drop incoming packets with destination port 53 that are addressed to the exterior network interface on the unit.

With respect to your question relating to follow-up research, I should reiterate that we found that some 17% of “visible” resolvers appear not to fail over to use TCP when they receive a UDP DNS response with the truncated bit set. Now the problem with the DNS is that there is very little in the way of fingerprinting of resolvers, so apart from their IP address it’s challenging to understand what is going on. A common assumption is that these units live behind a firewall that prevents TCP port 53 from traversing in either direction, but it’s an assumption I’ve not tested for explicitly. This assumption leads to the strong suspicion that it’s not DNS resolvers per se, but the environment into which they are deployed that is the problem here.

And with respect to using a small EDNS0 size field, then yes—if the response is larger than the offered size, then the responder will cram as much as it can into the offered size and set the truncate bit. The querier is meant to interpret this response as a signal to re-query over TCP, and if it fails for whatever reason, then they are unable to complete name resolution.

Thanks for your questions. I trust I’ve been able to answer them here.

—*Geoff Huston, APNIC*
gih@apnic.net

Geoff,

Regarding your article on DNS Protocols in IPJ Volume 17, No. 1, I have a few observations that may be of interest:

1. Blocking of TCP/Port 53 throughout the Internet, especially on endpoint networks, is a real issue. The security myth that blocking TCP/53 somehow makes DNS “more secure” by disallowing zone transfers originated sometime in the mid-1990s, and persists to this day (obviously, all that’s required to disallow unauthorized zone transfers is to configure one’s authoritative DNS servers properly). More than a few large-scale authoritative DNS hosters and DNS registrars who offer authoritative DNS hosting services incorrectly block TCP/53 queries to their authoritative DNS server arms. I run into all these issues with some regularity.
2. Many, many authoritative and recursive DNS servers are not scaled and tuned to support large numbers of simultaneous TCP connections, and will experience availability problems because they’re overwhelmed by a comparatively small number of TCP DNS queries versus an equivalent number of UDP DNS queries. This problem also holds true of load-balancing devices that are often placed in front of both recursive and authoritative DNS server farms. I run into this problem with some regularity, as well.

3. The relative latency of TCP connection setup times combined with the practices of dynamically assembling Web pages/app views from multiple named/numbered Web servers (whether they're actually separate servers or merely additional DNS records for the same actual server) in an attempt to speed page-load times via parallelism is also a significant challenge.

The second problem is potentially resolvable (pardon the pun), but would require a high degree of capital and operational expenditure committed across many organizations to make it practical.

The third problem is a real challenge, because the designers of Web servers and apps would have to be re-educated—and they are often completely siloed from any individuals or organizations with operational experience.

The first problem is well-nigh intractable, because after filters are put into place (in this case, out of misinformed ignorance), all too often they are never removed. It's a pretty safe bet that the networks that incorrectly filter TCP/53 at this late date are never going to “see the light,” so to speak.

So, while I agree that DNS over TCP would have many desirable characteristics, chief among them reducing the DNS reflection/amplification *Distributed Denial of Service* (DDoS) vector, I consider it unlikely to be practical. I first looked at this issue in 2005 when I was at Cisco, and all the problems mentioned before that applied then also apply now—in many cases with a much higher degree of prevalence than a decade ago.

Ultimately, the best solution from a number of standpoints may well be to move away from the DNS entirely towards something similar to the *Peer Name Resolution Protocol* (PNRP). I believe that more and more applications and services are going to end up being hyper-distributed among nodes we tend to think of today as “clients” (for example, mobile devices, CPE, all the various types of embedded systems)—and that because of its universal necessity and applicability, the migration of name resolution/directory services to such a model should be actively pursued.

—Roland Dobbins, Arbor Networks
rdobbins@arbor.net

The author responds:

Hi Roland,

Thanks for your comments. The speculation as to where next is certainly interesting, and the overriding consideration as to whether and how we can stay within a uniform and consistent name space while at the same time moving away from the existing DNS structure and the related resolution protocol is an aspect that greatly interests me. Again thanks for taking the time to note down your comments.

Kind regards,

—Geoff Huston, APNIC
gih@apnic.net

Book Review

Internet Peering Playbook

The Internet Peering Playbook—Connecting to the Core of the Internet, 2014 Edition, by William B. Norton, DrPeering Press, ISBN-13: 978-1937451110.

When I realized I needed to understand how Internet peering worked, it was timely that Bill Norton shared his book with me at the *North American Network Operators' Group* (NANOG) conference last summer. I read it on my 6-hour flight home and finished it the next night. Not only is it an easy read for a non-engineer in the telecom space, it is clear and concise on how peering actually works. Norton starts off with the basics on how the peering ecosystem works, who peers privately and who peers publicly, and why. He details the recent hoopla over Comcast, Netflix, and Level 3 with their public versus paid peering and simply delivers the facts and none of the emotion iterated on some of the players' blog sites.

Norton then reaches well beyond the basics of peering to include examples of “tricks of the trade.” He graphically lays out samples, and explains them so well the book makes for a great read. The topic of discussion today is that some ISPs are using their access network as a monopoly, and want to charge content providers for the traffic that runs on their network. Tricks of the trade start with simple bundling options that hide additional traffic, but quickly can end up “playing chicken,” where the network traffic becomes significantly asymmetric, usually resulting in one end dumping traffic on another peer's network. The network infrastructure is usually also asymmetric, resulting in a request to re-negotiate from what was a free peering situation to a pay-to-play requirement. Both sides believe the other side needs them more. Thus playing chicken is initiated.

Sometimes a new peering negotiation is made, sometimes not. If not, the result can be de-peering and possibly a severe traffic disruption. As Norton says “it really tests the assertion that both sides are receiving *equal value* from the relationship.” In most cases, additional connectivity is deployed and traffic is spread across more sites to even the load.

In order to peer, traffic volumes must meet a minimum to be worth the allocation of ports. Yet “bluffing” or claiming the traffic load is adequate when it's not is one way to get a peering transaction initiated. Another way is to claim performance problems that can be easily solved by simply peering when no peering had been in place previously. Because coordinators rarely have time for in-depth traffic analysis, this type of peering becomes another trick of the trade to gain free peering, at least in the short term.

Be open, be loud, be a friend, and be sweet are all positive routes for proper peering techniques. Negative approaches such as *Make It Long and Difficult* (MILD) are used, where discussions are prolonged and appear open but nothing really happens. Norton even brings up peering tactics that don't work, such as trying to dominate in a single foreign market, public badgering (I assume at NANOG events), holding content hostage, sending blind requests, or simply lacking knowledge in the peering backbone space is enough to be shunted in this tight community.

A recent report from *Measurement Lab* (M-Lab) shows “sustained performance degradation for access customers when traversing interconnections” and displays proof in numbers that peering degradation occurs. The report explains very clearly that the traffic degradation is due to the business relationships of the interconnections, and is not at all the fault of technical problems^[1].

What the Internet peering community does not want is for the *Federal Communications Commission* (FCC) to try to regulate this market. However, the FCC is beginning to take steps to criticize ISPs for “throttling” traffic. As our world revolves around increasing access to bandwidth, not everyone needs to understand the importance of Internet peering, but it surely is interesting!

—Eve Griliches
egriliches@btisystems.com

[1] http://www.measurementlab.net/static/observatory/M-Lab_Interconnection_Study_US.pdf

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don't have access to it. Contact us at ipj@protocoljournal.org for more information.

MANRS: Improving Global Routing Security & Resilience

Most end users don't give much—if any—thought to things like the Internet's global routing system because, for the most part, the Internet has *just worked* for years. However, in several instances vulnerabilities in the security and resilience of that routing system have manifested themselves: a 2008 incident that made *YouTube* temporarily unreachable around the globe, multiple cases of Internet traffic deflection by some Chinese Internet Service Providers, and an April 2014 incident in which an Indonesian network operator mistakenly claimed that it “owned” many of the world's networks, just to name a few.

Internet security, in general, is a difficult area when it comes to incentivizing network operators to act with the good of the whole Internet in mind, and security of the global Internet infrastructure, be it the *Domain Name System* (DNS) or routing, brings additional challenges because the utility of security measures depends on coordinated actions of many other parties. Thus, while technology is an essential element, technology alone is not sufficient. To stimulate visible improvements across the entire Internet, we need a culture of collective responsibility and action.

The good news is that throughout the history of the Internet, we've seen amazing feats of collaboration among participants and shared responsibility for its smooth operation. Collaboration and shared responsibility are two of the pillars supporting the Internet's tremendous growth and success, as well as its overall security and stability.

So, how can we collectively help improve the security and resilience of the global Internet routing system and prevent the kinds of vulnerabilities we mentioned earlier? One of the approaches is the *Routing Resilience Manifesto* initiative, which features the *Mutually Agreed Norms for Routing Security* (MANRS) document. This initiative of several leading network operators, supported and coordinated by the Internet Society, was publicly launched on November 6, 2014.

MANRS captures the collaborative spirit that has been a hallmark of the Internet's growth, and provides motivation and guidance to network operators in addressing issues of security and resilience of the global Internet routing system.

MANRS defines a compact and clear set of actions that network operators should implement to improve routing security on their own networks and across the Internet as a whole; specifically:

- Prevent propagation of incorrect routing information.
- Prevent traffic with spoofed source IP addresses.
- Facilitate global operational communication and coordination between network operators.
- Facilitate validation of routing information on a global scale.

The Routing Resilience Manifesto is more than just the MANRS document; it is a commitment to improve the global Internet. In order to become a participant of this initiative, a network operator has to implement one or more of the identified actions.

More than a dozen network operators around the world have already signed up, and we expect more to join. You can learn more about the effort and sign up to be part of this exciting movement to make the Internet a safer place for everyone at www.routingmanifesto.org

IAB Statement on Internet Confidentiality

The *Internet Architecture Board* (IAB) issued the following statement on November 14, 2014:

In 1996, the IAB and *Internet Engineering Steering Group* (IESG) recognized that the growth of the Internet depended on users having confidence that the network would protect their private information. RFC 1984^[1] documented this need. Since that time, we have seen evidence that the capabilities and activities of attackers are greater and more pervasive than previously known. The IAB now believes it is important for protocol designers, developers, and operators to make encryption the norm for Internet traffic. Encryption should be authenticated where possible, but even protocols providing confidentiality without authentication are useful in the face of pervasive surveillance as described in RFC 7258^[2].

Newly designed protocols should prefer encryption to cleartext operation. There may be exceptions to this default, but it is important to recognize that protocols do not operate in isolation. Information leaked by one protocol can be made part of a more substantial body of information by cross-correlation of traffic observation. There are protocols which may as a result require encryption on the Internet even when it would not be a requirement for that protocol operating in isolation.

We recommend that encryption be deployed throughout the protocol stack since there is not a single place within the stack where all kinds of communication can be protected.

The IAB urges protocol designers to design for confidential operation by default. We strongly encourage developers to include encryption in their implementations, and to make them encrypted by default.

We similarly encourage network and service operators to deploy encryption where it is not yet deployed, and we urge firewall policy administrators to permit encrypted traffic.

We believe that each of these changes will help restore the trust users must have in the Internet. We acknowledge that this will take time and trouble, though we believe recent successes in content delivery networks, messaging, and Internet application deployments demonstrate the feasibility of this migration. We also acknowledge that many network operations activities today, from traffic management and intrusion detection to spam prevention and policy enforcement, assume access to cleartext payload. For many of these activities there are no solutions yet, but the IAB will work with those affected to foster development of new approaches for these activities which allow us to move to an Internet where traffic is confidential by default.

[1] IAB and IESG, “IAB and IESG Statement on Cryptographic Technology and the Internet,” RFC 1984, August 1996.

[2] Stephen Farrell and Hannes Tschofenig, “Pervasive Monitoring Is an Attack,” RFC 7258, May 2014.

Upcoming Events

The *North American Network Operators’ Group* (NANOG) will meet in San Antonio, Texas February 2–4, 2015; and in San Francisco, California, June 1–3, 2015. See: <http://nanog.org>

The *Internet Corporation for Assigned Names and Numbers* (ICANN) will meet in Singapore, February 8–12, 2015; in Buenos Aires, Argentina, June 21–25, 2015; and in Dublin, Ireland, October 18–22, 2015. See: <http://icann.org>

The *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) will be held in Fukoka, Japan, February 24–March 6, 2015. See: <http://www.apricot.net>

The *Internet Engineering Task Force* (IETF) will meet in Dallas, Texas, March 22–27, 2015; in Prague, Czech Republic, July 19–24, 2015; and in Yokohama, Japan, November 1–6, 2015. See: <http://www.ietf.org/meeting/>

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

The Internet Protocol Journal (IPJ) is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. Publication of IPJ is made possible by:

<i>Supporters</i>  	<i>Diamond Sponsors</i>  
<i>Ruby Sponsor</i> 	<i>Sapphire Sponsors</i>  TEAM CYMRU INSIGHT THAT IMPROVES LIVES 
<i>Emerald Sponsors</i>            	
<i>Corporate Subscriptions</i>     	

Individual Sponsors

Lyman Chapin, Steve Corbató, Dave Crocker, Jay Etchings, Hagen Hultzsch, Dennis Jennings, Jim Johnston, Merike Kaeo, Bobby Krupczak, Richard Lamb, Tracy LaQuey Parker, Bill Manning, Mike O'Connor, Tim Pozar, George Sadowsky, Helge Skrivervik, Rob Thomas, Tom Vest, Rick Wesson.

For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

ADDRESS SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Fred Baker, Cisco Fellow
Cisco Systems, Inc.

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Steve Crocker, Chairman
Internet Corporation for Assigned Names and Numbers

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

March 2015

Volume 18, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Scaling the Root.....	2
Gigabit Ethernet	20
Fragments	33
Call for Papers.....	34
Supporters and Sponsors	35

FROM THE EDITOR

In the summer of 1984 I joined SRI International in Menlo Park, California, working at the *Network Information Center* (NIC). The NIC provided numerous services relating to the ARPANET and MILNET, including a telephone help line, various printed materials, login credentials for dialup users, and the all-important **HOSTS.TXT** file that mapped host names to their corresponding IP addresses. The **HOSTS.TXT** file was updated once a week and made available via FTP and from the NIC's dedicated name server. The original documents that described the *Domain Name System* (DNS) had been published in late 1983, and all of us at the NIC were keenly aware that a transition from a centrally maintained file to a distributed and hierarchical name resolution system would soon be underway. The original design of the DNS has proven itself to be both robust and scalable, and the protocol has been enhanced to support IPv6, as well as security (DNSSEC). In our first article, Geoff Huston gives an overview of the DNS and discusses possible ways in which to further scale the system.

Ethernet has been a critical component of *Local-Area Networks* (LANs) for many decades. As with most networking technologies, there have been several iterations of the Ethernet standards, each providing orders of magnitude faster transmission rates. In our second article, William Stallings gives an overview of recent developments and standardization efforts for *Gigabit Ethernet*.

If you received a printed copy of this journal in the mail, you should also have received a subscription activation e-mail with information about how to update and renew your subscription. If you didn't receive such a message, it may be because we do not have your correct e-mail address on file. To update and renew your subscription, just send a message to ipj@protocoljournal.org and include your subscription ID. Your subscription ID is printed on the back of your journal.

Let me remind you that IPJ relies on the support of numerous individuals and organizations. If you or your company would like to sponsor IPJ, please contact us for further details. Our website at protocoljournal.org contains all back issues, subscription information, a list of current sponsors, and much more. Our first blog entry "Notes from NANOG 63," was recently posted on the website.

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

Scaling the Root

by Geoff Huston, APNIC

The *Domain Name System* (DNS) of the Internet is a modern-day miracle that has proved to be exceptionally prodigious. This technology effectively supported the operation of the Internet from a scale of a few hundred thousand users to today's system of some 3 billion users and an estimated 8 to 10 billion devices. Not only has it supported that level of growth, it has done so without any obvious cracks within the basic protocol. But that point should not imply that nothing has changed in the DNS protocol over the past 30 years. While the basic architecture of the DNS as a simple query/response protocol has remained consistent over this period, we have adorned the base query/response protocol with various optional “bells and whistles” that are intended to improve its robustness, and we have constantly updated the platform infrastructure of the DNS to cope with ever-increasing query loads as the Internet expands. In this article we look at one aspect of this effort: the current considerations of how to scale the root of the DNS.

An Introduction to the DNS

The DNS is a hierarchically distributed naming system. The inherent utility of a name system in a digital network lies in an efficient and consistent mapping function between *names* and *IP addresses*. With the progenitor of the DNS, this mapping function took the form of a single file (**HOSTS.TXT**) that listed all the active host names and the corresponding IP address for each address that resided on each host. The problem with this approach was the coordination of entries in this hosts file so that all hosts had a consistent view of the name space of the network. As the network grew, the administrative burden of coordinating this burgeoning name space became unworkable. The DNS replaced this replicated file with a dynamic query system that allowed hosts to query a distributed name data base and retrieve the current value of the mapped IP address.

To achieve this goal, the name system uses a hierarchical structure. A name in the DNS is a sequence of labels that scan from left to right. This sequence of labels can be viewed in a pairwise fashion, where the label to the left is the “child” of the “parent” label to the right. The apex of this hierarchical name structure is the root, which is notionally defined as the trailing “.” at the rightmost part of a *Fully Qualified Domain Name* (FQDN).

As a name-resolution system, the DNS is constructed as a collection of agents that ask questions and receive answers, so called *Resolvers*, and a set of servers that can provide the authoritative answer for a certain set of questions, *Authoritative Name Servers*. A set of name servers that are configured as being authoritative for a given *zone* should provide identical answers in response to queries for names that lie within this zone.

The task of a resolver is to ask questions of servers. But if each server is able to provide answers for only specific zones, the question then becomes: which server to ask?

For example, to resolve the name `www.example.com.`, a resolver needs to find the authoritative name servers for the zone `example.com.` This information (the set of authoritative name servers for a zone) is stored as part of a zone delegation record that is loaded into the parent zone. In our example, to resolve `www.example.com.`, a resolver needs to query any of the servers for the `example.com.` zone. These servers can be found by querying any of the authoritative name servers for the `com.` zone. But to do that a resolver needs to know who are the authoritative name servers for the `com.` zone. This information is held in the Root Zone, and can be retrieved by sending a query to any of the *Root Zone Servers*.

The way this process is implemented in the DNS is that when an authoritative name server is queried for a name that lies within the scope of a delegated child of the zone of the server, the server will respond to the query not with the desired answer, but with the set of authoritative name servers for the immediate child zone that encompasses the name of the query.

Therefore, when a recursive resolver is passed a query relating to a name about which it has no knowledge, it will first send the query for this name to a root server. The response is not the desired information, but the name servers that are authoritative for the top-level domain name being queried. The recursive resolver will then query one of these name servers for the same name, and will receive in response the name servers that are authoritative for the next level of domain name, and so on. For example, a resolver with no a priori knowledge of the DNS other than a list of servers for the root zone, when attempting to resolve the name `www.example.com.`, will first query one of the root servers for this name. The response of the root server should be the set of authoritative name servers for the `.com` zone. The resolver will next query one of these servers with the same query. The response will be the set of servers for the `example.com.` zone. When one of these servers is queried, it should respond with the desired *Resource Record* (such as the mapped IP address of the name).

To make the DNS work efficiently, resolvers typically remember these responses in a local cache, and will not re-query for the same information unless the cached entry has timed out in the local cache. For example, a subsequent query for `ftp.example.com.` would use the local cache of the resolver to select the set of authoritative name servers for `example.com` and pose this query directly to one of these servers.

So resolvers can dynamically discover and cache all aspects of the DNS, with one critical exception. The root zone of the DNS cannot be dynamically discovered in this manner, because it has no parent. To get around this problem, DNS resolvers are configured with a *root hints* file, which contains a list of IP addresses of those name servers that are authoritative for the root zone. When a resolver starts up, it sends a *root priming query* to one of the servers listed in the hints file, requesting the current set of root name servers. In this way the resolver then is primed with the current set of root name servers.

What are root servers used for?

As explained previously, one intended role of the root servers is to respond to root priming queries. Secondly, as previously explained, the root servers respond to queries relating to labels that are defined in the root zone, and also respond negatively to queries relating to labels that are not defined in the root zone. Resolvers establish the identity of name servers for those names that are at the *top level* of the DNS name hierarchy by directing a query to a root server.

Obviously, if every name-resolution attempt involved a resolver making a query to the root name servers, then the DNS root servers would have melted under the consequent load years ago! Resolvers reduce this load by caching the answers they receive, so that the profile of queries set to the root are dominated by queries for “new” names that resolvers have not previously seen (and cached). Typically, the response from the root name server is one that indicates that the name does not exist in the root zone. Given that the overwhelming role of the root servers is to reply with a “no such name” response, then it would seem that the role of the root server is somewhat inconsequential. However, resolvers do not keep a permanent copy of the responses they receive from their queries for names that exist in the root zone. They use a timed local cache, and from time to time the resolver needs to repeat the query in order to refresh the cache entry. If the root servers disappeared, these refresh queries would fail and the resolver would be unable to answer further queries about this zone. So while the predominate load on the root servers is to respond to junk queries, their continuing availability is of paramount importance to the Internet. And if the resolvers of a network were isolated from the root server constellation, then the network would, over a short period of time, cease to have a working name system.

This situation would lead us to the thought that if root servers are so critical, then a single host serving the root zone for the entire Internet would be a poor design choice. Perhaps every network should run a root server. This thought touches on numerous considerations, not the least of which includes considerations of the underlying query and response protocol that the DNS uses.

DNS Protocol Considerations

The DNS is a simple query/response protocol. In order to allow a query agent to recognize the appropriate response, the response is a copy of the original query, with additional fields included.

The two mainstream IP transport protocols are the *User Datagram Protocol* (UDP) and the *Transmission Control Protocol* (TCP). TCP is ill-suited to simple query/response transactions, given that each TCP connection requires an initial packet exchange to open a connection, and a closing exchange, and while the transaction is open the server needs to maintain TCP session context. UDP eschews this overhead, and each query can be loaded into a single UDP packet, as can the corresponding response. Although TCP is a permitted transport protocol for DNS, the overwhelming operational preference is to use UDP. It's fast, efficient, and works well. But that's not quite the entire story. We need to go down one level of the protocol stack and look at the IP protocol.

The *IPv4 Host Requirements Specification*^[1] mandates that all compliant IP host systems be able to accept and process an IP packet that is at least 576 octets. Individual IP fragments may be smaller than this number, but the host must be able to reassemble the original IP packet if it is 576 bytes or less. The consequence of this requirement is that compliant hosts need not necessarily accept an IP packet larger than 576 octets, but they will all accept packets of this size or smaller.

Allowing for 20 octets of the IP header and a maximum of 40 octets of IP header options, 516 octets of payload remain. UDP headers are 8 octets, so we would expect that if we were to define a protocol that used UDP as its transport protocol, then a UDP payload of a maximum of 508 octets of payload would be assured to reach any IPv4-compliant host. However, the DNS specification is subtly mismatched, and the *DNS Specification*^[2] specifies a DNS payload size of 512 octets or less.

So how many distinct root name servers can be listed in a priming response if we want to keep the response size to 512 octets or less? The adoption of 13 distinct root servers is a compromise between these two pressures. The exact number was the outcome of the number of distinct root server labels, and their IPv4 addresses, that can be loaded into an unsigned IPv4 UDP response to a root server priming query that is less than 512 octets. An example of a root priming response is shown in Figure 1 on page 6. (The DNS is a binary protocol that uses field compression where possible. The *dig* utility^[3] generates specific DNS queries and produces a text representation of the response.)

This response is 503 octets, which will also just fit into the 512-byte limit as defined by the *Applications Requirements Specification*^[4].

Figure 1: DNS Root
Priming Response

```

$ dig +bufsize=512 +norecurse . NS @a.root-servers.net.

; <<>> DiG 9.9.6 <<>> +bufsize=512 +norecurse . NS @a.root-servers.net.
;; global options: +cmd
;; Got answer:
;; ->>HEADER<- opcode: QUERY, status: NOERROR, id: 38316
;; flags: qr aa; QUERY: 1, ANSWER: 13, AUTHORITY: 0, ADDITIONAL: 16

;; OPT PSEUDOSECTION:
; EDNS: version: 0, flags:; udp: 1472
;; QUESTION SECTION:
;.                               IN      NS

;; ANSWER SECTION:
.                               518400 IN      NS      a.root-servers.net.
.                               518400 IN      NS      b.root-servers.net.
.                               518400 IN      NS      c.root-servers.net.
.                               518400 IN      NS      d.root-servers.net.
.                               518400 IN      NS      e.root-servers.net.
.                               518400 IN      NS      f.root-servers.net.
.                               518400 IN      NS      g.root-servers.net.
.                               518400 IN      NS      h.root-servers.net.
.                               518400 IN      NS      i.root-servers.net.
.                               518400 IN      NS      j.root-servers.net.
.                               518400 IN      NS      k.root-servers.net.
.                               518400 IN      NS      l.root-servers.net.
.                               518400 IN      NS      m.root-servers.net.

;; ADDITIONAL SECTION:
a.root-servers.net. 518400 IN      A       198.41.0.4
b.root-servers.net. 518400 IN      A       192.228.79.201
c.root-servers.net. 518400 IN      A       192.33.4.12
d.root-servers.net. 518400 IN      A       199.7.91.13
e.root-servers.net. 518400 IN      A       192.203.230.10
f.root-servers.net. 518400 IN      A       192.5.5.241
g.root-servers.net. 518400 IN      A       192.112.36.4
h.root-servers.net. 518400 IN      A       128.63.2.53
i.root-servers.net. 518400 IN      A       192.36.148.17
j.root-servers.net. 518400 IN      A       192.58.128.30
k.root-servers.net. 518400 IN      A       193.0.14.129
l.root-servers.net. 518400 IN      A       199.7.83.42
m.root-servers.net. 518400 IN      A       202.12.27.33
a.root-servers.net. 518400 IN      AAAA    2001:503:ba3e::2:30
b.root-servers.net. 518400 IN      AAAA    2001:500:84::b

;; Query time: 145 msec
;; SERVER: 198.41.0.4#53(198.41.0.4)
;; WHEN: Sun Mar 01 10:45:46 UTC 2015
;; MSG SIZE rcvd: 503

```


Evolutionary Pressures

These days it's rare to see a host impose a maximum IP datagram size of 576 octets. A more common size of IP datagrams is 1,500 octets, as defined by the payload size of Ethernet frames. In theory, an IPv4 datagram can be up to 65,535 octets, and in IPv6 the jumbo payload option allows for an IPv6 datagram of some 4 billion octets, but both of these upper bounds are very much theoretical limits.

More practical limits can be found in the unstandardized work to support large packets in 802.3 networks, where the value of 9,000 octets for *Maximum Transmission Unit* (MTU) sizes is sometimes found on vendors' IEEE 802.3 equipment for Gigabit Ethernet^[11]. However, there are two problems with this 9,000 octet packet size. Not all networks support the transmission of 9,000 octet packets without resorting to packet fragmentation, and there are classes of security middleware that reject all packet fragments on the basis of their assumed security risk. So a 1,500-octet value appears to be a practical assumption for the maximum size of an unfragmented IP datagram that has a reasonable probability of being passed through IP networks. This is a useful assumption for the root priming response, as it is now larger than 508 (or even 512) octets by default.

IPv6

The transition of the Internet to use IPv6 implies a protracted period of support for both IP protocols, and the root server priming response is no exception to this implication. When we add the IPv6 addresses of the 13 root name servers to the packet, the size of the response expands to 755 octets, as shown in Figure 2 on page 8. (As is shown in the response, the E root server does not have an IPv6 address.)

DNSSEC

The next change has been in the adoption of *Domain Name System Security Extensions* (DNSSEC), used to sign the root zone^[5]. The priming response now needs to contain the digital signature of the *Name Server* (NS) records, which expands the root priming response by a further 158 octets (Figure 3 on page 9).

Figure 2: DNS Root
Priming Response with
IPv6 Addresses

```

dig +norecurse . NS @a.root-servers.net.

; <<>> DiG 9.9.6 <<>> +norecurse . NS @a.root-servers.net.
;; global options: +cmd
;; Got answer:
;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 56567
;; flags: qr aa; QUERY: 1, ANSWER: 13, AUTHORITY: 0, ADDITIONAL: 25

;; OPT PSEUDOSECTION:
; EDNS: version: 0, flags:; udp: 4096
;; QUESTION SECTION:
;.                               IN      NS

;; ANSWER SECTION:
.                               518400 IN      NS      e.root-servers.net.
.                               518400 IN      NS      h.root-servers.net.
.                               518400 IN      NS      b.root-servers.net.
.                               518400 IN      NS      j.root-servers.net.
.                               518400 IN      NS      c.root-servers.net.
.                               518400 IN      NS      a.root-servers.net.
.                               518400 IN      NS      g.root-servers.net.
.                               518400 IN      NS      l.root-servers.net.
.                               518400 IN      NS      i.root-servers.net.
.                               518400 IN      NS      m.root-servers.net.
.                               518400 IN      NS      f.root-servers.net.
.                               518400 IN      NS      d.root-servers.net.
.                               518400 IN      NS      k.root-servers.net.

;; ADDITIONAL SECTION:
e.root-servers.net. 3600000 IN      A      192.203.230.10
h.root-servers.net. 3600000 IN      A      128.63.2.53
h.root-servers.net. 3600000 IN      AAAA   2001:500:1::803f:235
b.root-servers.net. 3600000 IN      A      192.228.79.201
b.root-servers.net. 3600000 IN      AAAA   2001:500:84::b
j.root-servers.net. 3600000 IN      A      192.58.128.30
j.root-servers.net. 3600000 IN      AAAA   2001:503:c27::2:30
c.root-servers.net. 3600000 IN      A      192.33.4.12
c.root-servers.net. 3600000 IN      AAAA   2001:500:2::c
a.root-servers.net. 3600000 IN      A      198.41.0.4
a.root-servers.net. 3600000 IN      AAAA   2001:503:ba3e::2:30
g.root-servers.net. 3600000 IN      A      192.112.36.4
l.root-servers.net. 3600000 IN      A      199.7.83.42
l.root-servers.net. 3600000 IN      AAAA   2001:500:3::42
i.root-servers.net. 3600000 IN      A      192.36.148.17
i.root-servers.net. 3600000 IN      AAAA   2001:7fe::53
m.root-servers.net. 3600000 IN      A      202.12.27.33
m.root-servers.net. 3600000 IN      AAAA   2001:dc3::35
f.root-servers.net. 3600000 IN      A      192.5.5.241
f.root-servers.net. 3600000 IN      AAAA   2001:500:2f::f
d.root-servers.net. 3600000 IN      A      199.7.91.13
d.root-servers.net. 3600000 IN      AAAA   2001:500:2d::d
k.root-servers.net. 3600000 IN      A      193.0.14.129
k.root-servers.net. 3600000 IN      AAAA   2001:7fd::1

;; Query time: 144 msec
;; SERVER: 198.41.0.4#53(198.41.0.4)
;; WHEN: Sun Mar 01 10:50:01 UTC 2015
;; MSG SIZE rcvd: 755

```

Figure 3: DNS Root
Priming Response with
DNSSEC Signature

```
dig +norecurse +dnssec . NS @a.root-servers.net.

; <<>> DiG 9.9.6 <<>> +norecurse +dnssec . NS @a.root-servers.net.
;; global options: +cmd
;; Got answer:
;; ->>HEADER<- opcode: QUERY, status: NOERROR, id: 52686
;; flags: qr aa; QUERY: 1, ANSWER: 14, AUTHORITY: 0, ADDITIONAL: 25
;; OPT PSEUDOSECTION:
; EDNS: version: 0, flags: do; udp: 4096
;; QUESTION SECTION:
;.                               IN      NS

;; ANSWER SECTION:
.                               518400 IN      NS      e.root-servers.net.
.                               518400 IN      NS      m.root-servers.net.
.                               518400 IN      NS      a.root-servers.net.
.                               518400 IN      NS      k.root-servers.net.
.                               518400 IN      NS      b.root-servers.net.
.                               518400 IN      NS      h.root-servers.net.
.                               518400 IN      NS      d.root-servers.net.
.                               518400 IN      NS      f.root-servers.net.
.                               518400 IN      NS      l.root-servers.net.
.                               518400 IN      NS      j.root-servers.net.
.                               518400 IN      NS      i.root-servers.net.
.                               518400 IN      NS      g.root-servers.net.
.                               518400 IN      NS      c.root-servers.net.
.                               518400 IN      RRSIG   NS 8 0 518400 20150311050000
20150301040000 16665 . 1QPFbaT+1QHnYWO6yyFvLT2JD7qddTFcRxFaolGp+CysxaZSQ
LydQtPA q3PVaKCpIkYfaFgGrOyibkkMD+nFfBxFgh/0YZN9q984NUM6LBVjpfrA MVhLy6/
qDWssDn48HoO94RwdZPzdyz+T4/KIsyH5h2FL2kp9RF1tjKlE eUU=

;; ADDITIONAL SECTION:
e.root-servers.net. 3600000 IN      A      192.203.230.10
m.root-servers.net. 3600000 IN      A      202.12.27.33
m.root-servers.net. 3600000 IN      AAAA   2001:dc3::35
a.root-servers.net. 3600000 IN      A      198.41.0.4
a.root-servers.net. 3600000 IN      AAAA   2001:503:ba3e::2:30
k.root-servers.net. 3600000 IN      A      193.0.14.129
k.root-servers.net. 3600000 IN      AAAA   2001:7fd::1
b.root-servers.net. 3600000 IN      A      192.228.79.201
b.root-servers.net. 3600000 IN      AAAA   2001:500:84::b
h.root-servers.net. 3600000 IN      A      128.63.2.53
h.root-servers.net. 3600000 IN      AAAA   2001:500:1::803f:235
d.root-servers.net. 3600000 IN      A      199.7.91.13
d.root-servers.net. 3600000 IN      AAAA   2001:500:2d::d
f.root-servers.net. 3600000 IN      A      192.5.5.241
f.root-servers.net. 3600000 IN      AAAA   2001:500:2f::f
l.root-servers.net. 3600000 IN      A      199.7.83.42
l.root-servers.net. 3600000 IN      AAAA   2001:500:3::42
j.root-servers.net. 3600000 IN      A      192.58.128.30
j.root-servers.net. 3600000 IN      AAAA   2001:503:c27::2:30
i.root-servers.net. 3600000 IN      A      192.36.148.17
i.root-servers.net. 3600000 IN      AAAA   2001:7fe::53
g.root-servers.net. 3600000 IN      A      192.112.36.4
c.root-servers.net. 3600000 IN      A      192.33.4.12
c.root-servers.net. 3600000 IN      AAAA   2001:500:2::c

;; Query time: 145 msec
;; SERVER: 198.41.0.4#53(198.41.0.4)
;; WHEN: Sun Mar 01 10:51:04 UTC 2015
;; MSG SIZE rcvd: 913
```

When the 576-octet boundary is crossed, it seems that any response up to 1,432 octets would have an equal probability to be supported by all DNS resolvers, so there is a case to be made that the 13 root name servers could be expanded to a slightly larger set of 14 or 15 servers and have the root priming response sit within 1,432 octets. But perhaps that is not quite the case, because the root priming response may yet need to grow further. The consideration here is the question of how to perform a rollover of the keys used to sign the root zone, and it would be prudent to leave some additional space in the response to allow for the use of a second digital signature. It would also be prudent to allow for a slightly larger key size, given the overall shift to longer keys over the past couple of decades of cryptography. The consequence is that it's unlikely that a further one or two root name servers would really be feasible.

Even Larger Responses?

Why not go over this limit and allow the response to be fragmented? After all, *Extension Mechanisms for DNS* (EDNS0)^[6] allows for a resolver to inform the server of the maximum-size DNS payload that it can reassemble. The EDNS0 specification suggests a size of 4,096 as a “good compromise,” and it appears that many resolvers have followed this advice.

However, as we've already noted, the problem is that when a datagram exceeds 1,500 octets, it usually has to be fragmented within the network. In IPv4, fragmentation is not uncommon, but at the same time many security firewalls regard the admission of fragments as a security risk, and the silent discarding of fragments has been observed. In IPv6 the handling of UDP fragmentation is quite different. The gateway has to send an *Internet Control Message Protocol Version 6* (ICMPv6) message back to the packet originator, and the host will then inscribe an entry in its own IP forwarding table with the revised MTU size. Subsequent UDP responses to this destination address will then use this revised MTU size, but the original response is irretrievably lost. Many aspects of this behavior are prone to error, including the use of IPv6 privacy addresses, the filtering of incoming ICMPv6 messages, the finite size of the IPv6 forwarding table, and the vulnerability of the server to spoofed ICMPv6 messages.

So while responses larger than 1,500 octets are feasible, operationally it would be prudent to limit the size of the root zone priming response to be less than 1,432 octets in IPv4. Playing it cautiously with response to packet fragmentation in IPv6 would further reduce this upper bound to 1,232 octets, because the IPv6 specification defines 1,280 as the minimum packet size that will assuredly be passed through an IPv6 network without fragmentation.

More Root Servers?

Operating just 13 root name servers is not enough for the Internet, and it has not been for many years. Adding a further 1 or 2 new root server instances was never going to change that situation, given that the demand is for many thousands of new root server instances, even if the information for these additional servers could be packed into an unfragmented root priming response.

However, one aspect of the DNS service architecture can be usefully exploited. While every resolver has to reach at least one root name server at all times, no resolver has to reach every instance of the set of root name servers at all times.

What this reality implies is that the demands of scaling the root service in the face of an expanding network can be addressed through the adoption of *anycast* clouds for many of the root name servers.

What Is “Anycast?”

Normally it is considered to be a configuration error for two or more distinct hosts to share a common IP address. However, there are times when this feature can be usefully exploited to support a highly robust service with potentially high performance. The essential prerequisite is that every instance of an anycast constellation will respond in precisely the same manner to a given input. This way it does not matter which instance of a set of host servers receives the query; the response will be the same. In an anycast scenario, the routing system essentially segments the network so that all end points that are “close” to one instance of a host within the anycast service set will have their packets directed to one host, while other end points will be directed by the routing system to use other hosts.

Anycast is most effective when using a stateless simple query/response protocol, such as DNS over UDP. However, anycast can be supported when using a TCP transport protocol, although care should be taken to ensure that the TCP sessions are relatively short in duration and that routing instability is minimized, because the redirection of packets in the middle of a TCP session to a different host in the anycast set will cause the TCP session to fail. Some operational considerations of anycast services are documented in *Operation of Anycast Services*^[7].

The anycast structure has numerous major attributes that help the root server system. The first is that the multiple instances of the root server instance split up the query load against the IP address of that server into the localities served by each anycast instance. This scenario allows the root service instance to distribute its load, improving its service. Equally, it allows the root server instance to appear to be “close” to many disparate parts of the client base simultaneously, also contributing to an improvement in its service profile. This technique also allows the root server to cope with various forms of denial-of-service attacks. Wide-scale distributed attacks are spread across multiple server instances, implying that a greater server capacity is deployed to absorb the attack.

Point attacks from a small set of sources are pinned against a single server instance, minimizing the collateral damage from such an attack to a single instance of the anycast server set.

Although anycast has considerable capability and has enjoyed operational success in recent years, there are still some problems with its operational behavior. The major problem is that this environment is still “controlled,” and each anycast instance is, in effect, operated by the anycast service root name server operator. You can’t just spin up your own instance of a root server and expect that it will engender the same level of trust in the integrity of its operation as a duly controlled and managed instance of a root service.

But why not?

Why can’t we arbitrarily expand the root service in the Internet to a level well beyond these 13 root service operators? Can we admit the concept of an “uncontrolled” root zone where anyone can offer a root zone resolution service?

Scaling the Root

There have been a couple of recent Internet Drafts on potential ways to further scale the service of the root zone that does not require the explicit permission of an existing root zone operator.

The first of these drafts is a proposal to operate a root slave service on the local loopback interface^[8] of a resolver. This approach is not an architectural change to the DNS (or at least not intentionally). For recursive resolvers that implement this approach, this approach is a form of change in query behavior because a recursive resolver so configured will no longer query the root servers, but instead direct these queries to a local instance of a slave server that is listening on the recursive resolver loopback address. This slave server is serving a locally held instance of the root zone, and the recursive resolver would perform DNSSEC validation of responses from this local slave to ensure the integrity of responses received in this manner. For users of this recursive resolver, there is no apparent change to the DNS or to their local configurations.

The motivation behind this proposal is that a population of recursive resolvers is still too far away from all of the root servers, and this situation causes delays in the DNS resolution function. The caching properties of recursive DNS resolvers is such that the overall majority of queries directed to the root servers are for nonexistent top-level domains, so a pragmatic restatement of the problem space is that there are recursive resolvers that take too long to generate a *Non-existent Domain Name* (NXDOMAIN) response, and this approach would reduce this time delay.

However, given this particular formulation of the problem space, then the larger and more comprehensive the anycast constellations of the root servers, the less the demand for this particular approach.

Locales where there are adequately close DNS root services from the anycast root servers would find no particular advantage in operating a local slave DNS root server because the marginal speed differential may not be an adequate offset for the added complexity of configuration and operation of the local slave server.

The linking of the root zone information to the loopback is a point of fragility in the setup. Setting up a slave DNS server that is authoritative for the root zone would require using multiple root servers to ensure that it has access to a root zone from at least one of the anycast server constellations at any time. If at any time it cannot retrieve a master copy of the root zone, it should respond with a SERVFAIL (server failure) code, and the local recursive resolver should interpret this response as a signal to revert to conventional queries against the root servers.

This proposal provides integrity in the local root server through the mechanism of having the recursive resolver perform DNSSEC validation against the responses received from the local root slave. If the recursive resolver is configured as a DNSSEC-validating resolver, then it is configurable on current implementations of DNS recursive resolvers. However, if it is desired to limit DNSSEC validation to just the responses received from the local slave root server, then this configuration is not within the current capabilities of the more widely used DNS resolver implementations today.

The advantages of this approach is that the decision to set up a local slave root server is one that is entirely local to the recursive resolver, and the impacts of this decision affect only the clients of this recursive resolver. No coordination with the root server operators is required, nor is any explicit notification. The local slave server is only indirectly visible to the clients of this recursive resolver and no other.

A second proposal is slightly more conventional in that it proposes adding a new anycast root server constellation to the DNS root, but instead of adding a new entry to the existing root server set, it proposes a second root server hints file.^[9]

One possible motivation behind this proposal lies in the observation that before the root was DNSSEC-signed, we placed much reliance in the concept that the root zone was served from only a small number of IP addresses, but when the root zone is DNSSEC-signed, then the integrity of the responses generated from a root zone is based on the ability of the receiver of the response to validate the signed responses using its local copy of the root keys. Who serves the zone in a DNSSEC context is largely irrelevant.

This approach uses the same root zone signing key to sign a second root zone, where the root zone is served by an exclusively assigned anycast address set. This second set of addresses would not be exclusively assigned to any root server operators, but allowed to be used by any party, in a form of uncontrolled anycast.

In some ways this proposal is similar to the existing AS112 work^[10], where anyone can set up a server to respond to common queries for nonexistent top-level domains (such as `.local`) with NXDOMAIN responses. The implication is that if there is a perception that a locale is poorly served by the existing root server anycast constellations, then a local instance of this particular anycast root server can be set up. Because the address is specifically dedicated for unowned anycast, there is no need to coordinate with either the existing root server operators or with other operators of root zone servers on the same anycast address. A prospective operator of one of these root servers simply serves the unowned anycast root zone from one of the small pool (two address couplets, each linking an IPv4 and an IPv6 address) of reserved anycast addresses. Recursive resolvers would query this server by using a distinct unowned anycast root hints file.

Why would any recursive resolver trust the veracity of responses received from one of these unowned anycast root servers? We could well ask the same of recursive resolvers who query into any of the 13 anycast constellations for the root zone, and to some extent the risks of being led astray are similar. The one mitigation of the latter case is that hijacking an existing anycast constellation prefix requires the coercive corruption of the routing system to inject a false instance of an “owned” address, although there is no such concept as hijacking of the unowned anycast prefix. However, in both cases some skepticism on the part of the recursive resolver is to be encouraged, and recursive resolvers should be motivated to validate all such responses using DNSSEC, using the local copy of the DNS trust anchor material. This practice is still a part of “good housekeeping” recommended operational practice for recursive resolvers using the existing root servers, but is a more strongly worded requirement for resolvers using this unowned anycast service.

Although it could be regarded as a byproduct of a single hierarchical name space, the centralization of root zone information in the DNS is operationally problematical and does not cleanly fit within a distributed and decentralized peer model of a network architecture. The adoption of root server anycast constellations is an attempt to respond to this situation, to an extent, by overprovisioning of this critical service. A similar picture is emerging in the area of content provisioning, where cloud-based content is essentially an exercise in overprovisioning, where single points of distribution are replaced by a larger scale of multiple delivery points in order to improve the quality of the delivered service.

However, end users don’t enjoy the same level of control, and are dependent on external conditions that are effectively out of their direct control. Not only does this dependence result in highly variable service experiences, but it also leaves the user highly exposed to various forms of online surveillance. The distributing computing world has created external dependencies for users where access to local service is reliant on external availabilities.

The DNS is a good example of this scenario, in so far as resolution of a DNS name that is not already contained in local caches requires priming queries to external servers. Obviously these dependencies on external services highlight fragility where local services cannot be reliably provided using only local infrastructure.

Both of these proposals are incremental in nature, and propose a form of augmentation to the existing structure of recursive resolvers and the root name server, rather than any fundamental change to the existing structure. In so doing, there is a distinct possibility that this form of uncoordinated piecemeal expansion at a local level could prove to be more effective across the Internet, and the critical role of the existing 13 root server operators would diminish over time if it were.

Neither of these proposed approaches is entirely without some form of change. All these uncoordinated root server operators would mean that push notification of root zone changes via the NOTIFY message would not be not feasible, so it would be back to periodic zone transfers and timers in the root zone headers. There is no longer a quick mistake-correction capability in the root zone if served in this way, although it could be argued that the massive level of caching of DNS information actually implied that any changes in the root zone were subject to cache flushing in any case, irrespective of the speed of zone change at the level of the root server anycast constellations.

What is perhaps more worrisome is that the unowned anycast proposal is in effect a proposal to fork the root zone, and recursive resolvers are forced to position themselves within one regime or the other. The only common glue left in this environment is the root key, because the only way that a client of either regime can detect that it is receiving genuine answers is to perform response validation using the root key. This type of validation is placing a massive amount of invested trust in a security artifact that is used today by only a very small subset of recursive resolvers.

It also needs to be noted that our experience to date with unowned anycast has been very poor. At one stage the IPv6 transition experimented with a form of unowned anycast in the form of 6to4 tunnel servers, and the results were hardly reassuring. Anycast clouds follow routing, not geography, and diagnosing operational failures that occur within an uncoordinated anycast structure can range from the merely challenging to highly cryptic and insoluble. The relationship between a recursive resolver and the actual root server it is querying is then occluded, and instances of structural failure in DNS name resolution are far harder to diagnose and correct. Considering that the name translation function is an essential foundation for the Internet, adding operational opacity to the root zone query function is not a step that should be taken lightly.

But that reality does not imply that the other proposal is free from operational concerns, either. The complexity of the local slave resolver, with two concurrent DNS resolvers operating within the same host, should also be questioned. Although there is a current convention in DNS resolver and server deployment to avoid the model of a “mixed” mode resolver that is both an authoritative (or slave) server for some domains and a recursive resolver for all other domains, this avoidance is perhaps nothing more than a convention, and it seems overkill for a resolver to phrase a root query and reach through the loopback interface to ask a co-resident DNS resolver the queries that are being posed to the root.

Why not just operate the local resolver in mixed mode and allow the root query to simply become a memory lookup within a single DNS resolver instance? Perhaps the only justification in this case is the issue of root zone integrity. In the mixed mode of a single resolver instance, the question that arises is a reasonable one: How can the local resolver validate the contents of the transferred root zone? In the loopback model, the local resolver performs DNSSEC validation of the root zone responses and therefore does not necessarily need to separately validate the contents of the transferred root zone. In theory a mixed-mode resolver could DNSSEC-validate the responses retrieved from its local instance of a slave zone server before passing them to the recursive resolver function, but it’s not clear that any existing DNS resolver implementations perform this form of DNSSEC validation of internal queries in a mixed mode of operation. Alternate approaches of including a zone signature to ensure integrity are also a possibility to ensure that the recursive resolver is not placed into a position of inadvertently serving corrupt root zone data.

Why not take this thought a further step, and allow any recursive resolver to be a slave server for the root zone? If the zone transfer function included an integrity check across the entire transferred zone (such as a hash of the transferred zone, signed by the root zone signing key), then the recursive resolver could be assured that it was then serving an authentic copy of the root zone.

However, such an integrity check on the transferred zone only assists the local recursive resolver in assuring itself that it has obtained an authentic and current copy of the zone. Clients of that recursive resolver should be as skeptical as ever and ask for DNSSEC signatures, and perform DNSSEC validation over all signed responses that are received from the recursive resolver. The advantage of this approach is that it permits essentially an unlimited number of root name servers, where every recursive server that wants to can serve its own validated copy of the root zone. The disadvantage is that, like all large distributed systems, there is some introduced inertia into the system and updates take time to propagate. However, in a space that is already highly cached, the difference between what happens today and what would happen in such a scenario may well be very hard to see.

There are other ways that a recursive resolver can authoritatively serve responses from the root zone that avoids an explicit root zone transfer, yet still primes the recursive resolver with authoritative information about the contents of the root zone. The response to a query for a nonexistent domain provides NSEC responses that allow a resolver to construct a local cache of the entire root zone (Figure 4).

In the example shown in Figure 4, the response from the root server for the top-level name **nosuchdomain.** indicates in the signed NSEC record that was returned that all names that lie between **no.** and **np.** can be correctly interpreted to be nonexistent domains. As long as the resolver can successfully validate the digital signatures contained in this response, the resolver can cache this negative result and serve NX domain responses for all names in this range for the lifetime of the cache.

Figure 4: Example of an NSEC Response

```
dig +dnssec foo.bar.nosuchdomain @a.root-servers.net

; <<>> DiG 9.9.6 <<>> +dnssec foo.bar.nosuchdomain @a.root-servers.net
;; global options: +cmd
;; Got answer:
;; ->>HEADER<<- opcode: QUERY, status: NXDOMAIN, id: 18580
;; flags: qr aa rd; QUERY: 1, ANSWER: 0, AUTHORITY: 6, ADDITIONAL: 1
;; WARNING: recursion requested but not available

;; OPT PSEUDOSECTION:
; EDNS: version: 0, flags: do; udp: 4096
;; QUESTION SECTION:
;foo.bar.nosuchdomain.      IN      A

;; AUTHORITY SECTION:
.                86400    IN      SOA      a.root-servers.net. nstld.verisign-grs.com.
2015022200 1800 900 604800 86400
.                86400    IN      RRSIG     SOA 8 0 86400 20150304050000 20150222040000 16665 .
E7no0qtMyyVdVH/0t5LQOM+xV8VJB5GwWp6oaphV+63gi9Dj8LG71kb8 N00Sx0TaJAISa18NLa27/RPzoz3vvQAnIpyZxmhxzfkyk
fkLhXxaJtFCV 4hKWxqf0EymCzGCsBIRSMttl7fypf3aml5JF3ei0Cqmp/BHWjXjGs0mO te8=
.                86400    IN      NSEC      abogado. NS SOA RRSIG NSEC DNSKEY
.                86400    IN      RRSIG     NSEC 8 0 86400 20150304050000 20150222040000 16665 .
ojA/fqJKig89aw9+KtM2RswgMaxTVrogPiGeqoLUZgD9Rf3UN0n2tLtO VpDzzB45BquRpdfV+OludWo+L8lnRqjM5CAiZkaektUW
MmOcvqChFf9w RuiqMdAq4vVExSWZU4G/af6Y6WzoHVC5G1WS31PHfm9Ux2UeyeEMQ1o/ aac=
no.              86400    IN      NSEC      np. NS DS RRSIG NSEC
no.              86400    IN      RRSIG     NSEC 8 1 86400 20150304050000 20150222040000 16665 .
K7dypmjhxfdGtP5oWSvToH53qYdfcGi0MQ7xkF+/k89HfBZnFtOrhGd/ 8zJAdUtednJxvk55LjHY+uU4uiaNuNc+7jAilFEKL8BF
sgzvy98lvgkk 63YqAoRHJJ357q+hVil0UxVKcb8oCe3VWZsOtamap6ujSzZZQy4X3TKt jZ0=

;; Query time: 146 msec
;; SERVER: 198.41.0.4#53(198.41.0.4)
;; WHEN: Sun Feb 22 08:28:24 UTC 2015
;; MSG SIZE rcvd: 654
```

Conclusions

DNSSEC is indeed an extremely valuable asset for the DNS, and we can move forward with a larger and more robust system if we can count on various efforts to subvert the operation DNS being thwarted by all forms of clients of the DNS, even to the level of applications insisting that they are exposed to the DNS responses and their signatures, and performing their own validation on the received data.

There is a more general observation about scaling in the Internet. We are finding it increasingly challenging to react to inexorable pressures of scaling the infrastructure of the Internet while still maintaining basic backward compatibility with systems that conform only to technical standards of the 1980s. We need to move on. We have moved beyond the 512-octet packet limit, and these days it's reasonable to assume that useful resolvers can support EDNS0 options in DNS queries. Resolvers should perform DNSSEC validation. Wondering why there are 13 available slots for root name server operators, and wondering about how we could alter their composition, or change the interaction with the DNS root to support one or two additional root servers, are unproductive lines of investigation at so many levels. It may well be time to contemplate a different DNS that does not involve these arbitrary constraints over the number and composition of players that serve the signed root zone. Instead we should rely on using DNSSEC validation to ensure that the responses received from queries to the root zone are authentic.

The attraction of unconstrained systems is that local actors can respond to local needs, and respond by using local resources, without having to coordinate or cross-subsidize the activities of others. Much of the momentum of the Internet is directed by this loosely constrained model of interaction. If we can use the possibilities opened up by securing the DNS payload, where the question of who passed the DNS information to you is irrelevant but the question of whether the information is locally verifiable is critically important, then and only then can we contemplate what it would take to operate an unconstrained DNS system for serving the root zone.

References

- [1] R. Braden, “Requirements for Internet Hosts – Communication Layers,” RFC 1122, October 1989.
- [2] P. Mockapetris, “Domain names – Implementation and Specification,” RFC 1035, November 1987.
- [3] The DIG Command:
[http://en.wikipedia.org/wiki/Dig_\(command\)](http://en.wikipedia.org/wiki/Dig_(command))
- [4] R. Braden, “Requirements for Internet Hosts – Application and Support,” RFC 1123, October 1989.
- [5] R. Arends, et.al, “Protocol Modifications for the DNS Security Extensions,” RFC 4035, March 2005.
- [6] P. Vixie, J. Damas, and M. Graff, “Extension Mechanisms for DNS (EDNS0),” RFC 6891, April 2013.
- [7] K. Lindqvist and J. Abley, “Operation of Anycast Services,” RFC 4786, December 2006.
- [8] W. Kumhari and P. Hoffman, “Decreasing Access Time to Root Servers by Running One on Loopback,” work in progress, (Internet Draft: **draft-wkumari-dnsop-root-loopback-02.txt**), November 2014.
- [9] X. Lee and P. Vixie, “How to Scale the DNS Root System?” work in progress (Internet Draft: **draft-lee-dnsop-scalingroot-00.txt**), July 2014.
- [10] J. Abley and W. Maton, “AS112 Nameserver Operations,” RFC 6304, July 2011.
- [11] W. Stallings, “Gigabit Ethernet: From 1 to 100 Gbps and Beyond,” *The Internet Protocol Journal*, Volume 18, No. 1, March 2015.
- [12] G. Huston, “A Question of DNS Protocols,” *The Internet Protocol Journal*, Volume 17, No. 1, September 2014.
- [13] E. Feinler, “Host Tables, Top-Level Domain Names, and the Origin of Dot Com,” *IEEE Annals of the History of Computing*, July-September 2011, <http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=5986499>

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for the initial build of the Internet within the Australian academic and research sector. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001.

E-mail: gih@apnic.net

Gigabit Ethernet: From 1 to 100 Gbps and Beyond

by William Stallings

Since its first introduction in the early 1980s, Ethernet has been the dominant technology for implementing *Local-Area Networks* (LANs) in office environments. Over the years, the data-rate demands placed on LANs have grown at a rapid pace. Fortunately, Ethernet technology has adapted to provide ever-higher capacity to meet these needs. We are now in the era of the Gigabit Ethernet.

Ethernet began as an experimental bus-based 2.94-Mbps system^[1] using coaxial cable. In a shared bus system, all stations attach to a common cable, with only one station able to successfully transmit at a time. A *Medium Access Control* (MAC) protocol based on collision detection arbitrates the use of the bus. In essence, each station is free to transmit MAC frames on the bus. If a station detects a collision during transmission, it backs off a certain amount of time and tries again.

The first commercially available Ethernet products were bus-based systems operating at 10 Mbps^[2]. This introduction coincided with the standardization of Ethernet by the IEEE 802.3 committee. With no change to the MAC protocol or MAC frame format, Ethernet could also be configured in a star topology, with traffic going through a central hub, again with transmission limited to a single station at a time through the hub. To enable an increase in the data rate, a switch replaces the hub, allowing full-duplex operation. With the switch, the same MAC format and protocol are used, although collision detection is no longer needed. As the demand has evolved and the data rate requirement increased, some enhancements to the MAC layer have been added, such as provision for larger frame sizes.

Currently, Ethernet systems are available at speeds up to 100 Gbps. Table 1 summarizes the successive generations of IEEE 802.3 standardization.

Table 1: IEEE 802.3 Physical Layer Standards

Year Introduced	1983	1995	1998	2003	2010
Maximum data transfer speed	10 Mbps	100 Mbps	1 Gbps	10 Gbps	40 Gbps, 100 Gbps
Transmission media	Coax cable, unshielded twisted pair, optical fiber	Unshielded twisted pair, shielded twisted pair, optical fiber	Unshielded twisted pair, optical fiber, shielded twisted pair	Optical fiber	Optical fiber, backplane

Ethernet quickly achieved widespread acceptance and soon became the dominant technology for LANs. Its dominance has since spread to *Metropolitan-Area Networks* (MANs) and a wide range of applications and environments.

The huge success of Ethernet is due to its extraordinary adaptability. The same MAC protocol and frame format are used at all data rates. The differences are at the physical layer, in the definition of signaling technique and transmission medium.

In the remainder of this article, we look at characteristics of Ethernet in the gigabit range.

1-Gbps Ethernet

For many years the initial standard of Ethernet, at 10 Mbps, was adequate for most office environments. By the early 1990s, it was clear that higher data rates were needed to support the growing traffic load on the typical LAN. Key drivers in this evolution include:

- *Centralized Server Farms*: In many multimedia applications, there is a need for the client system to be able to draw huge amounts of data from multiple, centralized servers, called *Server Farms*. As the performance of the servers has increased, the network has become the bottleneck.
- *Power Workgroups*: These groups typically consist of a small number of cooperating users who need to exchange massive data files across the network. Example applications are software development and computer-aided design.
- *High-speed Local Backbone*: As processing demand grows, enterprises develop a configuration of multiple LANs interconnected with a high-speed backbone network.

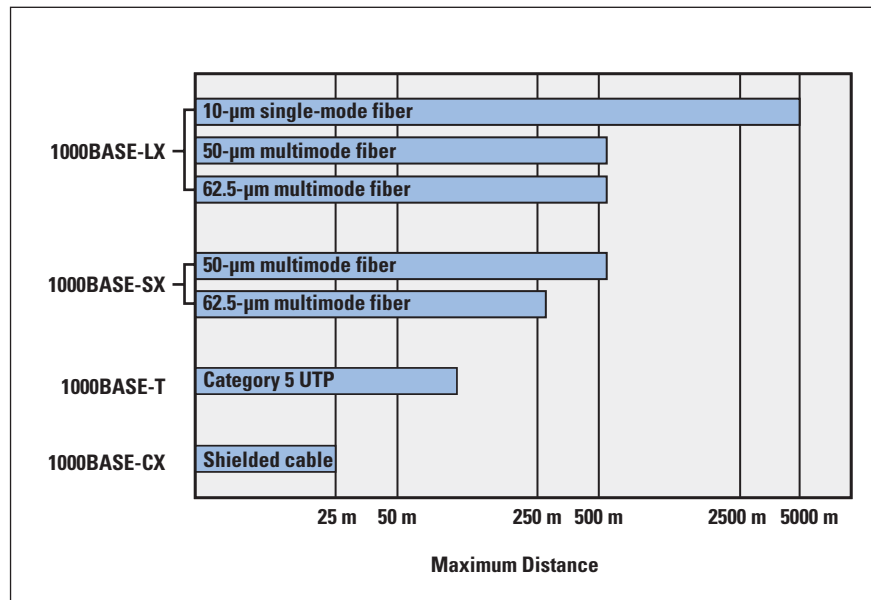
To meet such needs, the IEEE 802.3 committee developed a set of specifications for Ethernet at 100 Mbps, followed a few years later by a 1-Gbps family of standards. In each case, the new specifications defined transmission media and transmission encoding schemes built on the basic Ethernet framework, making the transition easier than if a completely new specification were issued.

The 1-Gbps standard includes a variety of transmission medium options^[3, 4] (Figure 1):

- *+1000BASE-SX*: This short-wavelength option supports duplex links of up to 275 m using 62.5- μ m multimode or up to 550 m using 50- μ m multimode fiber. Wavelengths are in the range of 770 to 860 nm.
- *1000BASE-LX*: This long-wavelength option supports duplex links of up to 550 m of 62.5- μ m or 50- μ m multimode fiber or 5 km of 10- μ m single-mode fiber. Wavelengths are in the range of 1270 to 1355 nm.

- **1000BASE-CX:** This option supports 1-Gbps links among devices located within a single room or equipment rack, using copper jumpers (specialized shielded twisted-pair cable that spans no more than 25 m). Each link is composed of a separate shielded twisted pair running in each direction.
- **1000BASE-T:** This option uses four pairs of Category 5 unshielded twisted pair to support devices over a range of up to 100 m, transmitting and receiving on all four pairs at the same time, with echo cancellation circuitry.

Figure 1: 1-Gbps Ethernet Medium Options (log scale)



The signal encoding scheme used for the first three Gigabit Ethernet options just listed is 8B/10B. With 8B/10B, each 8 bits of data is converted into 10 bits for transmission^[3]. The extra bits serve two purposes. First, the resulting signal stream has more transitions between logical 1 and 0 than an uncoded stream; it avoids the possibility of long strings of 1s or 0s that make synchronization between transmitter and receiver more difficult. Second, the code is designed in such a way as to provide a useful error-detection capability.

The signal-encoding scheme used for 1000BASE-T is 4D-PAM5, a complex scheme whose description is beyond our scope.

In a typical application of Gigabit Ethernet, a 1-Gbps LAN switch provides backbone connectivity for central servers and high-speed workgroup Ethernet switches. Each workgroup LAN switch supports both 1-Gbps links, to connect to the backbone LAN switch and to support high-performance workgroup servers, and 100-Mbps links, to support high-performance workstations, servers, and 100-Mbps LAN switches.

10-Gbps Ethernet

Even as the ink was drying on the 1-Gbps specification, the continuing increase in local traffic made this specification inadequate for needs in the short-term future. Accordingly, the IEEE 802.3 committee soon issued a standard for 10-Gbps Ethernet. The principal requirement for 10-Gbps Ethernet was the increase in intranet (local interconnected networks) and Internet traffic. Numerous factors contribute to the explosive growth in both Internet and intranet traffic:

- An increase in the number of network connections
- An increase in the connection speed of each end station (for example, 10-Mbps users moving to 100 Mbps, and analog 56k users moving to DSL and cable modems)
- An increase in the deployment of bandwidth-intensive applications such as high-quality video
- An increase in Web hosting and application hosting traffic

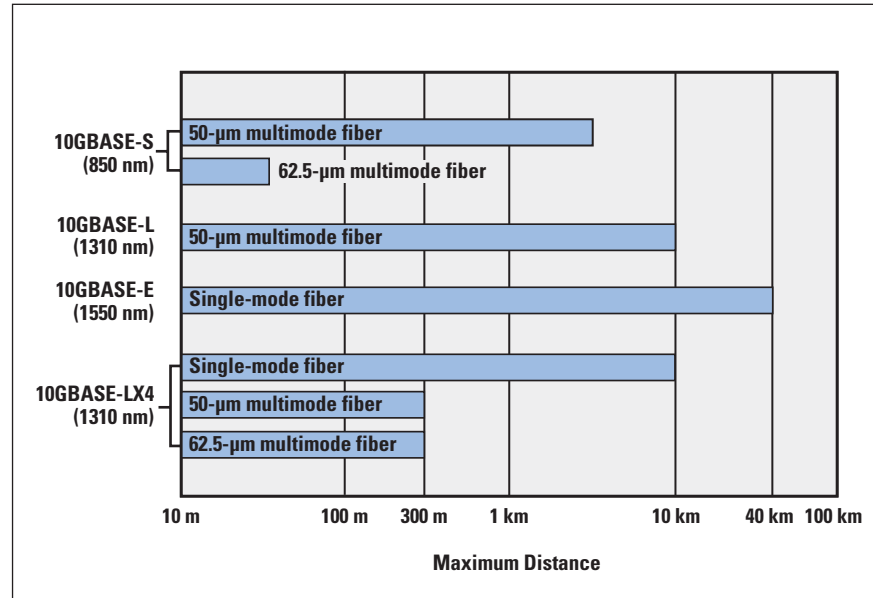
Initially, network managers used 10-Gbps Ethernet to provide high-speed, local backbone interconnection between large-capacity switches. As the demand for bandwidth increased, 10-Gbps Ethernet began to be deployed throughout the entire network, to include server farm, backbone, and campus-wide connectivity. This technology enables *Internet Service Providers* (ISPs) and *Network Service Providers* (NSPs) to create very-high-speed links at a very low cost, between co-located, carrier-class switches, and routers^[5].

The technology also allows the construction of MANs and *Wide-Area Networks* (WANs) that connect geographically dispersed LANs between campuses or *Points of Presence* (PoPs). Thus, Ethernet begins to compete with *Asynchronous Transfer Mode* (ATM) and other wide-area transmission and networking technologies. In most cases where the customer requirement is data and TCP/IP transport, 10-Gbps Ethernet provides substantial value over ATM transport for both network end users and service providers:

- No expensive, bandwidth-consuming conversion between Ethernet packets and ATM cells is required; the network is Ethernet, end to end.
- The combination of IP and Ethernet offers *Quality of Service* (QoS) and traffic policing capabilities that approach those provided by ATM, so that advanced traffic-engineering technologies are available to users and providers.
- A wide variety of standard optical interfaces (wavelengths and link distances) have been specified for 10 Gigabit Ethernet, optimizing its operation and cost for LAN, MAN, or WAN applications.

The goal for maximum link distances cover a range of applications is from 300 m to 40 km. The links operate in full-duplex mode only, using a variety of optical fiber physical media.

Figure 2: 10-Gbps Ethernet Distance Options (log scale)



Four physical layer options are defined for 10-Gbps Ethernet (Figure 2). The first three have two suboptions: an “R” suboption and a “W” suboption. The R designation refers to a family of physical layer implementations that use a signal encoding technique known as 64B/66B. With 64B/66B, each 64 bits of data is converted into 66 bits for transmission, resulting in substantially less overhead than 8B/10B but providing the same types of benefits. The R implementations are designed for use over *dark fiber* meaning a fiber-optic cable that is not in use and that is not connected to any other equipment. The W designation refers to a family of physical layer implementations that also use 64B/66B signaling but are then encapsulated to connect to SONET equipment.

The four physical layer options are as follows:

- **10GBASE-S (*short*):** Designed for 850-nm transmission on multimode fiber. This medium can achieve distances up to 300 m. There are 10GBASE-SR and 10GBASE-SW versions.
- **10GBASE-L (*long*):** Designed for 1310-nm transmission on single-mode fiber. This medium can achieve distances up to 10 km. There are 10GBASE-LR and 10GBASE-LW versions.
- **10GBASE-E (*extended*):** Designed for 1550-nm transmission on single-mode fiber. This medium can achieve distances up to 40 km. There are 10GBASE-ER and 10GBASE-EW versions.
- **10GBASE-LX4:** Designed for 1310-nm transmission on single-mode or multimode fiber. This medium can achieve distances up to 10 km. This medium uses *Wavelength-Division Multiplexing* (WDM) to multiplex the bit stream across four light waves.

40-/100-Gbps Ethernet

Ethernet is widely deployed and is the preferred technology for wired local-area networking. Ethernet dominates enterprise LANs, broadband access, and data center networking, and it has also become popular for communication across MANs and even WANs. Further, it is now the preferred carrier wire-line vehicle for bridging wireless technologies, such as Wi-Fi and WiMAX, into local Ethernet networks.

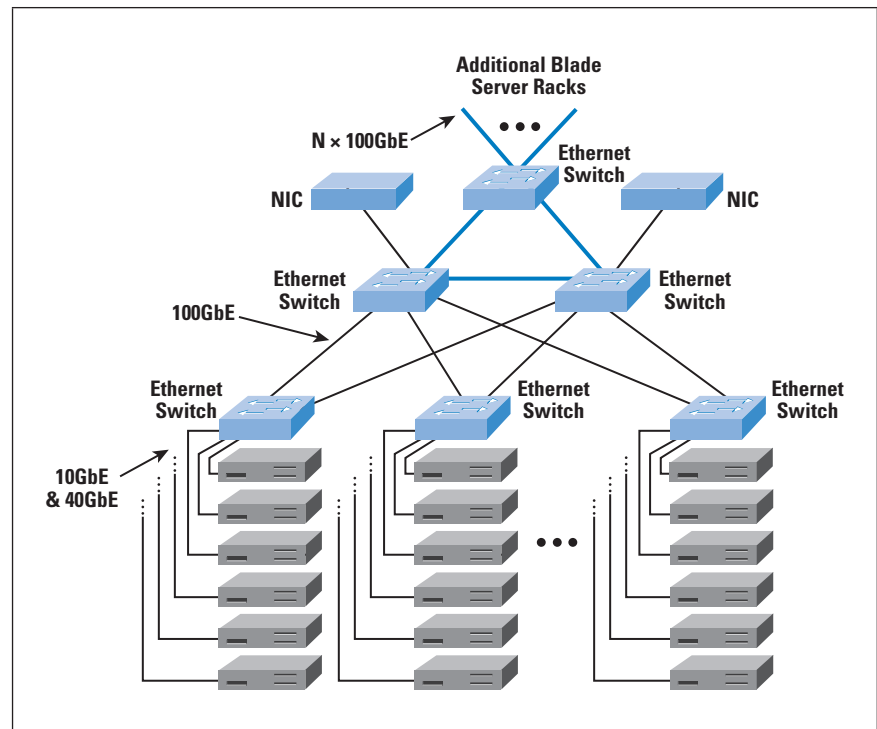
This popularity of Ethernet technology is due to the availability of cost-effective, reliable, and interoperable networking products from a variety of vendors. The development of converged and unified communications, the evolution of massive server farms, and the continuing expansion of *Voice over IP* (VoIP), *Television over IP* (TVoIP), and Web 2.0 applications have accelerated the need for ever-faster Ethernet switches. The following are market drivers for 100-Gbps Ethernet:

- *Data center/Internet Media Providers:* To support the growth of Internet multimedia content and Web applications, content providers have been expanding data centers, pushing 10-Gbps Ethernet to its limits. These providers are likely to be high-volume early adopters of 100-Gbps Ethernet.
- *Metro-Video/Service Providers:* Video on demand has been leading a new generation of 10-Gbps Ethernet metropolitan/core network build-outs. These providers are likely to be high-volume adopters in the medium term.
- *Enterprise LANs:* Continuing growth in convergence of voice/video/data and in unified communications is accelerating network switch demands. However, most enterprises still rely on 1-Gbps or a mix of 1-Gbps and 10-Gbps Ethernet, and adoption of 100-Gbps Ethernet is likely to be slow.
- *Internet exchanges/ISP Core Routing:* With the massive amount of traffic flowing through these nodes, these installations are likely to be early adopters of 100-Gbps Ethernet.

In 2007, the IEEE 802.3 working group authorized the *IEEE P802.3ba 40-Gbps and 100-Gbps Ethernet Task Force*. The 802.3ba project authorization request cited numerous examples of applications that require greater data-rate capacity than 10-Gbps Ethernet offers, including Internet exchanges, high-performance computing, and video-on-demand delivery. The authorization request justified the need for two different data rates in the new standard (40 Gbps and 100 Gbps) by recognizing that aggregate network requirements and end-station requirements are increasing at different rates.

An example of the application of 100-Gbps Ethernet is shown in Figure 3. The trend at large data centers, with substantial banks of blade servers, is the deployment of 10-Gbps ports on individual servers to handle the massive multimedia traffic provided by these servers. Typically, a single blade-server rack will contain multiple servers and one or two 10-Gbps Ethernet switches to interconnect all the servers and provide connectivity to the rest of the facility. The switches are often mounted in the rack and referred to as *Top-of-Rack* (ToR) switches. The term ToR has become synonymous with a server access switch, even if it is not located “top of rack.” For very large data centers, such as cloud providers, the interconnection of multiple blade-server racks with additional 10-Gbps switches is increasingly inadequate. To handle the increased traffic load, switches operating at greater than 10 Gbps are needed to support the interconnection of server racks and to provide adequate capacity for connecting off-site through *Network Interface Controllers* (NICs).

Figure 3: Example 100-Gbps Ethernet Configuration for Massive Blade-Server Cloud Site



The first products in this category appeared in 2009, and the IEEE 802.3ba standard was finalized in 2010. Initially, many enterprises are deploying 40-Gbps switches, but both 40- and 100-Gbps switches are projected to enjoy increased market penetration in the next few years^[6, 7, 8].

IEEE 802.3ba specifies three types of transmission media as shown in Table 2: copper backplane, twin axial (a type of cable similar to coaxial cable), and optical fiber. For copper media, four separate physical lanes are specified. For optical fiber, either 4 or 10 wavelengths are specified, depending on data rate and distance^[9, 10, 11].

Table 2: Media Options for 40- and 100-Gbps Ethernet

	40 Gbps	100 Gbps
1-m backplane	40GBASE-KR4	
10-m copper	40GBASE-CR4	100GBASE-CR10
100-m multimode fiber	40GBASE-SR4	100GBASE-SR10
10-km single-mode fiber	40GBASE-LR4	100GBASE-LR4
40-km single-mode fiber		100GBASE-ER4

Naming nomenclature:

Copper: K = Backplane; C = Cable assembly

Optical: S = Short reach (100 m); L - Long reach (10 km);

E = Extended long reach (40 km)

Coding scheme: R = 64B/66B block coding

Final number: Number of lanes (copper wires or fiber wavelengths)

Multilane Distribution

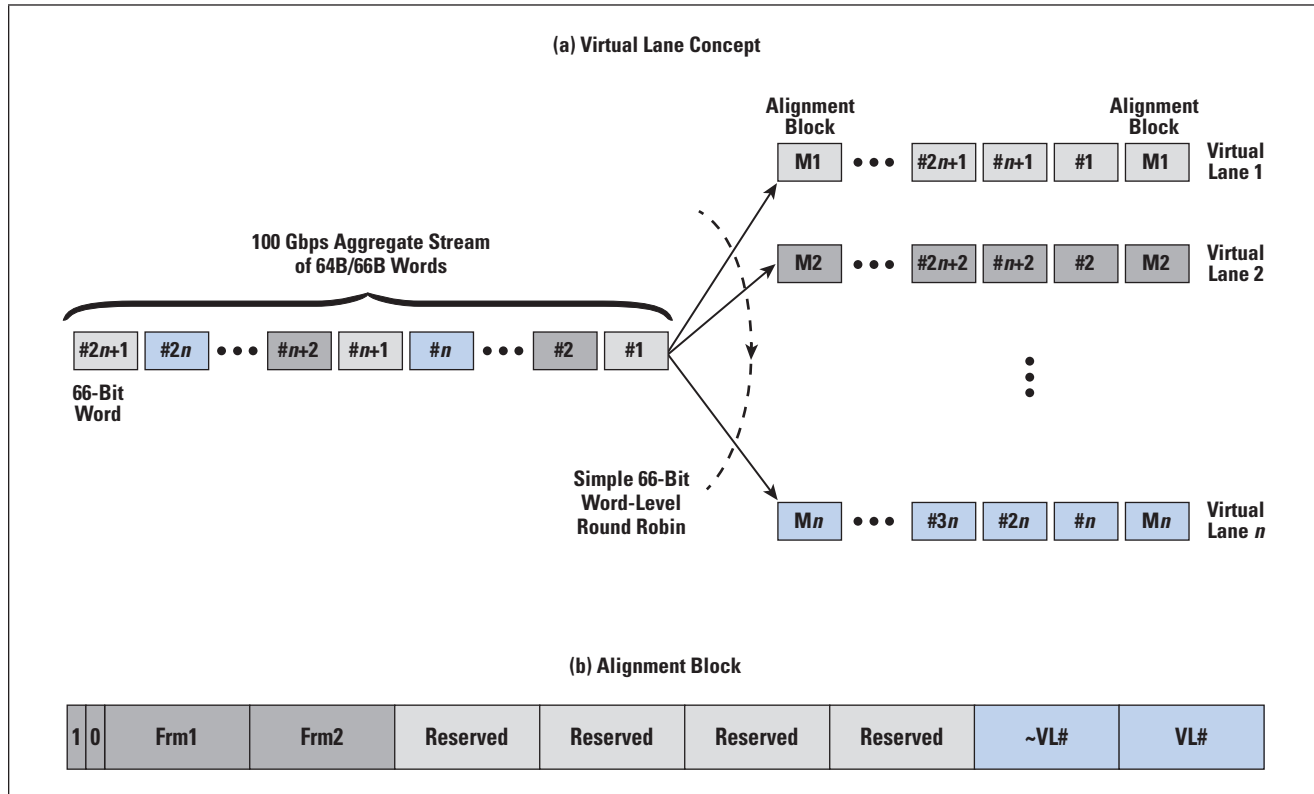
The 802.3ba standard uses a technique known as multilane distribution to achieve the required data rates. Two separate concepts need to be addressed: multilane distribution and virtual lanes.

The general idea of *multilane distribution* is that, in order to accommodate the very high data rates of 40 and 100 Gbps, the physical link between an end station and an Ethernet switch or the physical link between two switches may be implemented as multiple parallel channels. These parallel channels could be separate physical wires, such as four parallel twisted-pair links between nodes. Alternatively, the parallel channels could be separate frequency channels, such as provided by WDM over a single optical fiber link.

For simplicity and manufacturing ease, we would like to specify a specific multiple-lane structure in the electrical physical sublayer of the device, known as the *Physical Medium Attachment* (PMA) sublayer. The lanes produced are referred to as *virtual lanes*. If a different number of lanes are actually in use in the electrical or optical link, then the virtual lanes are distributed into the appropriate number of physical lanes in the *Physical Medium Dependent* (PMD) sublayer. This is a form of inverse multiplexing.

Figure 4a shows the virtual lane scheme at the transmitter. The user data stream is encoded using the 64B/66B, which is also used in 10-Gbps Ethernet. Data is distributed to the virtual lanes one 66-bit word at a time using a simple round robin scheme (first word to first lane, second word to second lane, etc.). A unique 66-bit alignment block is added to each virtual lane periodically. The alignment blocks are used to identify and reorder the virtual lanes and thus reconstruct the aggregate data stream.

Figure 4: Multilane Distribution for 100-Gbps Ethernet



The virtual lanes are then transmitted over physical lanes. If the number of physical lanes is smaller than the number of virtual lanes, then bit-level multiplexing is used to transmit the virtual lane traffic. The number of virtual lanes must be an integer multiple (1 or more) of the number of physical lanes.

Figure 4b shows the format of the alignment block. The block consists of 8 single-byte fields preceded by the 2-bit synchronization field, which has the value 10. The Frm fields contain a fixed framing pattern common to all virtual lanes and used by the receiver to locate the alignment blocks. The VL# fields contain a pattern unique to the virtual lane: one of the fields is the binary inverse of the other.

25-/50-Gbps Ethernet

One of the options for implementing 100 Gbps is as four 25-Gbps physical lanes. Thus, it would be relatively easy to develop standards for 25- and 50-Gbps Ethernet, using one or two lanes, respectively. Having these two lower-speed alternatives, based on the 100-Gbps technology, would give users more flexibility in meeting existing and near-term demands with a solution that would scale easily to higher data rates.

Such considerations have led to the formation of the *25 Gigabit Ethernet Consortium* by numerous leading cloud networking providers, including Google and Microsoft. The objective of the consortium is to support an industry-standard, interoperable Ethernet specification that boosts the performance and slashes the interconnect cost per Gbps between the NIC and ToR switch. The specification adopted by the consortium prescribes a single-lane 25-Gbps Ethernet and dual-lane 50-Gbps Ethernet link protocol, enabling up to 2.5 times higher performance per physical lane on twinax copper wire between the rack endpoint and switch compared to current 10- and 40-Gbps Ethernet links. The IEEE 802.3 committee is presently developing the needed standards for 25 Gbps, and it may include 50 Gbps^[12, 13].

It is too early to say how these various options (25, 40, 50, and 100 Gbps) will play out in the marketplace. In the intermediate term, the 100-Gbps switch is likely to predominate at large sites, but the availability of these slower and cheaper alternatives gives enterprises numerous paths for scaling up to meet increasing demand.

400-Gbps Ethernet

The growth in demand never lets up. IEEE 802.3 is currently exploring technology options for producing a 400-Gbps Ethernet standard, although no timetable is yet in place^[14, 15, 16, 17]. Looking beyond that milestone, there is widespread acknowledgment that a 1-Tbps (terabit per second, trillion bits per second) standard will eventually be produced^[18].

2.5-/5-Gbps Ethernet

As a testament to the versatility and ubiquity of Ethernet, and at the same time the fact that ever-higher data rates are being standardized, consensus is developing to standardize two lower rates: 2.5 and 5 Gbps^[19, 20]. These relatively low speeds are also known as *Multirate Gigabit BASE-T* (MGBASE-T). Currently, the *MGBASE-T Alliance* is overseeing the development of these standards outside of IEEE. It is likely that the IEEE 802.3 committee will ultimately issue standards based on these industry efforts.

These new data rates are intended mainly to support IEEE 802.11ac wireless traffic into a wired network. IEEE 802.11ac is a 3.2-Gbps Wi-Fi standard that is gaining acceptance where more than 1 Gbps of throughput is needed, such as to support mobile users in the office environment^[21]. This new wireless standard overruns 1-Gbps Ethernet link support but may not require the next step up, which is 10 Gbps. Assuming that 2.5 and 5 Gbps can be made to work over the same cable that supports 1 Gbps, then this standard would provide a much-needed uplink speed improvement for access points supporting 802.11ac radios with their high-bandwidth capabilities.

Conclusion

Ethernet is widely deployed and is the preferred technology for wired local-area networking. Ethernet dominates enterprise LANs, broadband access, and data center networking, and it has also become popular for communication across MANs and even WANs. Further, it is now the preferred carrier wire-line vehicle for bridging wireless technologies, such as Wi-Fi and WiMAX, into local Ethernet networks. Further, the Ethernet marketplace is now large enough to accelerate the development of speeds for specific use cases, such as 25/50 Gbps for data center ToR designs and 2.5/5 Gbps for wireless infrastructure backhaul. The availability of a wide variety of standardized Ethernet data rates allows the network manager to customize a solution to optimize performance, cost, and energy consumption goals^[22].

This popularity of Ethernet technology is due to the availability of cost-effective, reliable, and interoperable networking products from a variety of vendors. The development of converged and unified communications, the evolution of massive server farms, and the continuing expansion of VoIP, TVoIP, and Web 2.0 applications have accelerated the need for ever-faster Ethernet switches.

The success of Gigabit Ethernet and 10-Gbps Ethernet highlights the importance of network-management concerns in choosing a network technology. The 40- and 100-Gbps Ethernet specifications offer compatibility with existing installed LANs, network-management software, and applications. This compatibility has accounted for the survival of 30-year-old technology in today's fast-evolving network environment.

References

- [1] Metcalfe, R., and Boggs, D., “Ethernet: Distributed Packet Switching for Local Computer Networks,” *Communications of the ACM*, July 1976.
- [2] Shoch, J., Dalal, Y., Redell, D., and Crane, R., “Evolution of the Ethernet Local Computer Network,” *Computer*, August 1982.
- [3] Stallings, W., “Gigabit Ethernet,” *The Internet Protocol Journal*, Volume 2, No. 3, September 1999.
- [4] Frazier, H., and Johnson, H. “Gigabit Ethernet: From 100 to 1,000 Mbps,” *IEEE Internet Computing*, January/February 1999.
- [5] GadelRab, S., “10-Gigabit Ethernet Connectivity for Computer Servers,” *IEEE Micro*, May-June 2007.
- [6] McGillicuddy, S., “40 Gigabit Ethernet: The Migration Begins,” *Network Evolution E-Zine*, December 2012.
- [7] Chanda, G., and Yang, Y., “40 GbE: What, Why & Its Market Potential,” Ethernet Alliance White Paper, November 2010.
- [8] Nowell, M., Vusirikala, V., and Hays, R., “Overview of Requirements and Applications for 40 Gigabit and 100 Gigabit Ethernet,” Ethernet Alliance White Paper, August 2007.
- [9] D’Ambrosia, J., Law, D., and Nowell, M., “40 Gigabit Ethernet and 100 Gigabit Ethernet Technology Overview,” Ethernet Alliance White Paper, November 2008.
- [10] Toyoda, H., Ono, G., and Nishimura, S., “100 GbE PHY and MAC Layer Implementation,” *IEEE Communications Magazine*, March 2010.
- [11] Rabinovich, R., and Lucent, A., “40 Gb/s and 100 Gb/s Ethernet Short Reach Optical and Copper Host Board Channel Design,” *IEEE Communications Magazine*, April 2012.
- [12] Morgan, T., “IEEE Gets Behind 25G Ethernet Effort,” *Enterprise Tech*, July 27, 2014.
- [13] Merritt, R., “50G Ethernet Debate Brewing,” *EE Times*, September 3, 2014.
- [14] Nolle, T., “Will We Ever Need 400 Gigabit Ethernet Enterprise Networks?” *Network Evolution E-Zine*, December 2012.
- [15] D’Ambrosia, J., Mooney, P., and Nowell, M., “400 Gb/s Ethernet: Why Now?” Ethernet Alliance White Paper, April 2013.

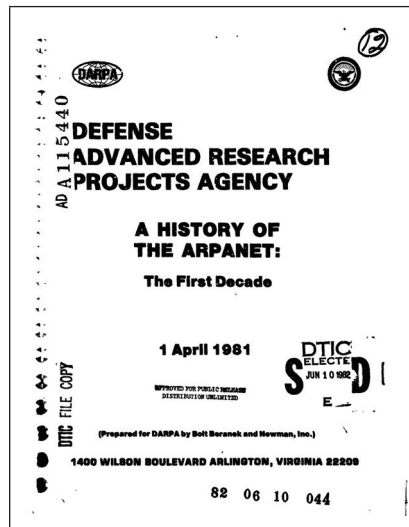
- [16] Hardy, S., “400 Gigabit Ethernet Task Force Ready to Get to Work,” *Lightwave*, March 28, 2014.
- [17] D’Ambrosia, J., “400GbE and High Performance Computing,” *Scientific Computing Blog*, April 18, 2014.
- [18] Duffy, J., “Bridge to Terabit Ethernet,” *Network World*, April 20, 2009.
- [19] D’Ambrosia, J., “TEF 2014: The Rate Debate,” *Ethernet Alliance Blog*, June 23, 2014.
- [20] Kipp, S., “5 New Speeds – 2.5, 5, 25, 50 and 400 GbE,” *Ethernet Alliance Blog*, August 8, 2014.
- [21] Stallings, W., “Gigabit Wi-Fi,” *The Internet Protocol Journal*, Volume 17, No. 1, September 2014.
- [22] Chalupsky, D., and Healey, A., “Datacenter Ethernet: Know Your Options,” *Network Computing*, March 28, 2014.

WILLIAM STALLINGS is an independent consultant and author of numerous books on security, computer networking, and computer architecture. His latest book is *Data and Computer Communications* (Pearson, 2014). He maintains a computer science resource site for computer science students and professionals at ComputerScienceStudent.com. He has a Ph.D. in computer science from M.I.T. He can be reached at ws@shore.net

ARPANET History

Bolt Beranek and Newman Report 4799 entitled “A History of the ARPANET: The First Decade,” is a fascinating document for anyone interested in early Internet History. First published in 1981, a scanned PDF version can be downloaded from the following link:

www.darpa.mil/WorkArea/DownloadAsset.aspx?id=2677



IANA Transition

Since the United States *National Telecommunications and Information Administration* (NTIA) announced its intent to “...transition Key Internet Domain Name Functions to the global multistakeholder community” last March, a flurry of activity has been taking place within the Internet technical and policy communities. The following website provides further information and links to the relevant working groups and documents: <https://www.icann.org/stewardship>

Upcoming Events

The *North American Network Operators’ Group* (NANOG) will meet in San Francisco, California, June 1–3, 2015 and in Montreal, Quebec, Canada, October 5–7, 2015. See: <http://nanog.org>

The *Internet Corporation for Assigned Names and Numbers* (ICANN) will meet in Buenos Aires, Argentina, June 21–25, 2015, and in Dublin, Ireland, October 18–22, 2015. See: <http://icann.org/>

The *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) will meet in Auckland, New Zealand, February 16–26, 2016. See: <http://www.apricot.net>

The *Internet Engineering Task Force* (IETF) will meet in Prague, Czech Republic, July 19–24, 2015, and in Yokohama, Japan, November 1–6, 2015. See: <http://www.ietf.org/meeting/>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

The Internet Protocol Journal (IPJ) is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. Publication of IPJ is made possible by:

<i>Supporters</i>	<i>Diamond Sponsors</i>
 	 
<i>Ruby Sponsor</i>	<i>Sapphire Sponsors</i>
	 
<i>Emerald Sponsors</i>	
   	
   	
   	
<i>Corporate Subscriptions</i>	
    	

Individual Sponsors

Lyman Chapin, Steve Corbató, Dave Crocker, Jay Etchings, Hagen Hultzs, Dennis Jennings, Jim Johnston, Merike Kaeo, Bobby Krupczak, Richard Lamb, Tracy LaQuey Parker, Bill Manning, Andrea Montefusco, Mike O'Connor, Tim Pozar, George Sadowsky, Helge Skrivervik, Rob Thomas, Tom Vest, Rick Wesson.

For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

ADDRESS SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Fred Baker, Cisco Fellow
Cisco Systems, Inc.

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Steve Crocker, Chairman
Internet Corporation for Assigned Names and Numbers

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

June 2015

Volume 18, Number 2

A Quarterly Technical Publication for
Internet and Intranet Professionals

In This Issue

From the Editor	1
Multipath TCP	2
TCP Protocol Wars	15
Fragments	23
Call for Papers	26
Supporters and Sponsors	27

FROM THE EDITOR

The *Transmission Control Protocol* (TCP) is one of the core protocols used in today's Internet. This issue of IPJ is almost entirely devoted to discussions about TCP. Anyone who has studied TCP/IP will have marveled at the "ASCII Art" state diagram for TCP on page 23 of RFC 793, published in 1981. This diagram is a good illustration of both the power and the limitations of using only text characters to draw a "picture." I am happy to report that efforts to define a new format for the RFC series of documents are nearing completion. We will report further on this new RFC format in a future issue.

Your mobile device contains several interfaces, such as USB, WiFi, Cellular Data, and Bluetooth. Most, if not all, of these interfaces can be used for Internet communications, specifically to carry TCP/IP datagrams. In our first article, Geoff Huston looks at an emerging standard, *Multipath TCP* (MPTCP), which allows TCP to operate several simultaneous connections using *different* interfaces.

Although TCP has not fundamentally changed since its introduction in 1981, much work has gone into improving TCP performance in the presence of network congestion and variations in network throughput. Our second article, entitled "TCP Protocol Wars," recalls a term from the late 1980s that referred to the battle between TCP/IP and the ISO/OSI Protocol Suite. This time, the term is used more humorously to compare the many special implementations and refinements to TCP.

If you received a printed copy of this journal in the mail, you should also have received a subscription activation e-mail with information about how to update and renew your subscription. If you didn't receive such a message, it may be because we do not have your correct e-mail address on file. To update and renew your subscription, just send a message to ipj@protocoljournal.org and include your subscription ID. Your subscription ID and expiration date are printed on the back of your journal.

Let me once again remind you that IPJ relies on the support of numerous individuals and organizations. If you or your company would like to sponsor IPJ, please contact us for further details. Our website at protocoljournal.org contains all back issues, subscription information, a list of current sponsors, and much more.

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

IP Multi-Addressing and Multipath TCP

by Geoff Huston, APNIC

The *Transmission Control Protocol* (TCP) is a core protocol of the Internet networking protocol suite. This protocol transforms the underlying unreliable datagram delivery service provided by the *Internet Protocol* (IP) into a reliable data stream protocol. For me this protocol was the single greatest transformative moment in the evolution of computer networks.

Prior to TCP, computer network protocols assumed that computers wanted a lossless reliable service from the network, and worked hard to provide it. The *Digital Data Communications Message Protocol* (DDCMP) in Digital Equipment Corporation's *DECnet* was a lossless data link control protocol. X.25 in the telecommunications world provided reliable stream services to the attached computers. Indeed, I recall that Ethernet was criticized when it was introduced to the world because of its lack of a reliable acknowledgement mechanism. TCP changed all of that. TCP pushed all of the critical functionality supporting reliable data transmission right out of the network and into the shared state of the computers at each end of the TCP conversation. TCP embodies the *End-to-End Principle* of the Internet architecture, where there is no benefit in replicating within the network functionality that can be provided by the end points of a conversation. What TCP required of the network was a far simpler service where packets were allowed to be delivered out of order, but packets could be dropped and TCP would detect and repair the problem and deliver to the far-end application precisely the same bit stream that was passed into the TCP socket in the first place.

The TCP protocol is now some 40 years old, but that doesn't mean that it has been frozen over all these years.

TCP is not only a reliable data stream protocol, but also a protocol that uses *Adaptive Rate Control*. TCP can operate in a mode that allows the protocol to push as much data through the network as it can. A common mode of operation is for an individual TCP session to constantly probe into the network to see what the highest sustainable data rate is, interpreting packet loss as the signal to drop the sending rate and resume the probing. This aspect of TCP has been a constant field of study, and much work has been done in the area of flow control. We now have many variants of TCP that attempt to optimize the flow rates across various forms of networks.

Other work has looked at the TCP data acknowledgement process, attempting to improve the efficiency of the algorithm under a broad diversity of conditions. *Selective Acknowledgments* (SACK) allowed a receiver to send back more information to the sender in response to missing data. *Forward Acknowledgment* (FACK) addresses data-loss issues during TCP *Slow Start*.

One approach to trying to improve the relative outcome of a data transfer, as compared to other simultaneously open TCP sessions, is to split the data into multiple parts and send each part in its own TCP session. This splitting effectively opens up numerous parallel TCP sessions. A variant of TCP, *MulTCP*, emulates the behavior of multiple parallel TCP sessions in a single TCP session. These behaviors assume the same endpoints for the parallel TCP sessions and assume the same end-to-end path through the network. An evolution of TCP that uses multiple parallel sessions but tries to spread these sessions across multiple *paths* through the network is *Multipath TCP* (MPTCP).

Multipath TCP had a brief moment of prominence when it was revealed that Apple's release of iOS 7 contained an implementation of Multipath TCP for the company's *Siri* application, but it has the potential to play a bigger role in the mobile Internet. In this article, I will explore this TCP option in a little more detail, and see how it works and how it may prove to be useful in today's mobile networks.

Multi-Addressing in IP

First we need to return to one of the basic concepts of networking, that of addressing and addresses. Addresses in the Internet Protocol were subtly different from many other computer communications protocols that were commonly used in the 1970s and 1980s. While many other protocols used the communications protocol-level address as the address of the host computer, the Internet Protocol was careful to associate an IP address with the *interface* to a network. This distinction was a relatively unimportant one in most cases because computers usually had only a single network attachment interface. But it was a critical distinction when the computer had two or more interfaces to two or more networks. An IP host with two network interfaces has two IP addresses—one for each interface. In IP it is the interface between the device and the network that is the addressed endpoint in a communication. An IP host accepts an IP packet as being addressed to itself if the IP address in the packet matches the IP addresses of the network interface that received the packet, and when sending a packet, the source address in the outgoing packet is the IP address of the network interface that was used to pass the packet from the host into the network.

As simple as this model of network addressing may be, it does present some operational problems. One implication of this form of addressing is that when a host has multiple interfaces, the application-level conversations using TCP are “sticky.” If, for example, a TCP session was opened on one network interface, the network stack in the host cannot quietly migrate this active session to another interface while maintaining the common session state. An attempt by one “end” of a TCP conversation to change the IP address for an active session would not normally be recognized at the other end of the conversation as being part of the original session. So having multiple interfaces and multiple addresses does not create additional resiliency of TCP connections.

The simplicity of giving each network interface a unique IP address does not suit every possible use case, and it was not all that long before the concept of *secondary addresses* came into use. The use of secondary addresses was a way of using multiple addresses to refer to a host by allowing a network interface to be configured with multiple IP addresses. In this scenario, an interface receives packets addressed to any of the IP addresses associated with the interface. Outgoing packet handling allows the transport layer to specify the source IP address, and this action overrides the default action of using the primary address of the interface on outgoing packets. Secondary addresses have their uses, particularly when you are trying to achieve the appearance of multiple application-level “personas” on a single common platform, but in IPv4 they were perhaps more of a specialized solution to a particular family of requirements, rather than a commonly used approach. Applications using TCP were still “sticky” with IP addresses that were used in the initial TCP handshake and could not switch the session between secondary IP addresses on the same interface.

IPv6 addressing is somewhat different. The protocol allows from the outset for an individual interface to be assigned multiple IPv6 unicast addresses without the notion of “primary” and “secondary” addresses. The IPv6 protocol introduces the concept of an address *scope*, so an address may be assuredly unique in the context of the local link-layer network, or it may have a global scope, for example. Privacy considerations have also introduced the concept of permanent and temporary addresses, and the efforts to support a certain form of mobility have introduced the concepts of *home addresses* and *care-of addresses*.

However, to some extent these IPv6 changes are cosmetic modifications to the original IPv4 address model. If an IPv6 host has multiple interfaces, each of these interfaces has its own set of IPv6 addresses, and when a TCP session is started using one address pair TCP does not admit the ability to shift to a different address pair in the life of the TCP session. A TCP conversation that started over one network interface is stuck with that network interface for the life of the conversation, whether it’s IPv4 or IPv6.

The Internet has changed significantly with the introduction of the mobile Internet, and the topic of multi-addresses is central to many of the problems with mobility. Mobile devices are adorned with many IP addresses. The cellular radio interface has its collection of IP addresses. Most of these “smart” devices also have a WiFi interface that also has its set of IPv4 and possibly IPv6 addresses. And there may be a Bluetooth network interface with IP addresses, and perhaps some USB network interface as well. When active, each of these network interfaces requires its own local IP address. We now are in an Internet where devices with multiple active interfaces and multiple usable IP addresses are relatively commonplace. But how can we use these multiple addresses?

For many scenarios there is little value in being able to use multiple addresses. The conventional behavior is where each new session is directed to a particular interface, and the session is given an outbound address as determined by local policies. However, when we start to consider applications in which the binding of location and identity is more fluid, network connections are transient, and the cost and capacity of connections differ (as is often the case in today's mobile cellular radio services and in WiFi roaming services), then having a session that has a certain amount of agility to switch across networks can be a significant factor.

If individual end-to-end sessions could use multiple addresses, and by inference could use multiple interfaces, then an application could perform a seamless handoff between cellular data and WiFi, or even use both at the same time. Given that the TCP interface to IPv4 and IPv6 is identical, it is even quite feasible to contemplate a seamless handoff between the two IP protocols. The decision as to which carriage service to use at any time would no longer be a decision of the mobile carrier or that of the WiFi carrier, or that of the device, or that of its host operating system. If applications could use multiple addresses, multiple protocols, and multiple interfaces, then the decision could be left to the application itself to determine how best to meet its needs as connections options become available or as they shut down. At the same time as the debate between traditional mobile operators in the licensed spectrum space and the WiFi operators in the unlicensed spectrum space heats up over access to the unlicensed spectrum, the very nature of how devices and applications implement "WiFi handoff" is changing. Who is in control of this handoff function is changing as a result. Multi-Addressing and Multipath TCP is an interesting response to this situation; it allows individual applications to determine how they want to operate in a multi-connected environment.

SHIM6

One of the first attempts to use multiple addresses in IP was the *Site Multihoming by IPv6 Intermediation* (SHIM6) effort in IPv6.

In this case the motivation was end-site resilience in an environment of multiple external connections, and the constraint was to avoid the use of an independently routed IPv6 address prefix for the site. So SHIM6 was an effort to support site multi-homing without routing fragmentation. To understand the SHIM6 model, we need to start with an end site that does not have its own provider-independent IPv6 address prefix, yet is connected to two or more upstream transit providers that each provide addresses to the end site. In IPv4 it's common to see this scenario approached with *Network Address Translators* (NATs). In IPv4 the site is internally addressed using a private address prefix, and the interface to each upstream provider is provisioned with a NAT. Outbound packets have their source address rewritten to use an address that is part of the provider's prefix as it transits the NAT.

Which provider is used is a case of internal routing policies toward each of the NATs. Although it is possible to configure a similar setup in IPv6 using an IPv6 *Unique Local Address* (ULA) prefix as the internal address and NAT IPv6-to-IPv6 devices connected to each upstream service provider, one of the concepts behind IPv6 and its massive increase in address space was the elimination of NATs. So how can an IPv6 end site be homed into multiple upstream service providers without needing to advertise a more specific routing entry in the interdomain routing tables and avoiding the use of any form of network address translation?

The conventional IPv6 architecture has the site receiving an end-site prefix delegation from each of its upstream service providers, and the interface routers each advertising its end-site prefix into the site. Hosts within the site see both router advertisements, and they configure their interface with multiple IPv6 addresses, one for each site prefix. Presumably, the end site chooses to multi-home in order to benefit from the additional resiliency that such a configuration should offer. When the link to one provider is down, there is a good chance that the other link will remain up, particularly if the site has been careful to engineer the multi-homed configuration using discrete components at every level. It would be even better if even when the link to the upstream provider is up and that provider can't reach a specific destination, another of the site's upstream providers could continue to support all active end-to-end conversations without interruption, in exactly the same manner as when this functionality is implemented in the routing system.

What SHIM6 attempted was a host-based approach to use the additional local IPv6 addresses in the host as indicators of potential backup paths to a destination. If a communication with a remote counterpart were to fail (that is, the flow of incoming packets from the remote host stopped), then the IP-level shim in the local host would switch to use a different source/destination address pair. To prevent the upper-level transport protocol from being fatally confused by these address changes in the middle of one or more active sessions, the local SHIM module also included a network address translation function. This function helped ensure that although the address pair on the wire may have changed, the address pair presented to the upper layer by the shim would remain constant, and the path change would not be directly visible at the transport layer of the protocol stack.

This approach essentially folds the NAT function into the host IP protocol stack. In terms of design it avoided altering either TCP or *User Datagram Protocol* (UDP), and endeavoured to preserve the IP addresses used by active transport sessions. What this approach implied was that if you wanted to change the routing path but not change the IP addresses used by transport, then address translation was an inevitable consequence.

Network-based NATs was the response in IPv4, and to avoid this problem in IPv6 the SHIM6 effort attempted to push the NAT functionality further “back,” implementing a NAT in each host.

SHIM6 was an approach that was less than entirely satisfactory.

Network operators expressed deep distrust of pushing decision-making functionality back into individual hosts (a distrust that network operators continue to hold when the same issue arises with WiFi handoff). The network operators wanted to control the connectivity structure for the hosts in their network, in precisely the same manner as the routing system provided network-level control over traffic flows. So although these network operators had some sympathy with the SHIM6 objective of avoiding further bloat in the routing table, which reduced the “independence” of attached end sites by using IPv6 address prefixes drawn from the upstream address block, they were unsupportive of an approach that pushed connectivity choice and control back to individual end host systems.

Outside of this issue of control over the end host was another multi-homing problem that SHIM6 did not address. Although the provision of backup paths in the case of failure of the primary path is useful, what is even more useful is the ability to use the backup paths in some form of load-sharing configuration. However, at this point the SHIM6 approach runs into problems. Because SHIM6 operates at the IP layer, it is not directly aware of packet sequencing. When a SHIM unit at one end of a conversation splays a sequence of packets across multiple paths, the corresponding SHIM unit at the remote end passes the packets into the upper transport layer in the order of their arrival, not in the original order. This out-of-order delivery can be a significant problem for TCP if SHIM6 leaves multiple paths open. The best SHIM6 can provide is a primary/backup model for individual sessions, where at any time all data traffic for a session is passed along the primary path.

Inexorably, we are drawn to the conclusion that the most effective place to insert functionality that allows a data flow to use multiple potential paths across the network is in the transport layer itself, and we need to jack ourselves further up the protocol stack from the IP level approach of SHIM6 and re-examine the space from the perspective of TCP.

Multipath TCP

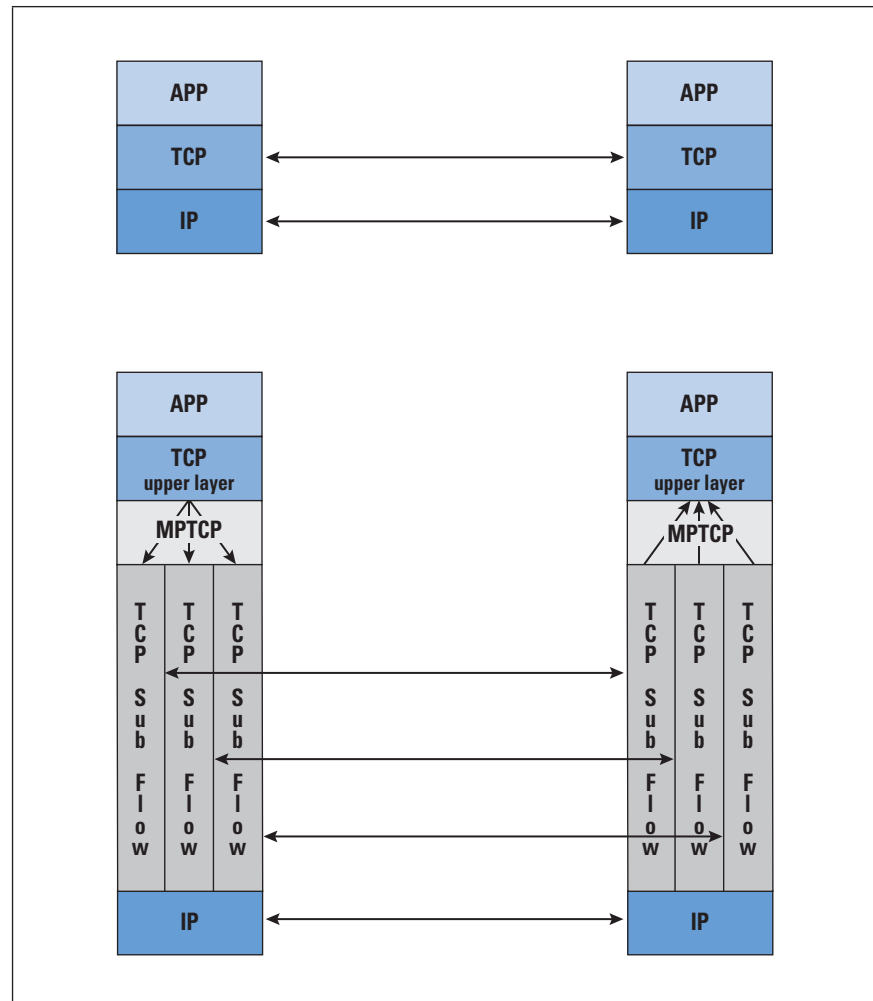
The approach of incorporating multiple IP addresses in the transport protocol is comparable to the efforts of SHIM6 one level further down in the protocol stack, in so far as this approach is an end-to-end mechanism with a shared multiplex state maintained in the two end hosts, and no state whatsoever in the network.

The basic mechanisms for MPTCP are also similar to that of SHIM6, with an initial capability exchange to confirm that both parties support the mechanism, allowing the parties to then open up additional paths, or channels. But at this point the functionality diverges. In SHIM6 these alternate paths are provisioned as backup paths if the primary path fails, whereas in the case of MPTCP these additional paths can be used immediately to spread the load of the communication across these paths, if the application so desires.

One of the most critical assumptions of MPTCP was drawn from SHIM6, in that the existence of multiple addresses in a host is sufficient to indicate the existence of multiple diverse paths within the network. Whether or not this assumption is, in fact, the case is perhaps not that critical, in that even in the case where the addresses are on the same path from end to end, the result is roughly equivalent to running multiple parallel sessions of TCP.

The basic approach to MPTCP is the division of the single outbound flow of the application into multiple *subflows*, each operating its own end-to-end TCP session, and the rejoining of multiple input subflows into a single flow to present to the remote counterpart application. This approach is shown in Figure 1.

Figure 1: Comparison of Standard TCP and MPTCP Protocol Stacks



This solution is essentially a “shim” inserted in the TCP module. To the upper-level application, MPTCP can operate in a manner that is entirely consistent with TCP, so that the opening up of subflows and the manner in which data is assigned to particular subflows is intentionally opaque to the upper-level application. The envisaged *Application Programming Interface* (API) allows the application to add and remove addresses from the local multipath pool, but the remainder of the operation of the MPTCP shim is not envisaged to be managed directly by the application. MPTCP also leaves the lower-level components of TCP essentially untouched, in so far as each MPTCP subflow is a conventional TCP flow. On the data sender’s side, the MPTCP shim essentially splits the received stream from the application into blocks and directs individual blocks into separate TCP subflows. On the receiver’s side, the MPTCP shim assembles the blocks from each TCP subflow and reassembles the original data stream to pass to the local application.

Operation of MPTCP

TCP has the ability to include 40 bytes of TCP *options* in the TCP header, indicated by the *Data Offset* value. If the Data Offset value is greater than 5, then the space between the final 32-bit word of the TCP header (*Checksum* and *Urgent Pointer*) and the first octet of the data can be used for options. MPTCP uses the *Option Kind* value of 30 to denote MPTCP options. All MPTCP signalling is contained in this TCP header options field.

The MPTCP operation starts when the initiating host passes a MP_CAPABLE capability message in the MPTCP options field to the remote host as part of the initial TCP SYN message when opening the TCP session. The SYN+ACK response contains a MP_CAPABLE flag in its MPTCP options field of the SYN+ACK response if the other end is also MPTCP-capable. The combined TCP and MPTCP handshake concludes with the ACK and MP_CAPABLE flag, confirming that both ends now have each other’s MPTCP session data. This capability negotiation exchanges 64-bit keys for the session, and each party generates a 32-bit hash of the session keys, which are subsequently used as a shared secret between the two hosts for this particular session to identify subsequent subjoin connection attempts.

Further TCP subflows can be added to the MPTCP session by a conventional TCP SYN exchange with the MPTCP option included. In this case the exchange contains the MP_JOIN values in the MPTCP options field. The values in the MP_JOIN exchange include the hash of the original receiver’s session key and the token value from the initial session, so that both ends can associate the new TCP session with the existing session, as well as a random value intended to prevent replay attacks.

The MP_JOIN option also includes the sender's address index value to allow both ends of the conversation to reference a particular address even when NATs on the path perform address transforms. MPTCP allows these MP_JOINS to be established on any port number, and by either end of the connection. Therefore, although a MPTCP web session may start using a port 80 service on the server, subsequent subflows may be established on any port pair, and it is not necessary for the server to have a LISTEN open on the new port. The MPTCP session token allows the 5-tuple of the new subflow (protocol number, source and destination addresses, and source and destination port numbers) to be associated with the originally established MPTCP flow. Two hosts can also inform each other of new local addresses without opening a new session by sending ADD_ADDR messages, and remove them with the complementary REMOVE_ADDR message.

Individual subflows use conventional TCP signalling. However, MPTCP adds a *Data Sequence Signal* (DSS) to the connection that describes the overall state of the data flow across the aggregate of all of the TCP subflows that are part of this MPTCP session. The sender sequence numbers include the overall data sequence number and the subflow sequence number that is used for the mapping of this data segment into a particular subflow. The DSS Data ACK sequence number is the aggregate acknowledgement of the highest in-order data that the receiver receives. MPTCP does not use SACK, because this acknowledgement is left to the individual subflows.

To prevent data loss that causes blockage on an individual subflow, a sender can retransmit data on additional subflows. Each subflow uses a conventional TCP sequencing algorithm, so an unreliable connection will cause that subflow to stall. In this case MPTCP can use a different subflow to resend the data, and if the stalled condition is persistent it can reset the stalled subflow with a TCP RST within the context of the subflow.

Individual subflows are stopped by a conventional TCP exchange of FIN messages, or through the TCP RST message. The shutting down of the MPTCP session is indicated by a data FIN message that is part of the data sequencing signalling within the MPTCP option space.

Congestion control appears still to be an open issue for MPTCP. An experimental approach is to couple the congestion windows of each of the subflows, increasing the sum of the total window sizes at a linear rate per *Round-Trip Time* (RTT) interval, and applying the greatest increase to the subflows with the largest existing window. In this way the aggregate flow is no worse than a single TCP session on the best available path, and the individual subflows take up a fair share of each of the paths they use. Other approaches are being considered that may reduce the level of coupling of the individual subflows.

MPTCP and Middleware

Today's Internet is not the Internet of old. It is replete with various forms of middleware that include NATs, load balancers, traffic shapers, proxies, filters, and firewalls. The implication of this reality is that any deviation from the most basic forms of the use of IP will run into various issues with various forms of middleware.

For MPTCP, the most obvious problem is that of middleware that strips out unknown TCP options.

However, more insidious issues come with the ADD_ADDR messages and NATs on the path. Sending IP addresses within the data payload of a NATed connection is always a failure-prone option, and MPTCP is no exception here. MPTCP contains no inbuilt NAT detection functions, and there is no way to determine the direction of the NAT. A host can communicate to the remote end its own IP address or additional available addresses, but if there is a NAT translating the local-host outbound connections, then the actual address will be unavailable for use until the host actually starts a TCP session using this local address as the source.

A simple approach that is effective where NATs are in place is to leave the role of initiation of new subflows to the host that started the connection in the first place. In a client-server environment this solution would imply that the role of setting up new subflows is best left to the client in such cases. However, no such constraints exist when there are no NATs, and in that case either end can initiate new subflows, and the ADD_ADDR messages can keep the other end informed about potential new parallel paths between the two hosts. Logically it makes little sense for MPTCP itself to define a NAT-sensing probe behavior, but it makes a lot of sense for the application using MPTCP to undertake such a test.

The Implications of MPTCP

MPTCP admits considerable flexibility in the way an application can operate when many connection options are available.

All TCP subflows carry the MPTCP option, so that the MPTCP shared state is shared across all active TCP subflows. No single subflow is the "master" in the MCTCP sense. Subflows can be created when interfaces come up, and removed when they go down. Subflows are also IP protocol agnostic: they can use a collection of IPv4 and IPv6 connections simultaneously. Subflows can be used to load-share across multiple network paths, or operate in a primary/backup configuration depending on the application and the flexibility offered in the API in particular implementations of MPTCP.

When applied to mobile devices, this behavior can lead to unexpected results. I always assumed that my device was incapable of “active handoff.” Any connections that were initiated across the cellular radio interface had to stay on that interface, and any connections established over the WiFi interface would also stay on that WiFi network. I always understood that active sessions could not be handed off to a different network. Although it was never an explicitly documented feature (or if it was I have never seen it), I had also assumed that when my mobile device was in an area with an active WiFi connection, then the WiFi would take precedence over its *fourth-generation* (4G) connection for all new connections. This assumption matched the factor of typical data tariffs, where the marginal cost of data over 4G is typically somewhere between 10 and 1,000 times higher than the marginal cost of the same data volume over the WiFi connection. But if applications use MPTCP instead of TCP, then how will they balance their network use across the various networks? The way MPTCP is defined it appears that applications simply open subflows on all available local interfaces, and then the fastest network, rather than the cheapest, will take on the greatest volume of traffic.

But, as usual, it can always get more complicated. What if the WiFi network is a corporate service, with NATs, split-horizon *Virtual Private Networks* (VPNs) and various secure servers? If my device starts to perform MPTCP in such contexts, then to what extent are the properties of my WiFi connection preserved in the cellular data connection? Have I exposed new vulnerabilities by doing this? How can a virtual interface, such as a VPN, inform an MPTCP-aware application that other interfaces are not in the same security domain as the VPN interface?

However, it does appear that MPTCP has a role to play in the area of seamless WiFi handoff. With MPTCP is it possible for a mobile handset to enter a WiFi-serviced area and include a WiFi subflow into the existing data transfer without stopping and restarting the data flow? The application may even shut down the cellular radio subflow when the WiFi subflow is active. This functionality is under the control of the application using MPTCP, rather than being under the control of the host operating system of the carrier.

Going Up the Stack

Of course it does not stop at the transport layer and with the use of MPTCP. Customized applications can perform handoffs themselves.

For example, the “mosh” application is an example of a serial form of address agility, where the session state is a shared secret, and the server will accept a reconnection from any client’s IP address, as long as the client can demonstrate its knowledge of the shared secret.

Extending the TCP data-transfer model to enlist multiple active TCP sessions at the application level in a load-balancing configuration is also possible, in a manner not all that different from MPTCP.

Of course one could take this further. Rather than use multiple TCP sessions between the same two endpoints, you could instead share the data from the same server across multiple endpoints, and use multiple TCP sessions to these multiple servers. At this point you have something that looks remarkably like the peer-to-peer data-distribution architecture.

Another approach is to format the data stream into “messages” and permit multiple messages to be sent across diverse paths between the two communicating systems. This approach, the *Stream Control Transmission Protocol* (SCTP), is similar to MPTCP in that it can take advantage of multiple addresses to support multiple paths. It combines the message transaction qualities of UDP with the reliable in-sequenced transport services of TCP. The problem of course in today’s network is that because it is neither TCP nor UDP, many forms of middleware, including NATs, are often hostile to SCTP and they drop SCTP packets. One additional cost of the escalation of middleware in today’s Internet. These days innovation in protocol models is limited by the rather narrow rules applied by network middleware, and the approximate general rule in today’s Internet is that it’s TCP, UDP, or middleware fodder!

It has been observed numerous times that the abstraction of a network protocol stack is somewhat arbitrary, and it’s possible to address exactly the same set of requirements at many different levels in the reference stack. In the work on multipath support in the Internet, we’ve seen approaches that exploit parallel data streams at the data link layer, at the IP layer, within routing, in the transport layer, and in the application layer. Each has its respective strengths and weaknesses. But what worries me is what happens if you inadvertently encounter a situation where you have all of these approaches active at the same time? Is the outcome one of amazing efficiency, or paralyzing complexity?

For Further Reading

- [1] Erik Nordmark and Marcelo Bagnulo, “Shim6: Level 3 Multihoming Shim Protocol for IPv6,” RFC 5533, June 2009.
- [2] Janardhan Iyengar, Costin Raiciu, Sebastien Barre, Mark Handley, and Alan Ford, “Architectural Guidelines for Multipath TCP Development,” RFC 6182, March 2011.
- [3] Mark Handley, Alan Ford, Costin Raiciu, and Olivier Bonaventure, “TCP Extensions for Multipath Operation with Multiple Addresses,” RFC 6824, January 2013.

- [4] Bit Torrent:
<https://wiki.theory.org/BitTorrentSpecification>
- [5] Geoff Huston, "Anatomy: A Look inside Network Address Translators," *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [6] M. Mathis, J. Mahdavi, S. Floyd, and A. Romanow, "TCP Selective Acknowledgment Options," RFC 1881, October 1996.
- [7] M. Mathis and J. Mahdavi, "Forward Acknowledgment: Refining TCP Congestion Control," Proceedings of SIGCOMM, August 1996.
- [8] Randall Stewart, "Stream Control Transmission Protocol," RFC 4960, September 2007.
- [9] Ethan Blanton and Mark Allman, "TCP Congestion Control," RFC 5681, September 2009.
- [10] Olivier Bonaventure, "Apple seems to also believe in Multipath TCP," Blog post,
<http://perso.uclouvain.be/olivier.bonaventure/blog/html/2013/09/18/mptcp.html>
- [11] Jonathan B. Postel, "Transmission Control Protocol," RFC 793, September 1981.
- [12] Geoff Huston, "TCP Performance," *The Internet Protocol Journal*, Volume 3, No. 2, June 2000.
- [13] Geoff Huston, "The Future for TCP," *The Internet Protocol Journal*, Volume 3, No. 3, September 2000.
- [14] Wesley M. Eddy, "Defenses Against TCP SYN Flooding Attacks," *The Internet Protocol Journal*, Volume 9, No. 4, December 2006.

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

TCP Protocol Wars

by Geoff Huston, APNIC

There are two end-to-end transport protocols in common use in today's Internet: the *User Datagram Protocol* (UDP) and the *Transmission Control Protocol* (TCP).

UDP is an abstraction of the basic IP datagram, in that UDP is an unreliable medium. Packets sent using UDP may or may not go to their intended destination. UDP packets may be reordered, duplicated, or lost. UDP has no flow control or throttling. The packet quantization in UDP is explicit: if the sender splits data into two UDP packets, then the receiver will collect the data using two distinct read operations.

TCP is a reliable end-to-end flow-controlled stream protocol. A stream of data passed into a TCP socket at one end will be read as a stream of data at the other end. The packet quantization is hidden from the application, as are the mechanics of flow control, loss detection, and retransmission and session establishment and teardown. TCP will not preserve any inherent timing within the data stream, but will preserve the integrity of the stream.

Critically, the Internet assumes that most of the network resources are devoted to passing TCP traffic, and it also assumes that the flow-control algorithms used by these TCP sessions all behave in approximately similar ways. If the switching and transmission resources of the network are seen as a common resource, then the assumption about the uniform behavior of TCP sessions implies that these end-to-end transport sessions will behave similarly under contention. The result is that, to a reasonable level of approximation, a set of concurrent TCP sessions will self-equilibrate to give each TCP session an equal share of the common resource. In other words, the network itself does not have to impose “fairness” on the TCP flows that pass across it—as long as all the flows are controlled by a uniform flow-control algorithm, the flows will interact with each other in a manner that is likely to allocate an equal proportion of the network resources to each active TCP flow. At least that's the theory.

This theory raises numerous questions of whether these assumptions are true in today's Internet and what may be changing with these assumptions.

Other Protocols?

Is it still a choice between UDP and TCP? Despite many technical efforts to specify new end-to-end transport protocols, there is little chance that any new protocol will gain acceptance in today's Internet. The network contains large numbers of intercepting “middleware,” and these units function as security firewalls by using rules that are very limited in the protocols that they admit.

The most common filters in middleware are configured to admit only IP protocols 6 and 17 (TCP and UDP, respectively), and drop all others. This setup has implied that more recent end-to-end transport protocols, such as the *Stream Control Transmission Protocol* (SCTP)^[10] or the *Datagram Congestion Control Protocol* (DCCP)^[11], for example, have very limited applicability in the public Internet, because they can be used only in environments where there is no such intercepting middleware.

TCP or UDP?

In this world where choice is limited to TCP or UDP, the conventional view was that the bulk of the traffic was carried in TCP, whereas UDP was used in limited contexts for *Domain Name System* (DNS) name resolution, running time, and network management. This view raises the question as to whether TCP still carries the bulk of the Internet traffic load.

It is not necessarily true that TCP still carries the bulk of the Internet traffic load, although reliable data sources that provide visibility into the actual traffic profile seen on end user-facing networks is not easily forthcoming. A recent study of traffic profiles between 2002 and 2009 by the *Center for Applied Internet Data Analysis* (CAIDA)^[1] points to a UDP/TCP ratio value of 0.11 when looking at the volume of data being transported by the two protocols. In other words, some 90% of the traffic was carried inside TCP sessions and 10% inside UDP sessions. For UDP this value is considerably higher than would be conventionally expected from the combination of only DNS and *Network Time Protocol* (NTP) payloads in UDP. The study points out: “A port-based analysis suggests that the recent increase in UDP flows on the traces analyzed stems mainly from *Peer-to-Peer* (P2P) applications using UDP for their overlay signalling traffic,” a result that corresponds to reports of the use of the *Low Extra Delay Background Transport* (LEDBAT) protocol for *BitTorrent*^[2]. More recently, video streaming applications have also turned to TCP, using local buffer management in the playback device to overcome TCP-induced signal jitter.

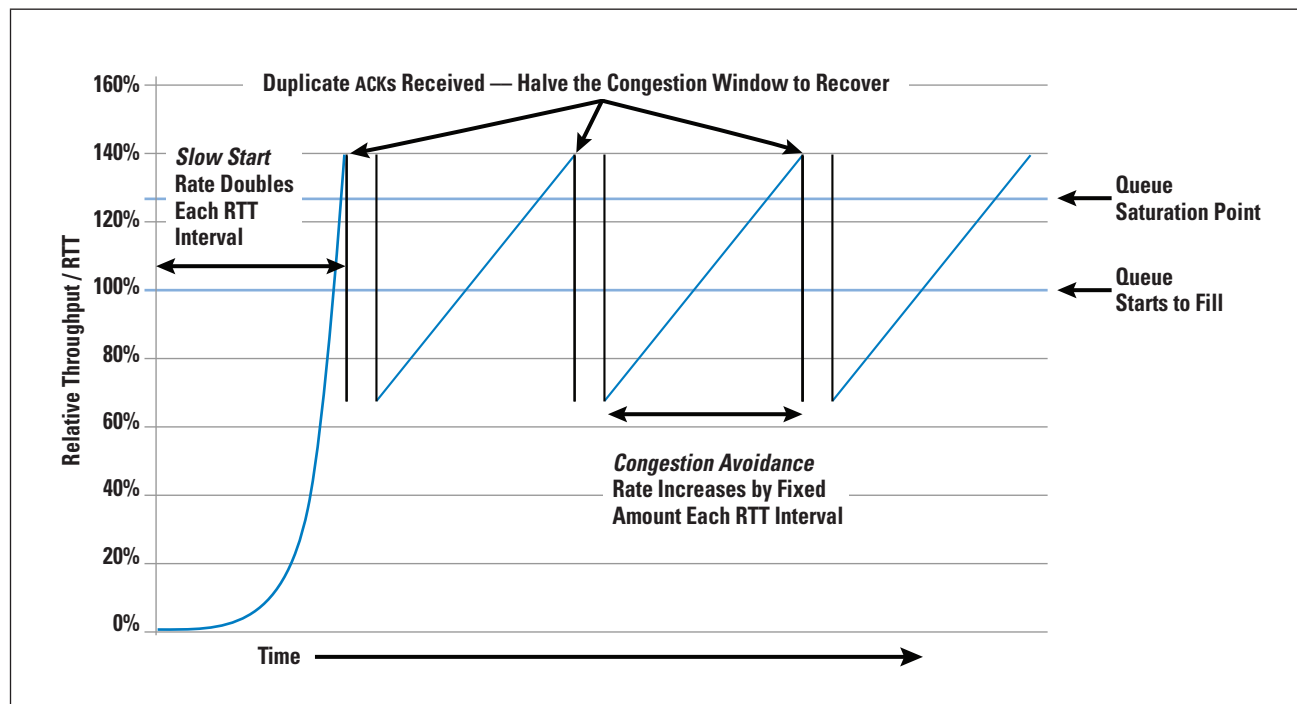
It is reasonable to assume that the overall majority of the Internet traffic load is carried in TCP, and therefore the behavior of the TCP flow-control algorithm is a matter of interest.

TCP Flow Control – TCP Reno

TCP does not have a single flow-control algorithm. Although the common TCP protocol specification defines how to establish and shut down a session, and defines the way in which received data is acknowledged back to the sender, the core protocol specification does not specify how the two ends negotiate the speed at which data is passed between them. This negotiation has been left to the various implementations of the TCP flow-control algorithm.

“Conventional” flow control in TCP is typified by the behavior of the TCP *Reno* algorithm (Figure 1).

Figure 1: Idealised TCP Reno Flow Control



There are two distinct phases of behavior: the *Slow Start* phase, where the sending rate is doubled every *Round-Trip Time* (RTT) interval, and a *Congestion Avoidance* phase, where the sending rate is increased by a fixed amount—one *Message Segment Size* (MSS)—in each RTT interval. When the sender is notified of packet loss—by receiving a duplicate *Acknowledgment* (ACK) message from the receiver—the actions of the sender vary according to its current phase. In *Slow Start* phase a duplicate ACK will shift the sender to *Congestion Avoidance* mode. In *Congestion Avoidance* mode a duplicate ACK will cause the sender to halve its sending rate and continue in this mode. Three duplicate ACKs in succession will cause the session to restart from scratch in *Slow Start* mode, because three duplicate ACKs signals a higher rate of congestion which means that the two ends of the TCP stream have lost their shared flow state assumption.

In steady state the TCP Reno flow-control algorithm increases the flow rate by a constant amount each round-trip time interval, and when a packet is dropped, because of buffer overflow in a switch, the algorithm halves the flow rate. The result is an *Additive Increase Multiplicative Decrease* (AIMD) algorithm, which tends to place high levels of pressure on the buffers in the network while there is still available buffer space, and react dramatically when the buffers eventually overfill and reach the packet drop point. Crudely, this process is a “boom and bust” form of feedback control.

Better than Reno

There have been strong motivations by application families to break out of this form of TCP flow-control behavior. One motivation is to use a more even packet flow across the network, and remove some of the “jerkiness” inherent in TCP Reno. There is also the motivation that a more sensitive flow-control application could achieve a superior outcome compared to TCP Reno. In other words, a different TCP flow-control algorithm could achieve better than its “fair share” when competing against a set of concurrent TCP Reno flows.

The first of these motivations is a simple change. In an attempt to double the pressure on other concurrent TCP sessions, the AIMD algorithm can be adjusted by increasing the sending speed a larger constant amount, and decreasing it by less following packet loss (*MulTCP* uses this model). For example, if the speed was increased by 2 MSS units each RTT interval and the sending rate was reduced by one-quarter rather than one-half upon receipt of a duplicate ACK, then the resultant behavior would, in an approximate sense, behave like two concurrent TCP sessions, and in a fair sharing scenario this form of flow control would attempt to secure double the network resources of an equivalent TCP Reno session.

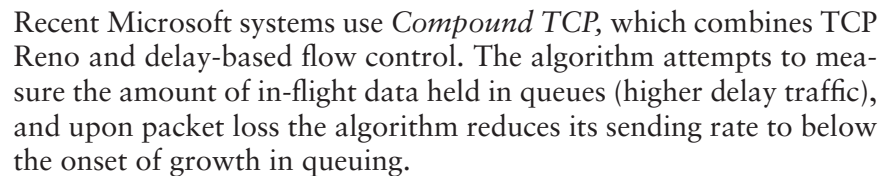
Another variant of this approach is *Highspeed TCP* which increases its frequency of probing into potentially claimable capacity by increasing its sending rate by a larger volume while keeping its reduction rate at a constant value. This protocol probes for the packet-loss onset at a far higher frequency than either TCP Reno or *MulTCP*, and is capable of accelerating to much higher flow speeds in a much shorter time interval.

Binary Increase Congestion Control (BIC) and its variant CUBIC use a nonlinear increase function rather than a constant rate increase function (Figure 2). Instead of increasing the speed by a fixed amount each RTT in Congestion Avoidance mode, BIC remembers the sending rate at the onset of packet drop, and each RTT interval increases its speed by one-half of the difference between the current sending rate and the assumed bottleneck rate.

BIC quickly drives the session towards the bottleneck capacity, and then probes more cautiously when the sending speed is close to the bottleneck capacity. Again, compared to a Reno flow session, CUBIC should produce a superior outcome.

Other flow-control algorithms move away from using packet loss as the control indication and tend to oscillate more frequently around the point of the onset of queuing in the routers in the network path. This form of feedback control is sensitive to the relative time differences between sent packets and received ACKs.

Figure 2: Idealized TCP BIC Flow Control



Crossing the Beams: TCP implemented in UDP

THE INTERNET PROTOCOL JOURNAL

This approach has been used in the widely deployed *BitTorrent* application (LEDBAT), and more recently by Google in its experiments with *Quick UDP Internet Connections* (QUIC)^[8] and *SPDY* (pronounced “speedy”).

Google’s QUIC uses a TCP emulation in UDP that has a data encoding that includes *Forward Error Correcting Codes* (FEC) as a way of performing a limited amount of repair of the data stream in the face of packet loss without retransmission. QUIC performs bandwidth estimation as a means of rapidly reaching an efficient sending rate. SPDY further assists QUIC by multiplexing application sessions within a single end-to-end transport protocol session. This approach avoids the startup overhead of each TCP session, and leverages the observation that TCP takes some time to establish the bottleneck capacity of the network. The use of UDP also avoids intercepting middleware that performs deep packet inspection on TCP flows and modifies their advertised window size to perform external moderation on TCP flow rate.

There is, however, one issue with the use of UDP as a substitute for TCP, and although public reports from Google on this topic have not been published, it is a source of concern. The problem relates to the use of UDP through *Network Address Translators* (NATs)^[12] and the issue of address binding times within the NAT. In TCP a NAT takes its directions from TCP. When the NAT sees an opening TCP handshake packet from the “inside,” it creates a temporary address binding and sends the packet to its intended destination (with the translated source address of course). The reception of the response part of the handshake at the NAT causes the NAT to confirm its binding entry and apply it to subsequent packets in this TCP flow. The NAT holds state until it sees a closing exchange or a reset signal that closes the TCP session, or until an idle timer expires. For TCP the NAT is attempting to hold the binding for as long as the TCP session is active. For NATs, UDP is different. Unlike TCP, there is no flow-status information in UDP. So when the NAT creates a UDP binding, it has to hold it for a certain amount of time. There is no clear technical standard here, so implementations vary. Some NATs use very short timers and release the binding quickly, matching the expectation of the use of UDP as a simple query/response protocol. The use of UDP as an ersatz packet-framing protocol for user-level TCP implementation requires the NAT to hold the UDP address binding for longer intervals, corresponding to the hidden TCP session. Some NATs will do so, while others will destroy the binding even though there are still UDP packets active, thus disturbing the hidden TCP session.

This example illustrates the level of compromise in today’s environment between end-to-end protocols and network middleware. TCP sessions are being modified by active middleware that attempts to govern the TCP flow rate by active modification of window sizes within the TCP session, negating some of the efforts of the TCP session to optimize its flow speed.

TCP in UDP passes control of the TCP flow management to the application, and hides the TCP flow parameters from the network. However, UDP sessions are susceptible to interruption by NAT intervention, because some NATs assume that UDP is used only for micro-sessions, and long-held UDP sessions are some form of anomalous behavior that should be filtered by removing the UDP port binding in the NAT.

The Transport Protocol Ecosystem

The Internet is somewhat unique in so far as there is no intrinsic network-level functionality that can allocate a certain amount of network resources to each active flow being carried across the network. The network is not actively “managed.” Network resources are allocated to traffic flows in a manner similar to fluid-flow equilibrium. Each active flow exerts pressure on all other concurrent flows. The higher the relative imbalance, the more the largest flows are pressured to reduce their flow rate by the smaller flows. The system reaches a meta-equilibrium point when all concurrent flows receive approximately equal amounts of network resource.

The underlying assumption here is that a fair result is achieved if all the concurrent flows are operating in a similar manner. What is happening in the network today is a fragmentation of the TCP flow-control algorithm as operating systems, and even applications, prefer to use a customized flow-control algorithm that attempts to optimize their position by exerting slightly more pressure on other TCP sessions, causing them to drop their flow rates in response. These techniques do not create additional network transmission capacity, they bias the way in which network capacity is available to individual traffic flows in their favor. So if a TCP session is able to secure better than its “fair share” of a laden network, then other sessions are necessarily affected and receive less than their “fair share.”

There is some relationship between these protocol-level efforts and the *Net Neutrality* policy debates. The proponents of a Net Neutrality position argue that the network should be a largely passive entity, and that the interaction of the various traffic flows produces a fair and efficient outcome. The network resources will be fully allocated to carrying traffic with relatively small levels of retransmission (efficiency), and the concurrent flows will interact with each other to produce an outcome where each flow gathers approximately equal network resource (fair). With the increasing level of diversity in approaches to packet-flow management, and the options of whether to use the flow-control services provided by the operating system platform or go the path of using UDP as the transport protocol and passing the flow-control algorithm to the application, what is being witnessed is some amount of escalation in competitive pressure between applications to secure network resources.

For Further Reading

- [0] Geoff Huston, “IP Multi-Addressing and Multipath TCP,” *The Internet Protocol Journal*, Volume 18, No. 2, June 2015.
- [1] CAIDA Traffic Analysis,
<http://www.caida.org/research/traffic-analysis/tcpudpratio/>
- [2] <http://en.wikipedia.org/wiki/LEDBAT>
- [3] Van Jacobson and Mike Karels, “Congestion Avoidance and Control,” 1988, <http://ee.lbl.gov/papers/congavoid.pdf>
- [4] W. Richard Stevens, “TCP Slow Start, Congestion Avoidance, Fast Retransmit, and Fast Recovery Algorithms, RFC 2001, January 1997.
- [5] Ethan Blanton and Mark Allman, “TCP Congestion Control,” RFC 5681, September 2009.
- [6] Peter Dodcal, “15 Newer TCP Implementations,”
<http://intronetworks.cs.luc.edu/current/html/newtcps.html>
- [7] FAST: https://en.wikipedia.org/wiki/FAST_TCP
- [8] Google’s QUIC:
<https://www.chromium.org/quic>
<http://blog.chromium.org/2015/04/a-quic-update-on-googles-experimental.htm>
- [9] IETF activity on TCP flow control: TCP Maintenance and Minor Extensions (tcpm)
<https://datatracker.ietf.org/wg/tcpm/charter/>
- [10] Randall Stewart, “Stream Control Transmission Protocol,” RFC 4960, September 2007.
- [11] E. Kohler, M. Handley, and S. Floyd, “Datagram Congestion Control Protocol (DCCP),” RFC 4340, March 2006.
- [12] Geoff Huston, “Anatomy: A Look inside Network Address Translators,” *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

IAB Statement on the Trade in Security Technologies

The *Internet Architecture Board* (IAB) published the following statement on June 15, 2015:

“The Internet Architecture Board is deeply sympathetic with the desire to enhance the security of Internet protocols, infrastructure, and Internet-connected systems. We believe, however, that efforts to enhance Internet security must proceed from a thorough knowledge of the threats against the network, its protocols, and the systems attached to it. Efforts to limit the export or transfer of Internet security technologies seem likely to limit that knowledge in ways that ultimately will frustrate the general goal of a secure and stable Internet.

The identification of vulnerabilities is a fundamental part of security practice. Restrictions on systems which perform that function will make it substantially more difficult for those performing that function to design and deploy secure systems.

Traffic analysis systems, though they may be used in other ways, are a similarly crucial part of the methods used to identify attacks and to analyze the success of remediations put in place. The Internet is a deeply interconnected set of networks that spans international borders, and attacks may occur in one part of the Internet that have extensive ramifications for the operation of the whole. Limiting traffic analysis technologies to specific territories seems likely to hinder efforts to detect and thwart both active threats and other network issues.

We note that in 1996 the IAB and *Internet Engineering Steering Group* (IESG) jointly published RFC 1984^[1], with the following comments on a similar matter, the export of encryption technology:

Export controls on encryption place companies in that country at a competitive disadvantage. Their competitors from countries without export restrictions can sell systems whose only design constraint is being secure, and easy to use.

Usage controls on encryption will also place companies in that country at a competitive disadvantage because these companies cannot securely and easily engage in electronic commerce.

Export controls and usage controls are slowing the deployment of security at the same time as the Internet is exponentially increasing in size and attackers are increasing in sophistication. This puts users in a dangerous position as they are forced to rely on insecure electronic communication.

We believe the same points to be fundamentally true for the export of traffic analysis, penetration testing, and similar security technologies.

While it may appear possible to narrowly circumscribe restrictions so that they target technologies that serve no possible purpose but attack, any modular system, including those intended solely for research, will like have some elements that, divorced from the system, would serve no other purpose. Efforts to target such systems will thus likely sweep up many other security technologies. We therefore recommend that export restrictions on security technologies be generally avoided.”

- [1] IAB and IESG, “IAB and IESG Statement on Cryptographic Technology and the Internet,” RFC 1984, August 1996.

A Primer on IPv4 Scarcity

The April 2015 Issue of the ACM SIGCOMM *Computer Communication Review* contained an excellent summary of the rise and fall of the IPv4 address space^[1]. The authors have managed to be wonderfully concise, packing into just a little over 8 pages a history of the initial address allocation practices, the evolution of needs-based address provisioning through the *Regional Internet Registry* (RIR) framework, and the onset of depletion and exhaustion in the last five years. The paper also reviews the routed address space, and explains the differences between *occupied*, *routed*, and *allocated* address space. It also explains the concept of efficiency of utilization of addresses. The authors consider IPv4 addresses as a resource and the long standing debate over whether addresses can be considered as conventional “property,” as well as the tension between the policies of the various registries and the perspectives of the holders of address space. The paper outlines recent efforts to augment the registry functions with a form of certification allowing third parties to use a *Public Key Infrastructure* (PKI) to validate the authenticity of attestations about addresses and their use, particularly in the context of the Internet’s routing system. The paper details current efforts in coping with an environment where the traditional source of IPv4 addresses has been exhausted, considers address *markets*, and the interplay between efforts to increase the address utilization efficiency in IPv4 and incentives to adopt IPv6.

- [1] Philipp Richter, Mark Allman, Randy Bush, Vern Paxson, “A Primer on IPv4 Scarcity,” ACM SIGCOMM *Computer Communication Review*, Volume 45, Number 2, April 2015.
<http://www.sigcomm.org/sites/default/files/ccr/papers/2015/April/0000000-0000002.pdf>

Corrections

While we are all looking forward to *Terabit* (1000G) Ethernet, the article in IPJ Volume 18, No.1 entitled “Gigabit Ethernet,” contained errors in the table on page 27. Thanks to reader Marcin Cieślak for pointing this out. Here is the corrected version:

Table 2: Media Options for 40- and 100-Gbps Ethernet

	40 Gbps	100 Gbps
1-m backplane	40GBASE-KR4	
10-m copper	40GBASE-CR4	100GBASE-CR10
100-m multimode fiber	40GBASE-SR4	100GBASE-SR10
10-km single-mode fiber	40GBASE-LR4	100GBASE-LR4
40-km single-mode fiber		100GBASE-ER4

Naming nomenclature:

Copper: K = Backplane; C = Cable assembly

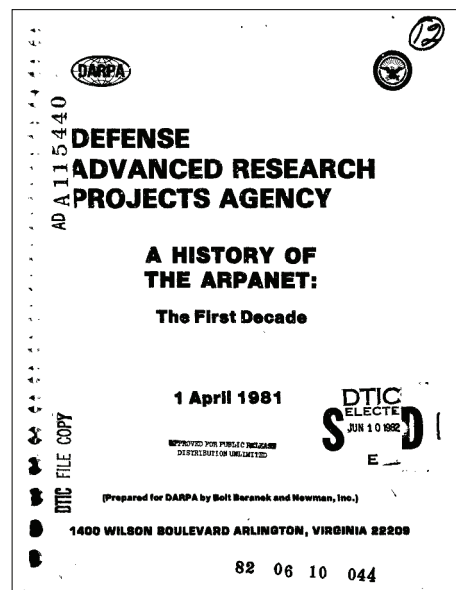
Optical: S = Short reach (100 m); L - Long reach (10 km);

E = Extended long reach (40 km)

Coding scheme: R = 64B/66B block coding

Final number: Number of lanes (copper wires or fiber wavelengths)

Also in Volume 18, No. 1, we told you about Bolt Beranek and Newman Report 4799 entitled “A History of the ARPANET: The First Decade.” It appears that this document is no longer available from the link we gave, so we have placed a copy in the “Downloads” section of our website at protocoljournal.org.



Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Publication of this journal is made possible by:

Supporters



Diamond Sponsors



Ruby Sponsor



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



Individual Sponsors

Lyman Chapin, Steve Corbató, Dave Crocker, Jay Etchings, Martin Hannigan, Hagen Hultzs, Dennis Jennings, Jim Johnston, Merike Kaeo, Bobby Krupczak, Richard Lamb, Tracy LaQuey Parker, Bill Manning, Andrea Montefusco, Tariq Mustafa, Mike O'Connor, Tim Pozar, George Sadowsky, Helge Skrivervik, Rob Thomas, Tom Vest, Rick Wesson.

For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

ADDRESS SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Fred Baker, Cisco Fellow
Cisco Systems, Inc.

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Steve Crocker, Chairman
Internet Corporation for Assigned Names and Numbers

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

September 2015

Volume 18, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
RIPE Atlas	2
Fragments	27
Call for Papers	30
Supporters and Sponsors	31

FROM THE EDITOR

Just over a year ago we relaunched *The Internet Protocol Journal* (IPJ) under a new funding model. Since that time we have published 5 issues of IPJ, designed and implemented a new subscription system, built a website, established new relationships with authors and contractors, and, most importantly, rebuilt our subscriber base with both previous and new subscribers. With some 15,000 print subscribers and 6,600 e-mail subscribers, we are recreating that all-important *community* that makes this publication unique. We've already received many messages of support and appreciation as well as suggestions for topics and article submissions. Please keep them coming!

None of the work behind IPJ would be possible without the support of numerous individuals and organizations. If you or your company would like to sponsor IPJ, please contact us for further details. Our website at protocoljournal.org contains all back issues, subscription information, a list of current sponsors, and much more.

The general topic of *Internet of Things* (IoT) has received much attention in recent years. One way to explain IoT is to describe it as a large collection of Internet-aware *sensors*—such as the temperature sensor in a thermostat, and associated *actuators*—such as the electronic switch that turns your heating or air conditioning unit on or off. However, Internet-aware sensors can also be used to measure details about the Internet itself, and this is the idea behind the *RIPE Atlas* project, which is described in our main article in this issue.

Just as I was finishing writing this editorial, the *American Registry for Internet Numbers* (ARIN) issued the final IPv4 addresses in its free pool. If your organization has not yet deployed IPv6, now would be a good time to start the process. You will find many articles and pointers to further information about IPv6 in back issues of IPJ, all available from our website.

Another reminder that if you received a printed copy of this journal in the mail, you should also have received a subscription activation e-mail. If you didn't receive such a message, it may be because we do not have your correct e-mail address on file. To update or renew your subscription, just send a message to ipj@protocoljournal.org and include your subscription ID. Your subscription ID is printed on the back of your journal.

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

RIPE Atlas: A Global Internet Measurement Network

by RIPE NCC Staff

R IPE Atlas is a global Internet measurement network that provides an understanding of the state of the IP Layer in real time. Here we describe the functionality and design of RIPE Atlas, which collects information about Internet connectivity and reachability via thousands of measurement devices around the world. Then it makes this data available to everyone by measuring the *Internet Protocol* (IP) layer of the Internet with real packets, from anywhere, at any time, by everyone, and for the benefit of all.

- *Real packets:* RIPE Atlas measures the IP layer of the whole Internet—not just the last mile or the network of a single provider. It sends real packets and observes responses so it can directly measure the performance of the IP layer. RIPE Atlas does not focus on measuring applications, which run on top of the IP layer. RIPE Atlas is not an observer of metadata like *Border Gateway Protocol* (BGP) routing traffic; there is no need to make inferences from metadata. And it does not observe any user traffic, thus avoiding many ethical concerns. RIPE Atlas supports five different types of measurements: *ping*, *traceroute*, *Domain Name System* (DNS), *Secure Sockets Layer* (SSL), and the *Network Time Protocol* (NTP).
- *From anywhere:* As of July 2015, RIPE Atlas comprises more than 8,400 active vantage points globally. These vantage points are both small hardware devices called probes that volunteers connect to their home or business networks, and *anchors*, which are generally installed in data centres and act as both enhanced probes with more measurement capacity and regional measurement targets within the greater RIPE Atlas network. Currently, more than 8,300 probes and 133 anchors are deployed in 173 countries and in 11% of all IPv6 *Autonomous System Numbers* (ASNs) and 6% of all IPv4 ASNs. Although this coverage is still far from “everywhere,” it is quite an extensive network—and the architecture allows for scaling up by another order of magnitude.
- *At any time:* Measurements can be started at any time from a chosen subset of these devices and can be performed quickly, or they can be set up to run for weeks, months, or even years. This scenario allows for both quick glances for troubleshooting purposes and deeper analyses of long-term trends. In order to make results comparable, we carefully control the experimental conditions.
- *By everyone:* RIPE Atlas was conceived to collaboratively build and share a huge measurement infrastructure, rather than building individual small ones for exclusive use. Everyone can contribute to RIPE Atlas by hosting a probe or anchor, by building tools to run measurements or analyse results, by providing sponsorship funding, or by helping us distribute probes to difficult-to-reach locations. Everyone contributing to RIPE Atlas can use the network to run their own measurements from virtually all devices.

- *For all:* RIPE Atlas data is openly available, and everyone can benefit from the tools and analyses that others provide. By default, measurement specifications are open for inspection, and all results are accessible to everyone and can be accessed for many years. This continued accessibility means that anyone can reproduce experiments and analyses using RIPE Atlas data, and that is the only way to do real science; it fosters development and the sharing of tools. Having all data openly available also allows others to reuse existing results. RIPE Atlas data consists of a large number of results from a large number of vantage points to a large number of destinations. The topology of the IP layer is far from flat or fully interconnected; thus, by design, the data contains a lot of information about “core” parts of the IP layer. This information can be used to produce maps about the IP layer topology, and that is why we chose the name, *RIPE Atlas*.

In general, RIPE Atlas can be used to:

- Continuously monitor the reachability of a network or host from thousands of vantage points around the globe
- Investigate and troubleshoot reported network problems by conducting ad hoc connectivity checks
- Test IPv6 connectivity
- Check the responsiveness of DNS infrastructure, such as root name servers
- Execute measurements from a large number of vantage points for use in academic research

We describe specific use cases for RIPE Atlas data in more detail in the “Use Cases” section later in this article.

History and Funding

Before RIPE Atlas, the *Réseaux IP Européens Network Coordination Centre* (RIPE NCC) ran numerous other measurement platforms. RIPE Atlas replaced the older *Test Traffic Measurement Service* (TTM)^[1], an active measurement tool geared towards optimising the use of transmission resources that was developed at a time when transmission capacity was much scarcer than it is now. In contrast, RIPE Atlas was developed for an environment where a stable and well-optimised network layer is more critical to the quality of service in the IP layer than squeezing the most out of scarce transmission resources.

The RIPE NCC^[2] began developing RIPE Atlas in 2010 to complement its array of measurement and data-gathering tools^[3] with membership-supported funding. The RIPE NCC membership continues to fund the bulk of RIPE Atlas operations, expansion, and development. Many external sponsors also support RIPE Atlas financially.

Our focus on the IP layer of the entire Internet is very much in line with the needs of our members, mainly Internet Service Providers, who provide the bulk of the funding. In recent years, RIPE Atlas has been extended to measure numerous core services, such as the DNS, in addition to the IP layer. However, this added measurement does not affect our main focus. We are also committed to keeping all measurement data accessible for as long as possible so that this data can be used for future purposes and so that anything based on it can be independently verified at any time.

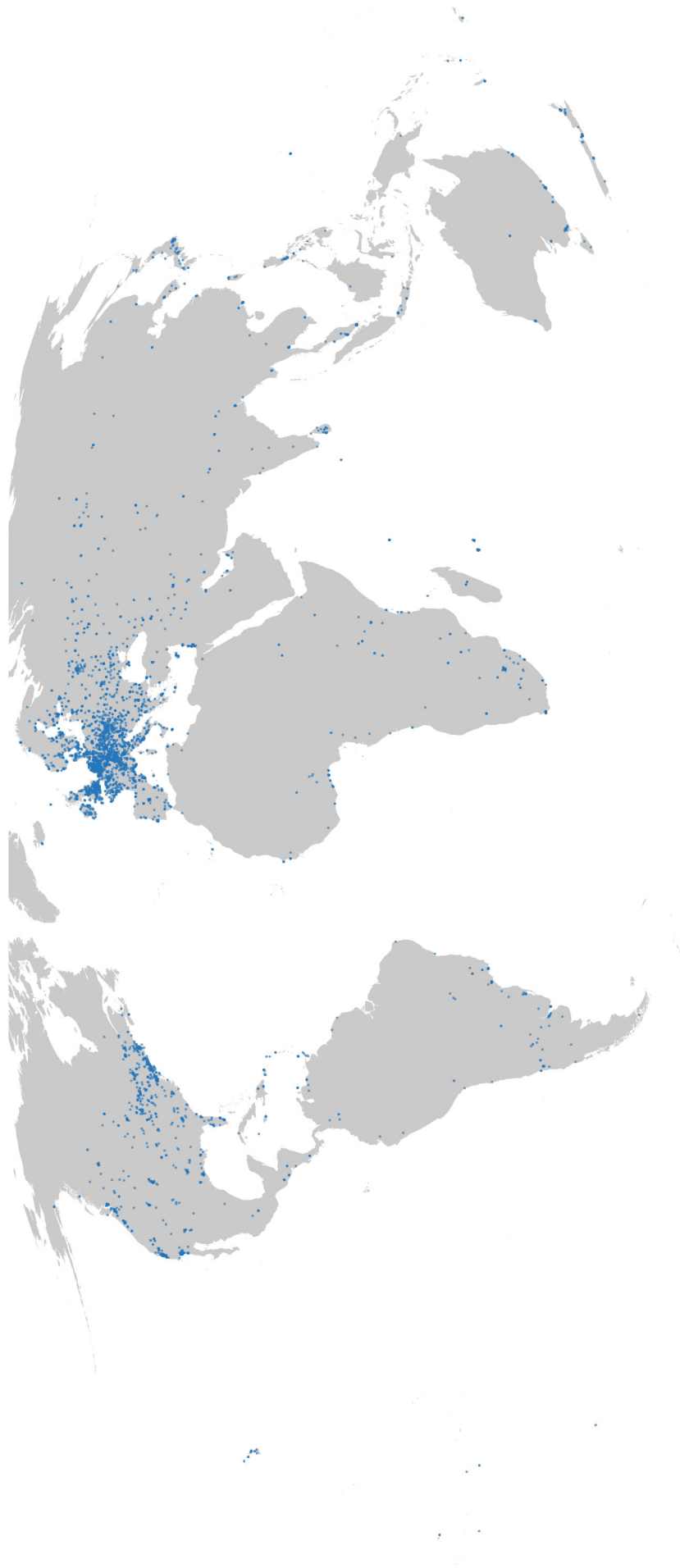
Use Cases

Since the beginning of RIPE Atlas, the network has produced interesting data that everyone can use, regardless of whether or not they host a probe or anchor themselves. Over the last few years especially, the increased reach of the network has meant that researchers, operators, and other RIPE Atlas users have used RIPE Atlas data to debug network problems, analyse network behaviour, and conduct their own interesting research. Many useful visualisations have also been developed based on RIPE Atlas data. Here we list just a few examples:

- A group of researchers from Africa measured inter-domain routing to determine the best possible locations to establish Internet exchange points in the region.^[4]
- A severe power outage and a variety of network outages were analysed and visualised, highlighting how the Internet routes around outages.^[5, 6]
- Engineers from the *Wikimedia Foundation* and the RIPE NCC collaborated on a project to measure the latency of Wikimedia sites and improve performance for users worldwide.^[7]
- Recently, the RIPE NCC developed a tool to analyse how much of a country's local traffic actually leaves the country, and the role that Internet exchange points play in keeping traffic local.^[8]
- A team of researchers investigated content-blocking incidents in Turkey and Russia, as well as outlining ethical considerations when using measurement networks that involve volunteers as hosts, and giving a comparative overview of RIPE Atlas and other measurement networks.^[9]

RIPE Atlas users also produce useful and stunning visualisations that make it easier to grasp the results of RIPE Atlas measurements. For example, one network operator visualised the measurements collected by the local RIPE Atlas anchor, and was able to analyse the quality of the local connectivity and topology changes and help debug network problems.^[10] *CartoDB* helped us visualise network outages and other interesting network behaviour based on RIPE Atlas.^[11] RIPE Atlas also helps make geolocation data of Internet infrastructure more accurate over time. We started a crowd-sourcing project in which operators can specify the geolocation of servers and other equipment based on RIPE Atlas *traceroute* data.^[12]

Figure 1: RIPE Atlas Probes Connected on 2015-07-20 at 05:42:26 UTC ^[23]



The first-ever *RIPE Atlas Hackathon*, held in March 2015, also produced a great variety of visualisations based on RIPE Atlas measurement data, including visualised *traceroute* consistencies over time, a display connecting large-scale probe disconnections to Twitter feeds and weather data, and a map of probe data by country.^[13]

Overall Design

We defined a few fundamental principles early on that influenced the system design: our choice to use dedicated hardware devices as vantage points, the ability to scale up, the distributed and collaborative nature of the deployment, and the related security considerations.

Hardware vs. Software Vantage Points

Even though it incurs measurable costs, we chose dedicated hardware devices as vantage points for many reasons.

The hardware devices are “install and forget,” and in this sense they are not prone to disappearing after an operating system upgrade or a home router replacement, for example. They are easy to deploy because they act as just another device on the network, and after installation they can run uninterrupted, 24/7. Hardware probes also provide a well-defined environment, so the results are much more comparable than widely varying platforms. All in all, these conditions support the goals of achieving datasets with long time series.

We also prefer to not assume responsibility for code that runs on host-provided systems. Any bug or vulnerability in a hardware probe does not affect or compromise another system. If necessary, a probe can easily be disabled by simply unplugging it.

The current generation of probes is well suited to support the functionality of the system. At the moment, the probes are much more constrained by bandwidth limitations, which the hosts can impose themselves, rather than their CPU, memory, or storage capabilities.

Scalability and Resilience

At the time of RIPE Atlas’ conception, about 35,000 *Autonomous Systems* (ASs), were active on the Internet. If we wanted to have vantage points in all of them, account for some of them not being active all the time, and aim for some additional redundancy, we needed to plan for three devices in each network—leading to a rough approximation of 100,000 probes. Of course it’s naïve to think we can install a device in every single AS, plus bigger and geographically more distributed networks would require even more probes. But this number is very useful as a potential upper limit, and it served as a good base for the infrastructure design.

Having this scalability goal does not mean that we immediately deployed enough infrastructure components to handle the full load. The system is in constant evolution in terms of supporting the hardware infrastructure, the redundancy and scalability of critical components, and the software components we use.

The scalability goal also pointed out early on the need for deploying controls on how the measurements and probes behave. Even with much lower deployed numbers, it's not reasonable to run all measurements on all probes all the time, so there is a need for controlling the amount, distribution, and scheduling of system resources. This need gave rise to numerous principles and components, like the credit system and the scheduler, which we describe in the next section.

The infrastructure is designed to operate in a distributed fashion. Most components have enough local knowledge to fulfill their role but don't necessarily know about the state of the whole system or even other components. This setup makes it possible to operate most functions even in the case of a network split. For example, probes keep on executing measurements even if they are not connected to the infrastructure, and all components (including probes) buffer results in case the "next hop" in the result processing chain is unavailable. The communication protocols are designed such that they can handle temporary disruptions. In this sense the system is self-healing; assuming that infrastructure problems can eventually be resolved, the probes and the associated processing pipeline will eventually converge on a stable state and all buffered data will be delivered. The probes are also "headless" (they don't have a console interface), and it's essential for the system to be able to support this kind of recovery.

Collaborative Approach and the Ecosystem

No global measurement network can grow beyond a trivial size without relying on collaborators and fostering an interested community around it.

Our approach is relatively simple: we ask network operators and users to help us reach new networks and deploy probes in them. By installing a probe in a network that previously didn't contain any, probe hosts increase the number of vantage points and allow RIPE Atlas to execute measurements from this new network. Beyond contributing to the larger network, probe hosts also enjoy tangible benefits:

- By keeping the probe active, the host accumulates "credits," a form of currency they can use to perform their own customised measurements using the RIPE Atlas network.
- As soon as a probe becomes active, it starts a set of "built-in" measurements that can point out potential local issues with the host's network.
- Probe hosts can use the credits they earn to execute measurements using any probe in the network, with numerous rules that regulate the use of these resources, as described later in this article.

The costs for measurements are determined such that they are proportional to the complexity of the measurement, the amount of traffic generated, and a few other factors that depend on the specific measurement type.

We also impose a daily limit on the number of credits that users can spend, the number of measurements each user can run, and the number of probes they can use in these measurements. The purpose of the credit system and these constraints is to prevent abuse and overload of the system itself or of popular measurement targets. However, a lack of credits should not prevent any reasonable use of RIPE Atlas. Credits can be transferred between users, and we provide additional credits for tool developers and specific measurement campaigns.

Security Considerations

Since the very idea of RIPE Atlas is to send and receive network traffic from vantage points across the globe, security and reliability have been very important from the beginning of the project.

In order to reduce the attack surface of the probes, we decided to make them as closed as possible. They don't offer any services that users or programs can connect to (in the TCP/IP sense). Instead, they make only outgoing connections, both towards the RIPE Atlas infrastructure and for measurements. They also don't share any authentication tokens among each other; each probe has its own key that it uses to authenticate itself to the infrastructure. Of course, since the probes are physically in the hands of hosts, it's virtually impossible to make them resilient to tampering or disassembly, but these actions should not affect other probes in the network.

All communication within the infrastructure (probes included) is done via encrypted, mutually authenticated channels. In the case of the probes, this communication is through a *Secure Shell* (SSH) protocol connection that is also used to channel control traffic and result traffic back and forth, whereas we use *Transport Layer Security* (TLS) for communication between the control components.

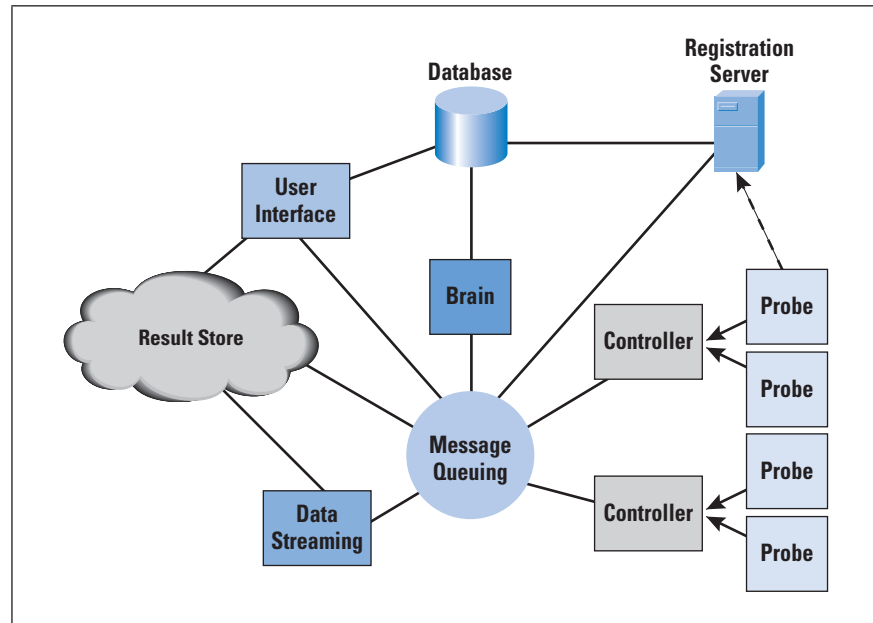
The software running on the probes is field-upgradeable, with the probes being able to verify the authenticity of new firmware via cryptographic signatures. This setup allows us to deploy new functionality as well as security enhancements, so the risk of having a security vulnerability that we are unable to fix is very low.

Of course, no system is 100% resilient against attacks. We run a custom firmware, which may have its flaws, on baseline operating systems (*CentOS*, *OpenWrt*, and *uClinux*). We believe the current security mechanisms provide adequate protection against casual attackers, and the potential for more serious disruptions to the network is limited to sufficiently resourceful attackers, making it quite low.

Architecture and Infrastructure

The RIPE Atlas architecture is designed with the previously mentioned key principles in mind: scalability, resilience, and so on. The general structure is shown in Figure 2.

Figure 2: The Overall System Structure



A central *Structured Query Language* (SQL) system database contains a lot of the key information: almost all of the information about the probes and their various properties; metadata about measurements; users; credits; etc.

A cluster of *message queue servers* acts as the central nervous system for the architecture. It provides connectivity between the various components, and ensures a one-time, guaranteed message delivery with very small delays. It also eliminates the need for each component to maintain knowledge about which other components are present, whether they are connected, etc. This cluster has enough capacity to both deal with control messages and pass on result data, though at a later point we may split these two tasks or change to a more direct data-reporting channel.

Whenever a probe starts up, it first connects to one of many pre-configured *registration servers*. These servers allow only connections from known probes using their specific public keys. When a probe connects, the servers determine which controller the probe should be directed to based on the location of the probe, the availability of controllers, and other factors. The servers then redirect the probe to the chosen controller and notify this controller about the probe. During this process, they inform the probe and the controller about each other's public keys.

The *controllers* accept probe connections from the list of probes they get from the registration servers; they do not have full knowledge of all probes. Once a probe connects, they keep the channel open to receive results and to notify the probes about measurement tasks. Although theoretically each controller can handle a large number of probes, we're trying not to overuse this capability in order to avoid too much fate-sharing between probes when individual controllers or data centres become disconnected.

The *brains* execute higher-level functions in the system, notably the measurement scheduling. This activity happens based on the requests received via the *User Interface* (UI) and *Application Program Interface* (API), and involves an elaborate set of steps about probe preselection, negotiations with the controllers about probe availability, etc. Brains also act as an interface and filtering service between controllers and the central database in order to propagate probe metadata to the UIs.

The UI takes care of all user interactions. It serves the pages for the web interface and handles the API calls. It also serves data download requests via the API; we will likely separate this activity into a separate component if the number of requests or the amount of data transfers creates a bottleneck.

The data streaming servers provide real-time access to the flow of result data, directly off the result queues, and include a few minutes of “short-term memory” to supply recent results. They can also replay results that have already been stored in the result store.

Finally, the *result store* is a Hadoop/HBase cluster for long-term storage of all results. It also handles periodic calculation of aggregates and various other tasks. The technology used allows for relatively cheap redundant storage of the dataset, as well as a simple way of executing *MapReduce*-type batch operations to extract information from datasets that are otherwise too large to download and analyse.

In addition, there are a number of auxiliary subsystems attached to the infrastructure that perform various functions, such as dealing with DNS (“SOS”) requests from probes, managing the administrative tasks, or monitoring and analysing the behaviour of the system in near real time. *RIPEstat*^[14] also closely cooperates with the UIs to visualise the stored datasets.

This architecture allows easy expansion if any one component becomes a bottleneck. For example, it’s very simple to add more controllers to the system to handle additional probes as needed.

Implementation

Over the years, RIPE Atlas has gone through numerous transitions, shaped both by interactions with the wider community and the challenge of managing so much data. At the start of the project, RIPE Atlas was quite static: result data was restricted to data collected by hosts’ own probes, with a handful of visualisations provided by the RIPE Atlas team based on data collected from the (then) small network of a few hundred probes.

Over time, however, we introduced the concept of user-defined measurements, which allowed individual probe hosts to use the growing network to do things we hadn’t initially considered. Slowly, and in consultation with the community, the options for these user-defined measurements grew, allowing for more protocols and more options within them, so that hosts can now take advantage of a global network to perform *ping*, *traceroute*, DNS, SSL, and NTP measurements with a plethora of options.

Interfaces

As we expanded the available features, we needed to change the ways in which users interfaced with the system. Where once we had a very data-heavy ExtJS-based website, we built a more user-friendly interface to access a variety of data visualisations. At the same time, we also gave access to our data to those wanting to create their own views by standardising on a *Representational State Transfer* (REST) ful API, streaming API, and result-parsing toolkit.

We also added a few new sections to the website to improve the user experience:

- A *credits* page with a chart to help users understand how many credits they were earning and spending
- A *keys* page to allow users to create API keys that they can use to create new measurements with the RESTful API or share existing result data
- Separate pages for different groups of users, including “ambassadors” (those who help us distribute probes), financial sponsors, and anchor hosts

The largest change was the way in which we display probe and measurement information. The first instances of this change were cluttered, grid-like interfaces, which we replaced with a searchable listing and highly visual renderings of measurement results and probe performance. Many of these new interfaces are powered by the publicly available APIs.

Lastly, we introduced an updated and much more user-friendly interface for creating measurements. Guided by feedback from users of the old ExtJS-based system, measurement creation became a single form with three sections, of which only the first required direct input. Users are left to decide whether they want to accept the defaults for the other sections or be more specific.

This new measurement form was coupled with numerous helper tools, including:

- A credit consumption visualiser that helps users gauge the long-term expense of their proposed measurement
- An elaborate probe selection wizard that allows users to search, filter, and select probes across networks and geographical space: This selector also implements a standard autocomplete that covers ASN, prefix, IP, address family, and current probe status.
- An API example panel that shows the equivalent request as if it were submitted via the API: This panel was designed to serve as a learning tool for those interested in working with the API for future measurement creation.

This new tool, combined with the availability of the RESTful API, made it much easier to create and manage measurements.

Visualisations

In addition to downloading results in *JavaScript Object Notation* (JSON) format or accessing them via the streaming API, users can also benefit from the visualisations we have created for some measurement types. These visualisations are usually made with JavaScript and *Scalable Vector Graphics* (SVG). The back end provides only the data, while the visualisation and the data correlation are done completely client-side and are accessible from a web browser without third-party plug-ins. They can be embedded in any HTML page or shared by means of permalinks. These visualisations include:

- Simple, tabulated layouts of key parameters extracted from measurement results: These layouts can give an immediate impression about how the target behaves according to the measurement.
- Maps of all kinds showing success rates, round-trip times, or measurement results, where each involved probe and its latest result are denoted by colour, shape, or size: These maps are constructed to work with practically all types of measurements that the system supports.
- Simple charts of packet loss and round-trip times for *ping* measurements, also featuring a shift and zoom function to explore past results.
- A “seismograph” that shows results of *ping* measurements from multiple probes: This tool allows users to compare results from multiple probes at the same time, making it really powerful to spot common behaviour across multiple ASNs, countries, or regions. The tool also supports zooming and exploring historical data.
- A complete rewrite of the RIPE NCC tool *DNS Monitoring* “(DNSMON)”^[15] to visualise long- and short-term behaviour and reachability of root DNS servers and important top-level DNS domains: In many aspects, the RIPE Atlas version of DNSMON goes well beyond the functionality of its predecessor, for example, by incorporating *traceroute* measurements alongside the DNS queries and offering this information to help users diagnose observed performance changes.
- A tool that shows common paths towards particular targets, based on *traceroute* measurements

There is a lot of further potential in providing more visualisations to our users. In particular, we are working on extracting more information from *traceroute* measurements. Combined with the ability to zoom in or out using ASN, prefix, or even geographical aggregations, this information can give an easy-to-understand, immediate explanation of where network bottlenecks or dysfunctional network paths are located.

Data Storage

A global measurement network as big as RIPE Atlas has unique requirements for storing and accessing the collected data.

As of July 2015, about 8,000 concurrently running measurements produce more than 2,700 results per second. Depending on the measurement type, the size of the result can vary from a few bytes to kilobytes. All measurements ever run since the start of RIPE Atlas have been kept. Hence it is not a surprise that the net size of all measurements until this point has grown to 24 TB. On top of this amount is another 4 TB worth of derived data, resulting in a total of 28 TB of active Internet measurement data. Keeping this data online, indexed, and accessible in real time is a challenge on its own.

Storage Solution

After a period of extensive testing and comparing available storage solutions, we decided on *Apache Hadoop*, because it seemed to be scalable and reliable. In addition, at the time of our decision in 2011, it was one of the few open-source solutions capable of dealing with data on a terabyte level. Related to the open-source nature of Hadoop was a choice of different distributions, and we selected the one created by Cloudera (CDH).

Hadoop, maintained as an Apache project, is a software framework that allows for the distributed processing of large datasets. Its design makes it possible to scale from a single server to thousands of machines, with each machine offering local computation and storage. High availability and data security are built into the library itself, so we didn't have to depend on expensive hardware and we could use inexpensive commodity hardware instead.

The Hadoop distributed file system allows for high-throughput access to the application data. Another key component is *MapReduce*, which is a simple programming model that breaks up large datasets and splits data-processing tasks so they can be run in parallel on multiple nodes in a cluster. On top of this stack we use *HBase*, another part of the Hadoop ecosystem, which provides a structured, table-like storage model for large, distributed datasets.

What we described earlier as the “result store” currently consists of one development cluster and two production clusters, one live and one standby. Each production cluster has 119 machines providing a total of 3.5 TB of memory and 952 CPUs. Of the 119 machines, 110 act as workers on which data is stored and jobs are executed. The other 9 machines schedule jobs and manage cluster resources. Each worker can store approximately 4 TB of data, providing a total of 400 TB of storage capacity (after formatting).

Each production cluster is roughly at two-thirds of its capacity, accounting for 260 TB. The cluster also stores other datasets, such as *Border Gateway Protocol* (BGP) routing data collected by *Routing Information Service* (RIS)^[16]. Nevertheless, the storage requirements for RIPE Atlas alone are significant because, in order to achieve reliable storage, Hadoop uses replication, which blows up the requirements for RIPE Atlas data to 75 TB using a replication factor of three.

It might be surprising that we run two production clusters, given that Hadoop has reliability built in, so one cluster might seem to be sufficient. In fact, only by using two clusters were we able to test new features and platform updates without running the risk of long service interruptions in case something went wrong. As with any software, Hadoop is not bug-free, and we especially noticed this reality when we used its first releases just as “big data” was becoming a buzzword. As it moves toward a more stable system, it might be an option in the future to combine the two clusters, resulting in double the current capacity.

Data Flow

From the message queue servers, the measurement data, in the form of *Advanced Message Queuing Protocol* (AMQP) messages, is loaded by daemons running on all worker nodes. The payload, in JSON format, is extracted and stored directly in HBase tables and saved as sequence files. Sequence files are a binary format of the collected data, which act as an efficient input for MapReduce jobs that create tables derived from raw data.

Data derivation is specifically necessary because Hadoop, like many other NoSQL solutions, is a key-value store and hence lacks elaborate concepts of data projection and selection. The key, which is the only available index, becomes very important, and different access use cases require different key layouts, resulting in data duplication.

An often-used derivation is time-based aggregation, which compacts data collected over a long period of time to a fraction of its original size. By extracting statistically relevant values like minimum, maximum, and median, users (via the UI) are still able to get an overview of the measurement result without having to download all of the measurement data.

Initially orchestrated by *cron* jobs, the execution of MapReduce jobs is now done with a software component called *Azkaban*. Azkaban can account for dependencies between datasets and—most importantly—data availability, and the result is much better data quality for the derived datasets.

Next to regularly running jobs, the system also supports ad hoc data processing, which enabled our data scientists to crunch volumes of data that would have never fit on anyone’s workstation.

Datasets

RIPE Atlas data consists of multiple tables. From the size value, one can see the difference between raw and aggregated tables, highlighting how aggregation saves a lot of space and makes overviews much faster.

- *Result blobs*: all messages fetched from the message queue servers; this data is the rawest form accessible as a table and acts as the input for most other tables (24 TB)

- *Latest results*: the latest measurement result per probe and measurement (50 GB)
- *Counters*: information about the amount of results delivered per probe and measurement, also used for credit billing (4 GB)
- *DNS details and aggregates*: all DNS measurements and their aggregates, also the ones used for DNSMON (800 + 200 GB)
- *Ping results and aggregates*: contains all *ping* measurements and all its aggregates used for *ping* visualisations (3.5 TB + 140 GB)
- *Traffic*: data on how much traffic a probe produced, upstream and downstream (12 GB)

The front-end servers (RESTful) gain access to the data via Apache *Thrift*. Thrift is a software framework that allowed us to build seamlessly integrated data types and service interfaces between the Hadoop Java environment and the front-end Python environment.

Measurements

RIPE Atlas performs two distinct sets of measurements. The first category is the “built-in measurements,” performed by all probes as soon as they are connected; it provides a baseline flow of generally useful results. The second category is the “user-defined measurements” that users specify themselves, which usually run on only a small subset of the probes.

Built-In Measurements

When a probe is plugged in, it initiates connection to the RIPE Atlas infrastructure and is connected to a controller, which sends a list of predefined, built-in measurements. These measurements are designed to provide useful data to the host for monitoring their connectivity and discovering potential local issues, but they also provide a wide array of generally useful results for network visualisations.

Built-in measurements include:

- *Pings* to first and second hops
- *Pings* and *traceroutes* to well-known destinations, such as DNS root servers and RIPE Atlas anchors
- DNS SOA, `version.bind`, `hostname.bind`, `id.server`, and `version.server` monitoring on the DNS root servers
- HTTP and SSL requests to some of the RIPE NCC servers

Built-in measurements are performed via both IPv4 and IPv6 if a probe seems capable of supporting both protocols. The measurements are ongoing, with most of them running every few minutes. The list of running built-in measurements is defined solely by the RIPE Atlas team; probe hosts have no influence over it. Data for the built-in measurements is publicly available.

User-Defined Measurements

The distinct feature that makes RIPE Atlas unique and provides value for all participants is its ability to let users schedule their own measurements using virtually all of the probes throughout the network with great flexibility. Users are free to choose:

- Measurement type (from the supported set)
- Type-specific options such as flags and parameters
- IP version
- Measurement target (for example, hostnames or IP addresses)
- Measurement source (that is, vantage points); this source is the set of probes defined by size and desired properties, such as geography, ASNs, prefixes, and tags. It is also possible to reuse the same sets of probes that were used in previous measurements.
- Measurement start/end time and the frequency of ongoing measurements

Each measurement can be ongoing (with a user-defined frequency) or a one-off measurement. One-off measurements are more expensive to handle, but they have a much faster reaction time, delivering results in a matter of seconds.

User-defined measurements can be created and maintained using the web interface, as well as the RESTful API. Users can stop measurements prematurely if they wish, and can also change the set of involved probes.

User-defined measurements are the most resource-consuming activity within RIPE Atlas, and we therefore need to balance the system capacity with our users' needs. To achieve this balance, we constantly increase the number of probes, throughput of the measurement result delivery, and storage. The credit system also plays a large role. It encourages users to keep their probes connected, so they earn credits to use on user-defined measurements, and to use the system with care. The more resource-intensive a measurement (number of probes, probe CPU, network activity), the more expensive it is.

Current Measurement Types

At the time of writing, RIPE Atlas supports the following measurement types:

- *Ping*: monitors network delays using *Internet Control Message Protocol* (ICMP) IP echo messages
- *Traceroute*: displays the route path and measuring transit delays
- *DNS*: queries Domain Name System servers or resolvers (provides users with similar functionality of well-known tools such as *dig* and *nslookup*)
- *SSL*: queries SSL/TLS certificates of remote servers
- *NTP*: queries Network Time Protocol servers, dumps the reply, and computes some statistics, such as reply time

Every measurement type allows the user to specify parameters for fine-tuning. For example, *ping* allows the user to change the number and size of packets, while DNS has a much wider set of options, allowing users, for example, to set query parameters, recursion, number of retries, and so on.

Future Measurement Types

We constantly work with the RIPE Atlas community to satisfy our users' demands for different measurement functionality. We have started to add new measurement types and continue to add new options for fine-tuning them, as long as they are consistent with the RIPE Atlas mission. The following measurements are likely to be introduced in the near future:

- *HTTP*: queries HTTP servers. The system is technically already capable of doing this task; it is used internally to deliver results, but it's not yet available for public use. Public availability will most likely be restricted to allowing users to perform GET and HEAD requests against a predefined list of targets (initially RIPE Atlas anchors).
- *WiFi*: an intentionally opt-in-only feature to verify the functionality of *Wireless LAN* (WLAN) access points: The intention is to check whether it's possible to connect to a requested, specific *Service Set Identifier* (SSID), and measure the IP connectivity of this interface. We do not intend to use WLAN connection for IP connectivity of the probe itself, and we don't plan to allow passive WiFi scans.

Day-to-Day Development

RIPE Atlas is developed by a team of software engineers with overlapping focuses. Although we don't follow a particular methodology, RIPE Atlas development is agile in the sense of iteratively delivering and improving upon working features, responding to feedback from the community, and continuously focusing on good design and architecture. Development is modular, with each type of component having its own *Git* repository, automated deployment mechanism, and release schedule. The frequency of these release cycles varies according to development intensity. For example, the UI has the highest rate of change and so has a new release almost every week—with larger stretches for fundamental changes that possibly include bug fixes and hot fixes in the interim—while the brain may go weeks between new features and releases.

In addition to the main RIPE Atlas system, two internal environments are fully functional, each made up of the various components necessary to make a working system. The first is a development environment, which is used to test new features, either directly or as a verification stage after modular, test-driven development. Test-driven development is particularly important and powerful when making changes that are cumbersome or time-consuming to test with a real RIPE Atlas network—even a scaled-down development version—where real components have to “talk” to each other at network speeds.

In such cases, it is helpful to write unit tests for each part that represents the desired or correct interaction between components, and to write code until the tests pass. The second environment is a test or preproduction network that runs the next release while it is being prepared for production. Depending on the magnitude of a release, this test can take anywhere from less than a day to a week or more.

Each of these environments, including the production system, has a suite of system tests that verify a range of behaviours, including logging in to the website, creating measurements, and downloading results. This suite of tests acts as an end-to-end sanity check and can catch things that manual testing and unit tests might miss, or issues that may crop up only after unpredictable real-world use. This suite is complemented by a dashboard, which enables interactive visualisation and analysis of patterns and spikes in error rates, and various bespoke statistical graphs.

Probe Hardware Experiences

The first version of RIPE Atlas probes was very limited. The probe was based on a Lantronix XPort Pro module.^[17] This module contains a FreeScale MCF5208 ColdFire processor, 8 MB of main memory, and 16 MB of flash storage. It is a 32-bit CPU with no *Memory Management Unit* (MMU). The module was complemented with a power board that takes power from USB (5V) and converts it to the 3.3V the module needs. The power board is then enclosed in a small black case.

Black-Box Model

The limitations of these first-generation probes led to a black-box model. The probes have no buttons; the Ethernet connector has link and activity LEDs and nothing else. Because of the limited amount of main memory, it would be difficult to run a web server on the probes for configuration.

The attractiveness of this model comes from the fact that a probe host has to connect only the USB connector to a USB port for power and plug in an Ethernet cable for network connectivity. The probe does not run any services, which is also good for security. Because the probes do not cost the hosts anything, we wanted to make them difficult to tamper with to avoid having people re-purpose them for their own uses.

Automatic Firmware Upgrades

Automatic firmware upgrades are another feature of the black-box model. Probe hosts do not have to care about upgrading the probe firmware. At first glance, firmware upgrades are easy enough. The flash memory of the probe is split into two parts; one part contains the firmware, and the other is used to temporarily store measurement results.

During a firmware upgrade, the new firmware is written to the partition that contains the measurement results, and the boot configuration is changed to boot the other partition.

Automatic firmware upgrades mean that almost all bugs can be fixed in later firmware versions, except those that prevent a firmware upgrade in the first place. We also want to use probes for as long as possible, so probes that may have been lying in a drawer for a long period of time should be able to upgrade to the latest firmware. This stipulation places rather tight restrictions on the backward compatibility of the firmware upgrade process.

One question that arises is whether it would be possible for an attacker to upgrade probes with malicious firmware. With the original firmware, this question meant attacking one of a few core servers or the SSH connection between the probe and the registration server. To further reduce the possibility of this attack, probes now verify a digital signature attached to the firmware. The digital signature is made offline and with secret splitting to make sure that multiple people are involved in signing the firmware.

Debugging a Black Box

A question that arises rather quickly is how to debug a black box. What do you do with a probe in a remote location that has no display and no buttons to press? Just about the only interaction possible is to power-cycle the probe. Because probes communicate over the network, we can distinguish two broad classes of problems: ones that prevent the probe from reporting results and ones that do not.

The second class can be solved most of the time by ensuring an appropriate amount of logging and taking advantage of the scale of the system. Over time, logging got more and more refined as a result of lessons learned trying to debug problems. For example, probes regularly log free memory and free disk space, making it possible to spot memory leaks.

For the first class, we have a few techniques and tricks. The first is that probes perform special DNS queries (called “SOS”) when they reboot; the queries include the message they want to send as part of the hostname they look up. In some cases, a probe can contact a registration server but nothing else. In that case, we can direct the probe, for example, to an IPv4 address literal if we suspect DNS problems or force the probe to use either IPv4 or IPv6 if we suspect that the problem is with one of the protocols.

All of these tasks can be done without involving the probe host, although they can assist in two ways. Some hosts can run *tcpdump* on their routers, in some cases possibly giving us a clue about what went wrong. The probe also logs details about what goes wrong when it tries to connect, so connecting the probe to a different network allows these logs to be reported. It should be noted that connecting a probe to a home router is more likely to work than most office or data centre setups, because home routers have rather simple dynamic configurations and a lack of firewalls compared to more strictly managed environments.

Static Network Configuration

Most probes obtain their network configuration dynamically, using *Dynamic Host Configuration Protocol* (DHCP) for IPv4 and *Router Advertisements* (RAs) for IPv6. In some networks, however, those services are not available. To support those networks, probes can also receive network configuration from the back-end servers, but this process creates a “Catch 22” because a probe needs network access to be able to obtain network configuration from the back-end servers.

To solve this problem, the probe first has to be connected to a network that dynamically configures the probe. The probe can then fetch the network configuration and store it in flash memory. After a static network configuration successfully communicates to the probe this way, the probe can be moved to the target network.

However, this process introduces two problems. The first is that static network configuration has to be carefully copied across firm-ware upgrades. The second is that if somehow the static network configuration does not work, the probe has to revert back to dynamic network configuration. This reversion is quite tricky to achieve in practice and, over time, quite a few bugs have caused problems with it. In many cases, they required carefully tweaking the affected probes to get them back on track. Even though this feature adds a lot of complications, it also allows deployment in networks that would otherwise be inaccessible for us, so we expect to support this model in the long run.

It’s worth noting that static configuration can also cause unexpected problems. We’ve seen many cases in which the addresses of DNS resolvers defined for the local network changed over time, but the probe was never informed about it. This change leaves the probe in an inconsistent state that requires manual intervention. For this reason, we do not recommend that hosts use static configuration unless there’s no other solution.

Network Problems

One challenge with the black-box model is that debugging problems typically requires the probe to have a working network connection. Unfortunately, misconfigured networks are among the most common problems.

One problem that is relatively easy to work around is misconfigured DNS. Probes have the names and public keys of two registration servers, which are used to bootstrap the connection to the back-end servers. By adding IPv4 and IPv6 addresses to the names, the probes can connect even if DNS resolution fails. The probes try addresses and names in random order until one succeeds. Obviously, DNS measurements are still doomed on such a probe.

It is possible that the probe simply did not get a lease from DHCP, or got the wrong address or default router, etc. In such a case, the infrastructure side can't detect any life signs from the probe, and it is entirely up to the probe host to resolve the problem. Something similar applies to firewalls; sometimes it is possible to spot a firewall if the DNS traffic from the probe is detected, but we see no SSH traffic.

Finally, there are *Path MTU Discovery* (PMTU) problems. IPv4 typically does not have PMTU problems, but IPv6 does. A quick hack to avoid the IPv6 PMTU problems then causes problems with some routers used by probe hosts that make mistakes in *TCP Maximum Segment Size* (MSS) clamping.

All of these issues make troubleshooting a probe that fails to connect a bit of a trial-and-error game. Again, it sometimes helps to ask the probe host to move the probe to a new network so that it can report what it logged about the first network.

Different Probe Versions

One problem caused by the lack of an MMU on the first-generation probes was memory fragmentation. On systems with an MMU, main memory fragmentation is no problem because it does not cause fragmentation of the virtual memory of a process. On systems without an MMU, fragmented main memory cannot be hidden and, essentially, the only way to deal with memory fragmentation is to reboot the probe. Certainly with early firmware versions, the 8 MB of memory would quickly fragment and probes would reboot, sometimes within one day.

Fortunately, Lantronix was responsive to this issue and released a version with 16-MB main memory. These devices became the second-generation probes.

However, the CPU in the version 1 and 2 probes remained slow, with it taking about 30 seconds to set up an SSH connection. The 8 MB that is available for storing measurement results was not a lot if the probes lost their connection for more than a few days. In addition, off-the-shelf travel routers are both more powerful and cheaper than the Lantronix modules, which are designed for industrial applications.

This situation led to the third-generation probes, which are based on the TP-Link TL-MR3020 travel routers (see Figure 3). They have a much faster CPU and 32-MB main memory, but only 4-MB flash memory. Fortunately, a USB connection was available where we could plug in a USB flash device, although this need required a redesign of the firmware. During normal operation, the probe runs from the external USB flash, which is encrypted to make the probe more tamper-resistant. On the built-in flash is a small amount of code that can switch to the external USB and can also fetch firmware over the network and write it to the USB flash device.

This approach has a disadvantage in that it is a lot more complex than what runs on the version 1 and 2 probes. However, an advantage is that, if the USB flash becomes corrupt or nonfunctional because of a bug, the built-in flash version can simply overwrite it with a fresh copy. Upgrading the built-in flash is also possible, but ensuring that all combinations work proves to be quite tricky.

Finally, a need for high-capacity probes that would be rack-mounted in data centres resulted in the anchors. Anchors are rack-mounted PCs running CentOS Linux. From a firmware point of view, the main differences are that the probe code runs as a regular application and all network management is left to normal CentOS configuration. To the host, an anchor is a black box just like the regular probes, but remote-management cards make it possible to configure the anchors directly. In contrast to normal probes, anchors also run various basic services, allowing them to act as measurement targets as well.

Figure 3: Atlas Probe Version 3



Hardware Problems

With the version 1 and 2 probes, hardware problems were, and remain, rare. A very small fraction of the devices have developed hardware failures, and those failures fall into three groups: power failures, Ethernet failures, and flash-memory failures. We have not looked into it deeply, but power failures are likely a failure of the power regulation board. A few probes developed problems with the Ethernet interface; typically, they can send packets but cannot receive them. Finally, with some probes, the flash memory has broken.

The situation is quite different with the version 3 probes. One of the first things that became apparent is that some USB ports used to power the probes did not supply enough power. The curious effect is that only the external USB flash device refuses to start, leaving the probe with a configuration as if no USB flash device is plugged in. We solved this problem by having the probe perform special DNS queries to indicate that it cannot find a USB flash device.

Some time later, the USB sticks once again proved to be a weak link, although we should note that the number of failures is a few dozen out of thousands of probes. The USB sticks sometimes fail to work in various ways: some become read-only, while others reset their configuration to a default state. Typically, they then show a capacity of only 64 MB and have simple patterns as serial numbers. Finally, actual data corruption has occurred in a few cases.

Far less frequently than broken USB flash devices, sometimes the TP-Link devices develop power failures or problems with the Ethernet interface.

IPv4 and IPv6

Probes support both IPv4 and IPv6; however, IPv4 and IPv6 differ in many subtle ways. One feature that does not occur in IPv6 (or at least we are not aware of any probe in the network that deploys it) is *Network Address Translation* (NAT)^[20]. When a probe reports a measurement result, the local IPv4 address is, in many cases, a local address defined by RFC 1918^[21]. The RIPE Atlas system tries to keep track of what the corresponding public IPv4 address is, and the back end inserts that address in the measurement results when they come in. For IPv6, that complexity is not implemented.

At first, all probes with a global IPv6 address were assumed to be IPv6-capable. However, it turns out that a significant fraction of those addresses have a *Unique Local Address* (ULA)^[22] and have no actual Internet connectivity. A second issue with IPv6 is that a probe may have multiple IPv6 addresses. The problem here is that, without special measures, it is not clear which address a probe will use in a measurement. A probe may use different addresses depending on the measurement target.

It is attractive to think of a probe being in a certain AS. However, each address can be in a different AS, meaning a probe can have different ASs for its IPv4 and IPv6 addresses, and a probe with multiple IPv6 addresses may, at least theoretically, also have different ASs for each address.

RIPE Atlas Community

From the beginning we knew that, in order to succeed, RIPE Atlas would need a strong community around it, both to help us grow the network by hosting probes and anchors and to help spread the word about the usefulness of the project for network operators, researchers, and interested users. We also rely on the Internet community and anyone using RIPE Atlas to give us feedback about new features we develop, useful functionality they would like to see, and the future direction of the project.

We regularly update the community about the development of RIPE Atlas via articles on RIPE Labs^[18] that explain new features and functionality and ask for feedback.

We also publish different use cases and analyses that employ RIPE Atlas data, and have a special collection on RIPE Labs^[19] that includes articles and blog posts written by external RIPE Atlas users about their own experiences, as well as scientific papers based on RIPE Atlas data, and presentations about RIPE Atlas given by members of the community at various conferences.

The RIPE Atlas network is essentially a volunteer-based project; it relies almost entirely on a community of probe and anchor hosts to install the hardware devices in their own networks. Most probe hosts initially heard about RIPE Atlas from the RIPE NCC directly, either via our website, mailing lists, or at various conferences, and applied online for their probe.

Initially, we distributed probes via post, free of charge, to anyone who applied, with the goal of reaching critical mass. As we've come closer to reaching our goal of 10,000 probes in the past year, we've started being more selective in distributing probes and have employed checks to ensure that probes are distributed to ASNs that don't already have a probe connected within them (although RIPE NCC members can receive a probe regardless of this stipulation). Anchor hosts also play an important role in strengthening the RIPE Atlas network and boosting its capacity, and, unlike regular probe hosts, they contribute to the project financially by purchasing the required hardware.

We also rely heavily on our ambassadors—currently more than 240 volunteers—who help us distribute probes and give presentations about RIPE Atlas at conferences all over the world. The RIPE NCC also has partnerships with the other *Regional Internet Registries* (RIRs) to help distribute probes in their service regions.

Each year, numerous organisations support RIPE Atlas by contributing to the project financially, and we are grateful to them for their support.

The RIPE Atlas website contains numerous “Community” pages that highlight new probe hosts, those hosts with the largest number of measurements and credits spent, anchor hosts, sponsors, and the different conferences where ambassadors will be available to answer questions and hand out probes.

RIPE NCC members enjoy several special benefits, including receiving extra credits for their own use and having access to a recently developed webinar on advanced use of RIPE Atlas measurements. During the webinar they have the opportunity to learn from and interact directly with RIPE Atlas developers.

In March 2015, we also hosted the first *RIPE Atlas Hackathon*, a three-day event in which developers, designers, computer science students, and open data enthusiasts were invited to use RIPE Atlas data to develop useful, creative, and stunning visualisations for the benefit of the entire Internet community.

The hackathon produced some very promising results, and we hope to host more such events in the future.

The past five years have been very exciting for the RIPE Atlas development team. It has been interesting to see the system grow from an idea into more fully realised concepts, then a prototype, and finally into a full service. RIPE Atlas is certainly still evolving, and its continued development is very much based on a steady stream of ideas and suggestions from the community.

It has also been nice to see how more and more people—network operators and researchers, but also regular Internet users—have started using RIPE Atlas to measure to their heart’s content. We are happy to see people build tools for their specific uses and share them. We are very grateful to all our users, probe and anchor hosts, sponsors, ambassadors, contributors and everyone who has helped build this network—and hope we can count on your continued involvement for the benefit of the entire Internet community. You can get involved with RIPE Atlas by visiting this website:

<https://atlas.ripe.net/get-involved/>

Conclusion

RIPE Atlas is a globally deployed tool to actively measure the IP layer of the Internet. It is a collaboration of many, led by the development team at the RIPE NCC. It is useful for both ad hoc observations and long-term data collection. The number of its permanently operating vantage worldwide points stands at more than 8,000 and is constantly increasing. Everyone can contribute, and anyone who does contribute can use the entire system worldwide. Measurement results are stored for many years and thus everything based on RIPE Atlas results can be independently verified.

References

- [1] “Test Traffic Measurement Service (TTM),”
<http://ttm.ripe.net/>
- [2] RIPE Home Page, <https://www.ripe.net/>
- [3] RIPE Analyse Start Page, <https://www.ripe.net/analyse>
- [4] Roderic Fanou, “On the Diversity of Interdomain Routing in Africa,”
https://labs.ripe.net/Members/fanou_roderick/on-the-diversity-of-interdomain-routing-in-africa
- [5] Andreas Strikos, “Amsterdam Power Outage as Seen by RIPE Atlas,”
https://labs.ripe.net/Members/andreas_strikos/amsterdam-power-outage-as-seen-by-ripe-atlas
- [6] Robert Kisteleki, “The AMS-IX Outage as Seen with RIPE Atlas,”
<https://labs.ripe.net/Members/kistel/the-ams-ix-outage-as-seen-with-ripe-atlas>

THE RIPE NCC is the Regional Internet Registry for Europe, the Middle East, and parts of Central Asia. As such, it allocates and registers blocks of Internet number resources to its membership in this region, which consists mainly of Internet service providers, telecommunication organisations, and large corporations. It is a not-for-profit organisation that works to support the RIPE (*Réseaux IP Européens*) community and the wider Internet community. The following RIPE NCC staff have contributed to this article: Daniel Karrenberg (overall vision and probe firmware); Robert Kistelevi (architecture and team leader); Antony Antony and Philip Homburg (probe firmware); Christopher Amin, Massimo Candela, Viktor Naumov, Daniel Quinn, Andreas Strikos, Johan ter Beest and Christian Teuschel (infrastructure); Emile Aben and René Wilhelm (research); Suzanne Taylor Muzzin (communications); Mirjam Kühne and Vesna Manojlovic (community building). The authors can be contacted at: atlas-ipj@ripe.net

- [7] Emile Aben, “How RIPE Atlas Helped Wikipedia Users,” <https://labs.ripe.net/Members/emileaben/how-ripe-atlas-helped-wikipedia-users>
- [8] Emile Aben, “Measuring Countries and IXPs with RIPE Atlas,” <https://labs.ripe.net/Members/emileaben/measuring-ixps-with-ripe-atlas>
- [9] Collin Anderson and Philipp Winter, “Global Network Interference Detection Over the RIPE Atlas Network,” <https://www.usenix.org/conference/foci14/workshop-program/presentation/anderson>
- [10] Salim Gasmi, “Visualising RIPE Atlas Anchor Measurements,” https://labs.ripe.net/Members/salim_gasmi/visualising-ripe-atlas-anchor-measurements
- [11] Emile Aben, “Visualising Network Outages With RIPE Atlas,” <https://labs.ripe.net/Members/emileaben/visualising-network-outages-with-ripe-atlas>
- [12] Emile Aben, “Infrastructure Geolocation – Plan of Action,” <https://labs.ripe.net/Members/emileaben/infrastructure-geolocation-plan-of-action>
- [13] Vesna Manojlovic, “RIPE Atlas Hackathon Results,” <https://labs.ripe.net/Members/becha/ripe-atlas-hackathon-results>
- [14] RIPEstat, <https://stat.ripe.net/>
- [15] RIPE DNSMON, <https://atlas.ripe.net/dnsmon/>
- [16] RIPE Routing Information Service (RIS), <http://ris.ripe.net/>
- [17] Lantronix XPort Pro, <http://www.lantronix.com/device-networking/embedded-device-servers/xport-pro.html>
- [18] RIPE Labs, <http://labs.ripe.net>
- [19] RIPE Atlas User Experiences, <http://labs.ripe.net/atlas/user-experiences>
- [20] Geoff Huston, “Anatomy: A Look Inside Network Address Translators,” *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [21] Daniel Karrenberg, Yakov Rekhter, Eliot Lear, and Geert Jan de Groot, “Address Allocation for Private Internets,” RFC 1918, February 1996.
- [22] Brian Haberman and Robert M. Hinden, “Unique Local IPv6 Unicast Addresses,” RFC 4193, October 2005.
- [23] The map on page 5 shows the world using a *Mollweide Projection*, https://en.wikipedia.org/wiki/Mollweide_projection

Fragments

Rob Blokzijl Receives 2015 Postel Service Award

The Internet Society recently announced that its prestigious *Jonathan B. Postel Service Award* was presented to Rob Blokzijl for his pioneering work, 25 years of leadership at *Réseaux IP Européens* (RIPE), and for enabling countless others to spread the Internet across Europe and beyond. Dr. Blokzijl was selected by an international award committee, comprised of former Jonathan B. Postel award winners, which placed particular emphasis on candidates who have supported and enabled others in addition to their own specific actions.

During the 1980s, Dr. Blokzijl was active in building networks for the particle physics community in Europe. Through his experience at the *National Institute for Nuclear and High Energy Physics* (NIKHEF) and *The European Organization for Nuclear Research* (CERN), he recognized the power of collaborating with others building networks for research and travelled worldwide to promote cooperation across networkers. In the 1990s, Dr. Blokzijl was influential in the creation of the Amsterdam Internet Exchange, one of the first in Europe. His most widely recognized contribution is as founding member and 25-year chairman of RIPE, the European open forum for IP networking. Dr. Blokzijl was also instrumental in the creation of the *RIPE Network Coordination Centre* (RIPE NCC) in 1992, the first Regional Internet Registry in the world.

“Rob’s technical expertise and tireless work had a profound impact on the development of the Internet as we know it today,” said Kathy Brown, President and Chief Executive Officer of the Internet Society. “Beyond the breadth of his technical contributions, Rob is known across the Internet community for his strong leadership and unwavering commitment to collaboration and cooperation, exemplifying the spirit of this award.”



*Rob Blokzijl and his wife
Lynn at the IETF 93 Plenary*

© 2015 Stonehouse Photographics/
Internet Society

The Postel Award was established by the Internet Society to honor individuals or organizations that, like Jon Postel, have made outstanding contributions in service to the data communications community. The award is focused on sustained and substantial technical contributions, service to the community, and leadership.

The Internet Society presented the award to Dr. Blokzijl, including a US\$20,000 honorarium and a crystal engraved globe, during the 93rd meeting of the *Internet Engineering Task Force* (IETF) held in Prague, Czech Republic, July 19–24, 2015.

The Internet Society (www.internetsociety.org) is the trusted independent source for Internet information and thought leadership around the world. It is also the organizational home for the IETF. With its principled vision, substantial technological foundation and its global presence, the Internet Society promotes open dialogue on Internet policy, technology, and future development among users, companies, governments, and other organizations. Working with its members and Chapters around the world, the Internet Society enables the continued evolution and growth of the Internet for everyone.

IESG Statement on Maximizing Encrypted Access To IETF Information

The *Internet Engineering Task Force* (IETF) has recognised that the act of accessing public information required for routine tasks can be privacy sensitive and can benefit from using a confidentiality service, such as is provided by *Transport Layer Security* (TLS).^[1] The IETF in its normal operation publishes a significant volume of public data (such as Internet-drafts), to which this argument applies. The IETF also handles non-public data (such as comments to *NomCom*, the IETF's nominating committee) that requires confidentiality due to the nature of the data concerned.

The *Internet Engineering Steering Group* (IESG) and the broader community^[3] have further concluded that there can be other harmful effects in continuing to access public data as clear-text. Recent massive-scale man-on-the-side intermediary attackers have been seen to take advantage of the absence of security to mount active attacks that would be more difficult had a transport security mechanism such as TLS been used.^[2, 4]

The IESG has therefore agreed that all IETF information must, by default, be made available in a privacy-friendly form that matches relevant best current practices. Further, all future embedded interactions with the IETF (such as `<a>` tags in HTML) should default to causing access via that privacy-friendly form. For content currently accessed using the HTTP protocol, using HTTPS URIs and appropriate TLS cipher-suites^[5] will be the preferred access mechanism, however this direction encompasses more than HTTP traffic alone.

However, as there may be tools affected by this, and recognising that there are a number of IETF participants who prefer to continue to access materials via cleartext, or who have issues with using standard confidentiality services, the IESG are also requiring that public information continue to be made available in clear, for example via HTTP without TLS.

The changes caused by this statement should only need operational systems work and should be transparent to almost all consumers of IETF information. There are a small number of cases where these changes might cause some issues, for example, the current Internet-Draft boilerplate text, which uses the `http:` URI scheme. The IESG will work with the broader community, tools teams, and IETF Secretariat to make these adjustments while minimising disruption to the community.

Note that the “secure/privacy-friendly as the default according to best practices” principle set out in this statement applies to all IETF information, regardless of the protocol used to access that information.

References

- [1] Stephen Farrell and Hannes Tschofenig, “Pervasive Monitoring Is an Attack,” RFC 7258, BCP 188, May 2014.
- [2] Bill Marczak, Nicholas Weaver, Jakub Dalek, Roya Ensafi, David Fifield, Sarah McKune, Arn Rey, John Scott-Railton, Ronald Deibert, and Vern Paxson, “China’s Great Cannon,” <https://citizenlab.org/2015/04/chinas-great-cannon/>
- [3] Richard Barnes, “Deprecating Non-Secure HTTP,” Mozilla Security Blog, <https://blog.mozilla.org/security/2015/04/30/deprecating-non-secure-http/>
- [4] Nicholas Weaver, “A Close Look at the NSA’s Most Powerful Internet Attack Tool,” *WIRED*, March 13, 2014. <https://www.wired.com/2014/03/quantum/>
- [5] Yaron Sheffer, Peter Saint-Andre, and Ralph Holz, “Recommendations for Secure Use of Transport Layer Security (TLS) and Datagram Transport Layer Security (DTLS),” RFC 7525, BCP 195, May 2015.

IANA Transition Update

As announced in our March issue (Volume 18, No. 1), the United States *National Telecommunications and Information Administration* (NTIA) announced its intent to “...transition Key Internet Domain Name Functions to the global multistakeholder community” in March 2014. For the latest information on this process, please visit <https://www.icann.org/stewardship>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Publication of this journal is made possible by:

Supporters



Diamond Sponsors



Ruby Sponsor



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



Individual Sponsors

Lyman Chapin, Steve Corbató, Dave Crocker, Jay Etchings, Martin Hannigan, Hagen Hultzs, Dennis Jennings, Jim Johnston, Merike Kaeo, Bobby Krupczak, Richard Lamb, Tracy LaQuey Parker, Bill Manning, Andrea Montefusco, Tariq Mustafa, Mike O'Connor, Tim Pozar, George Sadowsky, Scott Seifel, Helge Skrivervik, Rob Thomas, Tom Vest, Rick Wesson.

For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

ADDRESS SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Fred Baker, Cisco Fellow
Cisco Systems, Inc.

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Steve Crocker, Chairman
Internet Corporation for Assigned Names and Numbers

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

December 2015

Volume 18, Number 4

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Internet of Things: Network and Security Architecture.....	2
RFCs Beyond ASCII	25
Fragments	28
Call for Papers	30
Supporters and Sponsors	31

FROM THE EDITOR

The Internet continues to grow at an amazing pace. WiFi access to the Internet is now almost standard in every hotel room, business, home, and even onboard many airplanes. Mobile operators are offering improved Internet access for smartphones and tablets as well as “personal hotspots” through deployment of emerging *Long Term Evolution* (LTE) standards. Several of the *Regional Internet Registries* (RIRs) have now depleted their available IPv4 address pool, and IPv6 deployment is taking place in most parts of the world.

This growth of the Internet is not only the result of more *people* using the network, but also the result of more Internet-aware *things* being connected. These things can be anything from traditional computer systems and peripherals to home security systems, vehicles, sensors of all kinds, and even light bulbs. Collectively referred to as *The Internet of Things* (IoT), this emerging area of technology is receiving much attention, and efforts are underway to standardize the communication protocols used in IoT. William Stallings gives an overview of these efforts in our first article.

As with most emerging technologies, the potential for failure, misuse, and just “bad design” is always present. The article “The Internet of Stupid Things” by Geoff Huston illustrates some of the challenges in IoT. The article is available from APNIC’s website:

<https://labs.apnic.net/?p=620>

According to the RFC Editor website, the *Request for Comments* (RFC) series “...contains technical and organizational documents about the Internet, including the specifications and policy documents produced by four streams: the *Internet Engineering Task Force* (IETF), the *Internet Research Task Force* (IRTF), the *Internet Architecture Board* (IAB), and *Independent Submissions*.” Since 1969 these documents have been produced in ASCII-only format without the ability to include drawings (other than so-called “ASCII Art”) or other typographical refinements. This situation is about to change with the introduction of a new format for RFCs. We asked Heather Flanagan, the RFC Series Editor, to give us an overview of this effort.

We would love your feedback on anything you read in this journal. With your permission we can include your comments in the form of a Letter to the Editor, or you may consider writing a Book Review. Send your message to ipj@protocoljournal.org

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

The Internet of Things: Network and Security Architecture

by William Stallings, Independent Consultant

The *Internet of Things* (IoT) is the latest development in the long and continuing revolution of computing and communications. Its size, ubiquity, and influence on everyday lives, business, and government dwarf any technical advance that has gone before. IoT is a term that refers to the expanding interconnection of smart devices—ranging from appliances to tiny sensors. A dominant theme is the embedding of short-range mobile transceivers into a wide array of gadgets and everyday items, enabling new forms of communication between people and things, and between things themselves. The Internet now supports the interconnection of billions of industrial and personal objects, usually through cloud systems. The objects deliver sensor information, act on their environment, and in some cases modify themselves, to create overall management of a larger system, like a factory or city^[1].

The “things” in IoT are primarily deeply embedded devices, characterized by narrow bandwidth, low-repetition data capture, low-volume data usage. These devices communicate with each other and provide data via user interfaces. Some embedded appliances in the IoT, such as high-resolution video security cameras, video *Voice over IP* (VoIP) phones, and a handful of others, require high-bandwidth streaming capabilities. But countless products simply require packets of data to be intermittently delivered.

This article provides an overview of IoT, and then looks at IoT network and security architectures that will help guide the design, implementation, and deployment of IoT.

Background

The evolving Internet involves billions of objects that use standard communications architectures to provide services to end users. This evolution provides new interactions between the physical world and computing, digital content, analysis, applications, and services. The resulting IoT provides unprecedented opportunities for users, manufacturers, and service providers in a wide variety of sectors. Areas that will benefit from IoT data collection, analysis, and automation capabilities include health and fitness, healthcare, home monitoring and automation, energy savings and smart grid, farming, transportation, environmental monitoring, inventory and product management, security, surveillance, education, and many others.

Technology development is occurring in many areas. Not surprisingly, wireless networking research is being conducted and actually has been conducted for quite a while now, but under previous titles such as mobile computing, pervasive computing, wireless sensor networks, and cyber-physical systems.

Many proposals and products have been developed for low-power protocols, security and privacy, addressing, low-cost radios, energy-efficient schemes for long battery life, and reliability for networks of unreliable and intermittently sleeping nodes. These wireless developments are crucial for the growth of IoT. In addition, areas of development have also involved giving IoT devices social networking capabilities, taking advantage of machine-to-machine communications, storing and processing large amounts of real-time data, and application programming to provide end users with intelligent and useful interfaces to these devices and data.

Many have provided a vision for the IoT. Stankovic^[2] suggests personal benefits such as digitizing daily life activities; patches of bionic skin to communicate with surrounding smart spaces for improved comfort, health, and safety; and smart watches and body nodes that optimize access to city services. Citywide benefits could include efficient, delay-free transportation with no traffic lights and 3-D transportation vehicles. Smart buildings could not only control energy and security, but also support health and wellness activities. In the same ways people have been provided new ways of accessing the world through smartphones, the IoT will create a new paradigm in the ways we have continuous access to needed information and services.

Cisco estimates that over the next decade the value at stake (net profit) for the IoT economy is \$14.4 trillion^[3]. The company's research indicates that five main drivers of this value are at stake:

- *Asset use* (\$2.5 trillion): IoT reduces selling, general, and administrative expenses and cost of goods sold by improving business-process execution and capital efficiency.
- *Employee productivity* (\$2.5 trillion): IoT creates labor efficiencies that result in fewer or more productive man-hours.
- *Supply chain and logistics* (\$2.7 trillion): IoT eliminates waste and improves process efficiencies.
- *Customer experience* (\$3.7 trillion): IoT increases customer lifetime value and grows market share by adding more customers.
- *Innovation, including reducing time to market* (\$3.0 trillion): IoT increases the return on R&D investments, reduces time to market, and creates additional revenue streams from new business models and opportunities.

Similarly, a 2015 report from McKinsey Global Institute^[4] estimates that the IoT has a total potential economic impact of \$3.9 trillion to \$11.1 trillion per year by 2025. On the top end, the value of this impact—including consumer surplus—would be equivalent to about 11 percent of the world economy in 2025.

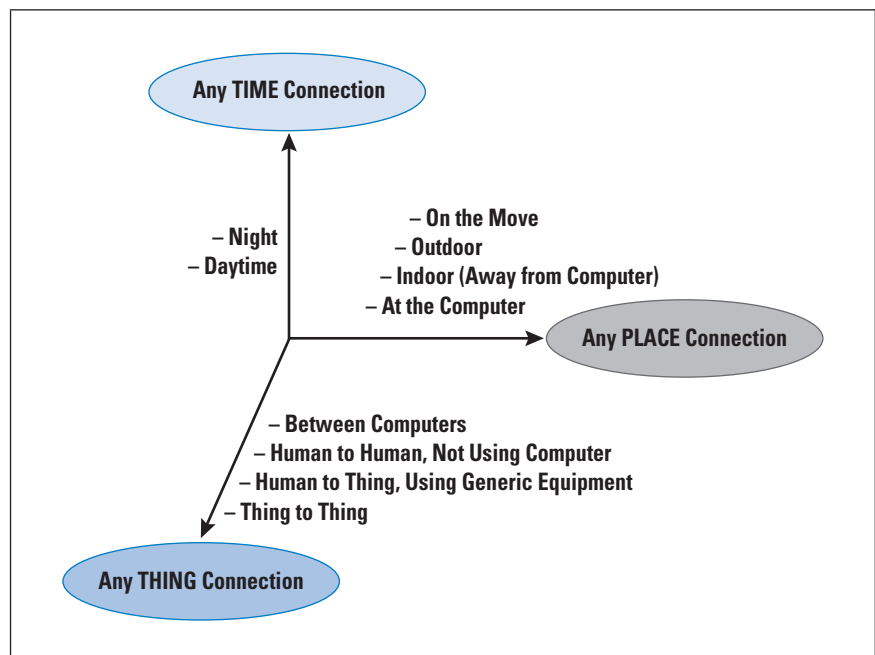
The Scope of the Internet of Things

The *Telecommunication Standardization Sector of the International Telecommunication Union* (ITU-T) has published Recommendation Y.2060, entitled “Overview of the Internet of Things.”^[5] The document provides the following definitions that suggest the scope of IoT:

- *Internet of Things* (IoT): A global infrastructure for the information society, enabling advanced services by interconnecting (physical and virtual) things based on existing and evolving interoperable information and communication technologies.
- *Thing*: With regard to the Internet of Things, this is an object of the physical world (*physical things*) or the information world (*virtual things*), which is capable of being identified and integrated into communication networks.
- *Device*: With regard to the Internet of Things, this is a piece of equipment with the mandatory capabilities of communication and the optional capabilities of sensing, actuation, data capture, data storage, and data processing.

Most of the literature views the IoT as involving intercommunicating smart objects. Recommendation Y.2060 extends this concept to include virtual things, a topic examined subsequently. Recommendation Y.2060 characterizes the IoT as adding the dimension “Any THING communication” to the information and communication technologies that already provide “any TIME” and “any PLACE” communication (Figure 1).

Figure 1: The New Dimension Introduced in the Internet of Things



In the book *Designing the Internet of Things*^[6], the elements of the IoT are condensed into a simple equation:

$$\text{Physical Objects} + \text{Controllers, Sensors, Actuators} + \text{Internet} = \text{IoT}$$

This equation neatly captures the essence of the Internet of Things. An instance of the IoT consists of a collection of physical objects, each of which:

- Contains a microcontroller that provides intelligence;
- Contains a sensor that measures some physical parameter and/or an actuator that acts on some physical parameter;
- Provides a means of communicating via the Internet or some other network.

One item not covered in the equation, and referred to in the Y.2060 definition, is a means of identification of an individual thing, usually referred to as a tag.

Note that although the phrase *the Internet of Things* is always used in the literature, a more accurate description would be an *Internet of Things*, or a *Network of Things*. A smart-home installation, for example, consists of numerous things in the home that are interconnected via Wi-Fi or Bluetooth with some central controller. In a factory or farm setting, a network of things may be enabling enterprise applications to interact with the environment and run applications to exploit the network of things. In these examples, remote access over the Internet is usually, but not invariably, available. Whether or not such Internet connection is available, the collection of smart objects at a site, plus any other local compute and storage devices, can be characterized as a network or an internet of things.

Table 1, on page 6, based on a graphic from Beecham Research^[7], gives an idea of the scope of IoT.

IoT Interoperability Standards

In the near term, disparate islands of solutions are likely to outpace deployment of interoperable standards-based solutions for IoT. This situation is common when any new technology or application area emerges. For example, Sutaria and Govindachari^[8] point out that two characteristics of networked IoT devices that pose challenges are the presence of low-power devices (which need to function for months or years without power recharge) and frequent data exchanges over lossy networks. Existing Internet standard protocols are suboptimal in this context. In a broader sense, there is a mismatch between the vast number of devices generating data at a rapid rate over a dispersed area and using a variety of network technologies and cloud-based systems that store vast amounts of data in a small number of locations with a relatively slow rate of data update. Integrating these two classes of systems to meet user needs requires specific protocol capabilities along the whole network/protocol architecture, from physical through middleware to application levels.

Table 1: The Internet of Things

Service Sectors	Application Groups	Locations	Example Devices
IT and Networks	Public	Services, e-commerce, data centers, mobile carriers, fixed carriers, ISPs	Servers, storage, PCs, routers, switches, PBXs
	Enterprise	IT/data center, office, private nets	
Security/Public Safety	Surveillance Equipment, Tracking	Radar/satellite, military security, unmanned, weapons, vehicles, ships, aircraft, gear	Tanks, fighter jets, battlefield comms, jeeps
	Public Infrastructure	Human, animal, postal, food/health, packaging, baggage, water treatment, building environmental, general environmental	Cars, breakdown-lane worker, homeland security, fire, environmental monitor
	Emergency Services	Equipment and personnel, police, fire, regulatory	Ambulances, public security vehicles
Retail	Specialty	Fuel stations, gaming, bowling, cinema, discos, special events	POS terminals, tags, cash registers, vending machines, signs
	Hospitality	Hotels, restaurants, bars, cafes, clubs	
	Stores	Supermarkets, shopping centers, single sites, distribution centers	
Transportation	Non-vehicular	Air, rail, marine	Vehicles, lights, ships, planes, signage, tolls
	Vehicles	Consumer, commercial, construction, off-road	
	Transportation Systems	Tolls, traffic management, navigation	
Industrial	Distribution	Pipelines, materials handling, conveyance	Pumps, valves, vats, conveyers, pipelines, motors, drives, converting, fabrication, assembly/packing, vessels, tanks
	Converting, Discrete	Metals, paper, rubber, plastic, metalworking, electronics assembly, test	
	Fluid/Processes	Petro-chemical, hydrocarbon, food, beverage	
	Resource Automation	Mining, irrigation, agricultural, woodland	
Healthcare and Life Science	Care	Hospital, ER, mobile PoC, clinic, labs, doctor office	MRIs, PDAs, implants, surgical equipment, pumps, monitors, telemedicine
	In-vivo, Home	Implants, home monitoring systems	
	Research	Drug discovery, diagnostics, labs	
Consumer and Home	Infrastructure	Wiring, network access, energy management	Digital camera, power systems, dishwashers, eReaders, desktop computers, washer/dryer, meters, lights, TVs, MP3, games console, lighting, alarms
	Awareness and Safety	Security/alert, fire safety, environmental safety, elderly, children, power protection	
	Convenience and Entertainment	HVAC/climate, lighting, appliance, entertainment	
Energy	Supply/Demand	Power generation, transportation and distribution, low voltage, power quality, energy management	Turbines, windmills, UPS, batteries, generators, meters, drills, fuel cells
	Alternative	Solar, wind, co-generation, electro-chemical	
	Oil/Gas	Rigs, derricks, well heads, pumps, pipelines	
Buildings	Commercial, Institutional	Office, education, retail, hospitality, healthcare, airports, stadiums	HVAC, transport, fire and safety, lighting, security, access
	Industrial	Process, clean room, campus	

To address these issues, several industry bodies and standards forums are working on extending or adopting the Internet protocols to the IoT devices. To provide for a common frame of reference and categorize needed functions and their location in the protocol stack, several of these groups are also addressing the issue of a formal architecture for IoT. While existing standards and the Internet make IoT possible, a suite of widely expected new standards that adapt or augment existing ones for IoT is likely not possible in the near term. Like many other developments made possible by the Internet, IoT will evolve in the wild for a while and pass through Darwinistic processes, with sensible technologies and protocol mechanisms gradually becoming visible. In this article, we look at two efforts at developing overall frameworks that may be useful in this ongoing standardization process.

ITU-T IoT Reference Model

Given the complexity of an IoT, it is useful to have an architecture that specifies the main elements and their interrelationship. An IoT architecture can have the following benefits:

- It provides the IT or network manager with a useful checklist with which to evaluate the functionality and completeness of vendor offerings.
- It provides guidance to developers as to which functions are needed in an IoT and how these functions work together.
- It can serve as a framework for standardization, promoting interoperability and cost reduction.

This section presents an overview of the IoT architecture developed by ITU-T. The next section looks at one developed by *IoT World Forum*. The latter architecture, developed by an industry group, offers a useful alternative framework for understanding the scope and functionality of IoT.

The ITU-T IoT Reference Model is defined in Recommendation Y.2060^[5]. Unlike most of the other IoT reference models and architectural models in the literature, the ITU-T model goes into detail about the actual physical components of the IoT ecosystem. This treatment is useful because it makes visible the elements in the IoT ecosystem that must be interconnected, integrated, managed, and made available to applications. This detailed specification of the ecosystem drives the requirements for the IoT capability.

An important insight the model provides is that the IoT is in fact not a network of physical things. Rather, it is a network of *devices* that interact with physical things, together with application platforms—such as computers, tablets, and smartphones—that interact with these devices. Thus, we begin our overview of the ITU-T model with a discussion of devices.

Terminology

The following is a list of definitions of key terms used in Recommendation Y.2060:

Communication Network: An infrastructure network that connects devices and applications, such as an IP-based network or internet.

Thing: An object of the physical world (*physical things*) or the information world (*virtual things*) that is capable of being identified and integrated into communication networks.

Device: A piece of equipment with the mandatory capability of communication and the optional capabilities of sensing, actuation, data capture, data storage, and data processing.

Data-carrying Device: A device attached to a physical thing to indirectly connect the physical thing with the communication networks. Active RFID tags are examples.

Data-capturing Device: A reader/writer device with the capability to interact with physical things. The interaction can happen indirectly via data-carrying devices, or directly via data carriers attached to the physical things.

Data Carrier: A battery-free data-carrying object attached to a physical thing that can provide information to a suitable data-capturing device. This category includes bar codes and *Quick Response* (QR) codes attached to physical things.

Sensing Device: A device that detects or measures information related to the surrounding environment and converts it into digital electronic signals.

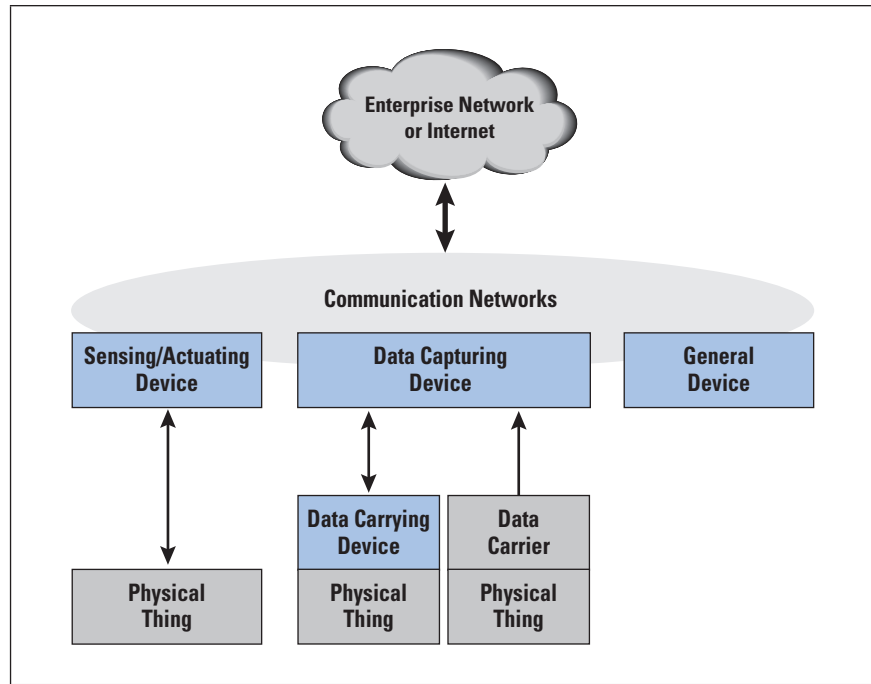
Actuating Device: A device that converts digital electronic signals from the information networks into operations.

General Device: A general device that has embedded processing and communication capabilities and may communicate with the communication networks via wired or wireless technologies. General devices include equipment and appliances for different IoT application domains, such as industrial machines, home electrical appliances, and smartphones.

Gateway: A unit in the IoT that interconnects the devices with the communication networks. It performs the necessary translation between the protocols used in the communication networks and those used by devices.

The unique aspect of an IoT, compared to other network systems, of course, is the presence of numerous physical things and devices other than computing or data processing devices. Figure 2, adapted from one in Recommendation Y.2060, shows the types of devices in the ITU-T model. The model views an IoT as functioning as a network of devices that are tightly coupled with things. Sensors and actuators interact with physical things in the environment. Data-capturing devices read data from and/or write data to physical things via interaction with a data-carrying device or a data carrier attached or associated in some way with a physical object.

Figure 2: Types of Devices and Their Relationship with Physical Things



The model makes a distinction between data-carrying devices and data carriers. A data-carrying device is a device in the Recommendation Y.2060 sense. A device at minimum is capable of communication and may include other electronic capabilities. An example of a data-carrying device is an RFID tag. By contrast, a data carrier is an element attached to a physical thing for the purpose of identification or providing some other sort of information.

Y.2060 notes that technologies used for interaction between data-capturing devices and data-carrying devices or data carriers include radio frequency, infrared, optical, and galvanic driving. Examples of each include:

- *Radio Frequency:* A *Radio-Frequency Identification* (RFID) tag is an example.
- *Infrared:* Infrared badges are used in military, hospital, and other settings where the location and movement of personnel need to be tracked. Examples include infrared reflective patches used by the military and battery-operated badges that emit identifying information. The latter can include a button that must be pressed so that the badge can be used as a means of passing through a portal, and a badge that automatically repeats the signal as a means of tracking personnel. Remote-control devices used in the home or other settings to control electronic devices can also easily be incorporated into an IoT.
- *Optical:* Bar codes and QR codes are examples of identifying data carriers that can be read optically.

- *Galvanic Driving*: An example is implanted medical devices that use the conductive properties of the body^[9]. In implant-to-surface communication, galvanic coupling sends signals from an implanted device to electrodes on the skin. This scheme uses very little power and reduces the size and complexity of the implanted device.

The final type of device shown in Figure 2 is the general device. These devices have processing and communications capabilities that can be incorporated into an IoT. A good example is smart-home technology that can integrate virtually every device in the home into a network for central or remote control.

Figure 3 provides an overview of the elements of interest in an IoT. The various ways that physical devices can be connected are shown on the left side of the figure. It is assumed that one or multiple networks support communication among the devices.

Figure 3: Technical Overview of the IoT (Recommendation Y.2060)

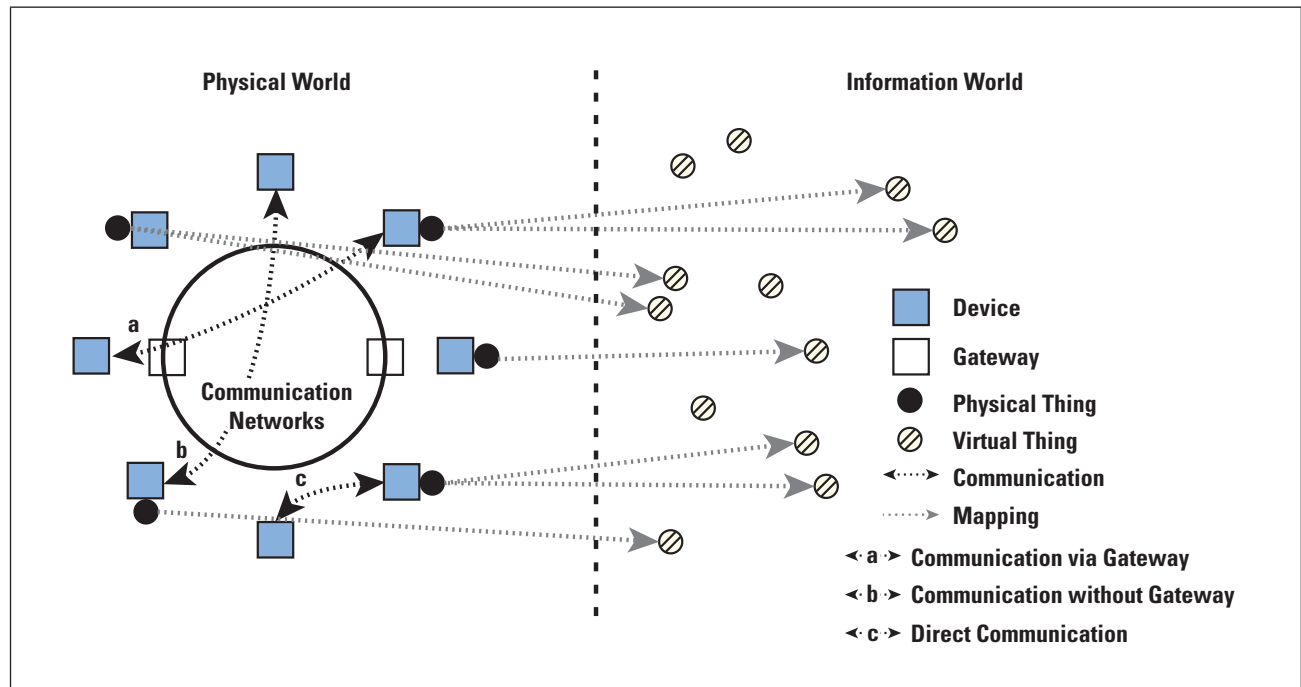


Figure 3 introduces one additional IoT-related device: the *gateway*. At minimum, a gateway functions as a protocol translator. Gateways address one of the greatest challenges in designing an IoT, which is connectivity, both among devices and between devices and the Internet or enterprise network. Smart devices support a wide variety of wireless and wired transmission technologies and networking protocols. Further, these devices typically have limited processing capability.

Recommendation Y.2067^[10] lays out the requirements for IoT gateways, which generally fall into three categories:

- The gateway supports a variety of device access technologies, enabling devices to communicate with each other and across an Internet or enterprise network with IoT applications. The access schemes could include, for example, ZigBee, Bluetooth, and Wi-Fi.
- The gateway supports the necessary networking technologies for both local and wide-area networking. These technologies could include Ethernet and Wi-Fi on the premises, and cellular, Ethernet, DSL, and cable access to the Internet and wide-area enterprise networks.
- The gateway supports interaction with application, network management, and security functions.

The first two requirements involve protocol translation between different network technologies and protocol suites. The third requirement is generally referred to as an *IoT agent* function. In essence, the IoT agent provides higher-level functionality on behalf of IoT devices, such as organizing and/or summarizing data from multiple devices to pass on to IoT applications, implementing security protocols and functions, and interacting with network management systems.

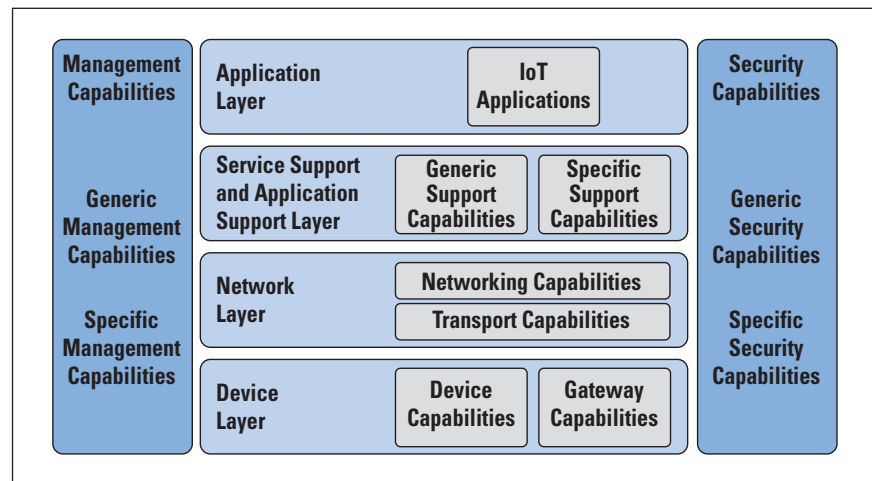
At this point, it should be noted that the term *Communication Network* is not directly defined in the Y.206x series of IoT standards. The communication network or networks support(s) communication among devices and may directly support application platforms. This may be the extent of a small IoT, such as a home network of smart devices. More generally, the device network(s) connect to enterprise networks or the Internet for communication with systems that host apps and servers that host databases related to the IoT.

We can now return to the left side of Figure 3, which illustrates the communication possibilities among devices. The first possibility is for communication between devices via the gateway. For example, a sensor or actuator with Bluetooth capability could communicate with a data-capturing device or general device that uses Wi-Fi by means of the gateway. The second possibility is communication across the communication network without a gateway. For example, all of the devices in a smart-home network may use Bluetooth and could be managed from a Bluetooth-enabled computer, tablet, or smartphone. The third possibility is devices that communicate directly with each other through a separate local network and then (not shown in the figure) communicate through the communication network via a local network gateway. An example of this third possibility follows: Numerous low-power sensor devices could be deployed in an extended area, such as farmland or a factory. These devices could communicate with one another to pass data on toward a device connected to a gateway to the communication network.

The right side of Figure 3 emphasizes that each physical thing in an IoT may be represented in the information world by one or more virtual things, but a virtual thing can also exist without any associated physical thing. Physical things are mapped to virtual things stored in databases and other data structures. Applications process and deal with virtual things.

Figure 4 depicts the ITU-T IoT Reference Model, which consists of four layers as well as management capabilities and security capabilities that apply across layers. We have so far been considering the device layer. In terms of communications functionality, the device layer includes, roughly, the OSI physical and data link layers. We now look at the other layers.

Figure 4: ITU-T Recommendation Y.2060 IoT Reference Model



The *Network Layer* performs two basic functions. Networking capabilities refer to the interconnection of devices and gateways. Transport capabilities refer to the transport of IoT service- and application-specific information as well as IoT-related control and management information. Roughly, these capabilities correspond to those of the OSI network and transport layers.

The *Service Support and Application Support Layer* provides capabilities that applications use. Many different applications can use generic support capabilities. Examples include common data processing and database management capabilities. Specific support capabilities are those that cater for the requirements of a specific subset of IoT applications.

The *Application Layer* consists of all the applications that interact with IoT devices.

The *Management Capabilities Layer* covers the traditional network-oriented management functions of fault, configuration, accounting, and performance management.

Recommendation Y.2060 lists the following as examples of generic management capabilities:

- *Device Management*: Examples include device discovery, authentication, remote device activation and de-activation, configuration, diagnostics, firmware and/or software updating, and device working-status management.
- *Local Network Topology Management*: An example is network configuration management.
- *Traffic and Congestion Management*: Examples include the detection of network overflow conditions and the implementation of resource reservation for time- and/or life-critical data flows.

Specific management capabilities are tailored to specific classes of applications. An example is smart-grid power-transmission-line monitoring.

The *Security Capabilities Layer* includes generic security capabilities that are independent of applications. Y.2060 lists the following as examples of generic security capabilities:

- *Application Layer*: authorization, authentication, and application data confidentiality and integrity protection, privacy protection, security audit, and anti-virus.
- *Network Layer*: authorization, authentication, user data, and signaling data confidentiality, and signaling integrity protection.
- *Device Layer*: authentication, authorization, device-integrity validation, access control, data confidentiality, and integrity protection.

Specific security capabilities relate to specific application requirements, such as mobile payment security requirements.

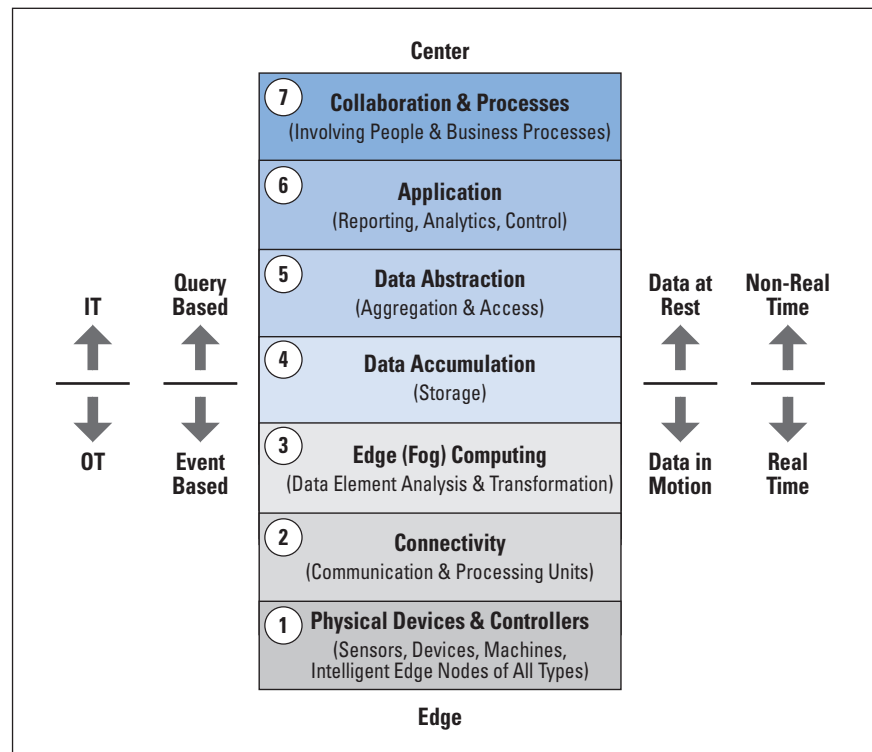
IoT World Forum Reference Model

The *IoT World Forum* (IWF) is an industry-sponsored annual event that brings together representatives of business, government, and academia to promote the market adoption of IoT. The IoT World Forum *Architecture Committee*, made up of industry leaders including IBM, Intel, and Cisco, released an IoT reference model in October 2014. This model serves as a common framework to help the industry accelerate IoT deployments. The reference model is intended to foster collaboration and encourage the development of replicable deployment models.

This reference model is a useful complement to the ITU-T reference model. The ITU-T documents focus on the device and gateway level with only a broad depiction of the upper layers. Indeed, Recommendation Y.2060 describes the application layer with a single sentence. The ITU-T Recommendation Y.206x series seems most concerned with defining a framework to support development of standards for interaction with IoT devices.

The IWF is concerned with the broader issue of developing the applications, middleware, and support functions for an enterprise-based IoT. Figure 5 depicts the seven-level model.

Figure 5: IoT World Forum Reference Model



The white paper on the IWF model issued by Cisco^[11] indicates that the model is designed to have the following characteristics:

- *Simplifies*: It helps break down complex systems so that each part is more understandable.
- *Clarifies*: It provides additional information to precisely identify levels of the IoT and to establish common terminology.
- *Identifies*: It identifies where specific types of processing are optimized across different parts of the system.
- *Standardizes*: It provides a first step in enabling vendors to create IoT products that work with each other.
- *Organizes*: It makes the IoT real and approachable, instead of simply conceptual.

Level 1 comprises physical devices and controllers that might control multiple devices. Level 1 of the IWF model corresponds approximately to the device level of the ITU-T model (Figure 4). As with the ITU-T model, the elements at this level are not physical things as such, but rather devices that interact with physical things, such as sensors and actuators. Among the capabilities that devices may have are analog-to-digital and digital-to-analog conversion, data generation, and the ability to be queried and/or controlled remotely.

From a logical point of view, this level enables communication between devices and between devices and the low-level processing that occurs at level 3. From a physical point of view, this level consists of networking devices such as routers, switches, gateways, and firewalls that are used to construct local and wide-area networks and provide Internet connectivity. This level enables devices to communicate with one another and to communicate, via the upper logical levels, with application platforms such as computers, remote-control devices, and smartphones.

Level 2 of the IWF model corresponds approximately to the network level of the ITU-T model. The main difference is that the IWF model includes gateways in level 2, whereas the ITU-T model puts the gateway at level 1. Because the gateway is a networking and connectivity device, its placement at level 2 seems to make more sense.

In many IoT deployments, massive amounts of data may be generated by a distributed network of sensors. For example, offshore oil fields and refineries can generate a terabyte of data per day. An airplane can create multiple terabytes of data per hour. Rather than store all of that data permanently (or at least for a long period) in central storage accessible to IoT applications, it is often desirable to do as much data processing close to the sensors as possible. Thus, the purpose of the edge computing level is to convert network data flows into information that is suitable for storage and higher-level processing. Processing elements at these levels may deal with high volumes of data and perform data-transformation operations, resulting in the storage of much lower volumes of data. The Cisco white paper on the IWF model^[11] lists the following examples of edge computing operations:

- *Evaluation*: Evaluating data for criteria as to whether it should be processed at a higher level.
- *Formatting*: Reformatting data for consistent higher-level processing.
- *Expanding/decoding*: Handling cryptic data with additional context (such as the origin).
- *Distillation/reduction*: Reducing and/or summarizing data to minimize the impact of data and traffic on the network and higher-level processing systems.
- *Assessment*: Determining whether data represents a threshold or alert; this process could include redirecting data to additional destinations.

Processing elements at this level corresponds to general devices in the ITU-T model (Figure 2). Generally, they are deployed physically near the edge of the IoT network; that is, near the sensors and other data-generating devices. Thus, some of the basic processing of large volumes of generated data is offloaded and outsourced from IoT application software located at the center.

Processing at the edge computing level is sometimes referred to as *Fog Computing*. Fog computing and fog services are expected to be a distinguishing characteristic of the IoT. Figure 6 illustrates the concept. Fog computing represents an opposite trend in modern networking from cloud computing. With cloud computing, massive, centralized storage and processing resources are made available to distributed customers over cloud networking facilities to a relatively small number of users. With fog computing, massive numbers of individual smart objects are interconnected with fog networking facilities that provide processing and storage resources close to the edge devices in an IoT. Fog computing addresses the challenges raised by the activity of thousands or millions of smart devices, including security, privacy, network-capacity constraints, and latency requirements. The term “Fog Computing” is inspired by the fact that fog tends to hover low to the ground, whereas clouds are high in the sky.

Figure 6: Fog Computing

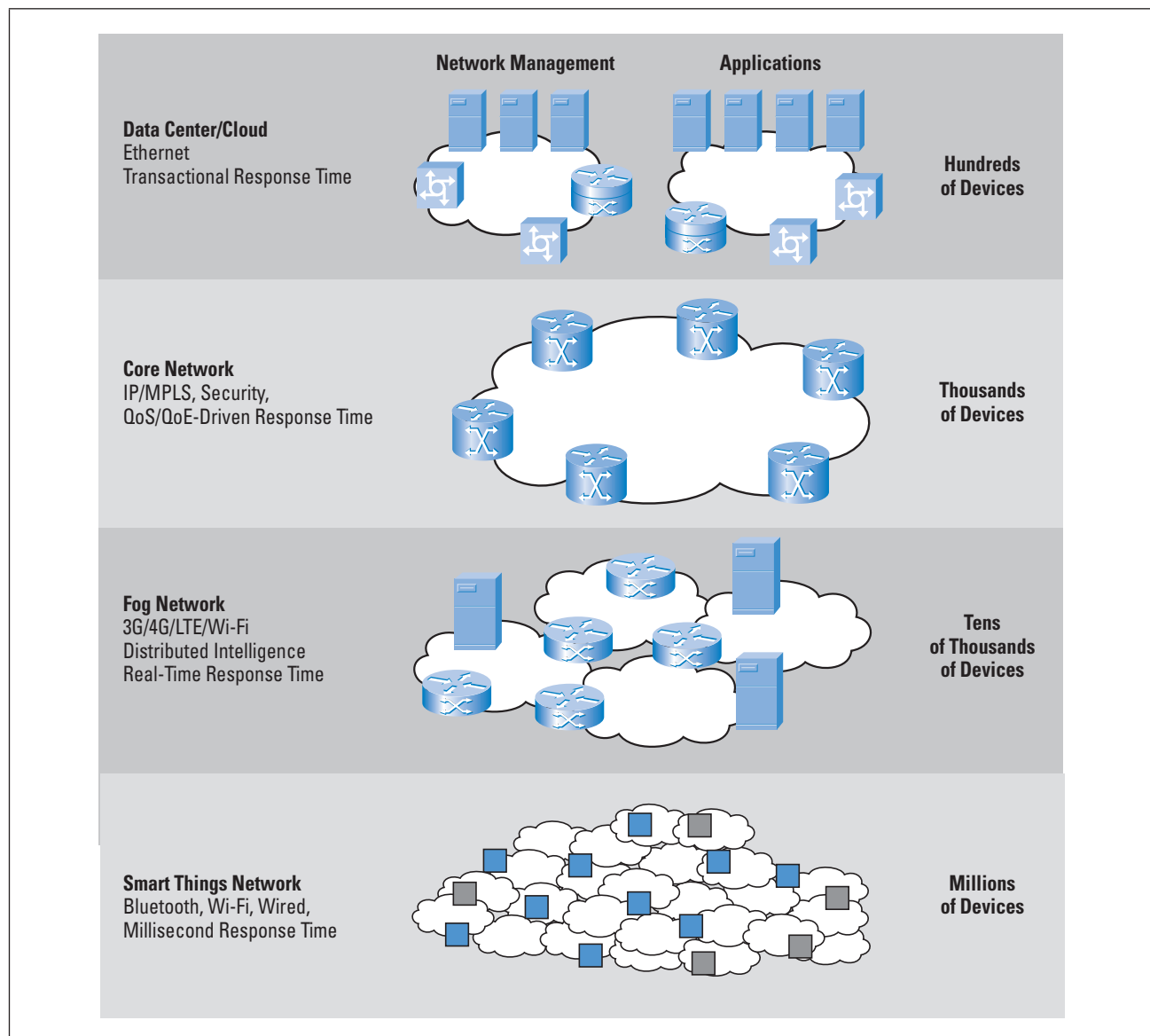


Table 2, based on one in [12], compares cloud and fog computing.

Table 2: Comparison of Cloud and Fog Features

	Cloud	Fog
Location of processing/ storage resources	Center	Edge
Latency	Low to high	Low
Access	Fixed or wireless	Mainly wireless
Support for mobility	Not applicable	Yes
Control	Centralized/hierarchical (full control)	Distributed/hierarchical (partial control)
Service access	Through core	At the edge/on handheld device
Availability	99.99%	Highly volatile/highly redundant
Number of users/devices	Tens/hundreds of millions	Tens of billions
Main content generator	Humans and devices	Devices/sensors
Content generation	Central location	Anywhere
Content consumption	End device	Anywhere
Software virtual infrastructure	Central enterprise servers	User devices

Level 4, the data accumulation level, is where data coming from the numerous devices, and filtered and processed by the edge computing level, is placed in storage that will be accessible by higher levels. This level marks a clear distinction in the design issues, requirements, and method of processing between lower-level (fog) computing and upper-level (typically cloud) computing.

Data moving through a network is referred to as *data in motion*. The rate and organization of the data in motion is determined by the devices generating the data. Data generation is event-driven, either periodically or by an event in the environment. To capture the data and deal with it in some fashion, it is necessary to respond in real time. By contrasts, most applications do not need to process data at network transfer speeds. As a practical matter, neither the cloud network nor the application platforms would be able to keep up with data volume generated by a huge number of IoT devices. Instead, applications deal with *data at rest*, which is data in some readily accessible storage facility. Applications can access the data as needed, on a non-real-time basis. Thus, the upper levels operate on a query or transaction basis, whereas the lower three levels operate on an event basis.

The following are listed as operations performed at the data-accumulation level in [13]:

- Converts data in motion to data at rest
- Converts format from network packets to database relational tables
- Achieves transition from event-based to query-based computing
- Dramatically reduces data through filtering and selective storing

Another way of viewing the data-accumulation level is that it marks the boundary between *Information Technology* (IT), which is the common term for the entire spectrum of technologies for information processing, including software, hardware, communications technologies and related services, and *Operational Technology* (OT), which refers to hardware and software that detects or causes a change through the direct monitoring and/or control of physical devices, processes, and events in the enterprise.

The data-accumulation level absorbs large quantities of data and places them in storage, with little or no tailoring to specific applications or groups of applications. Numerous different types of data in varying formats and from heterogeneous processors may be coming up from the edge computing level for storage. The data-abstraction level can aggregate and format this data in ways that make access by applications more manageable and efficient. Tasks involved could include:

- Combining data from multiple sources, including reconciling multiple data formats.
- Performing necessary conversions to provide consistent semantics of data across sources.
- Placing formatted data in an appropriate database; for example, high-volume repetitive data may go into a big data system such as Hadoop. Event data would be steered to a relational database management system, which provides faster query times and an appropriate interface for this type of data.
- Alerting higher-level applications that data is complete or has accumulated to a defined threshold.
- Consolidating data into one place (with ETL (*extract, transform, load*), ELT (*extract, load, transform*), or data replication) or providing access to multiple data stores through data virtualization.
- Protecting data with appropriate authentication and authorization.
- Normalizing or denormalizing and indexing data to provide fast application access.

The application level contains any type of application that uses IoT input or controls IoT devices. Generally, applications interact with level 5 and the data at rest, and so do not have to operate at network speeds. Provision should be available for streamlined operation that allows applications to bypass intermediate layers and interact directly with Layer 3 or even Layer 2. The IWF model does not strictly define applications, considering it beyond the scope of IWT model discussion.

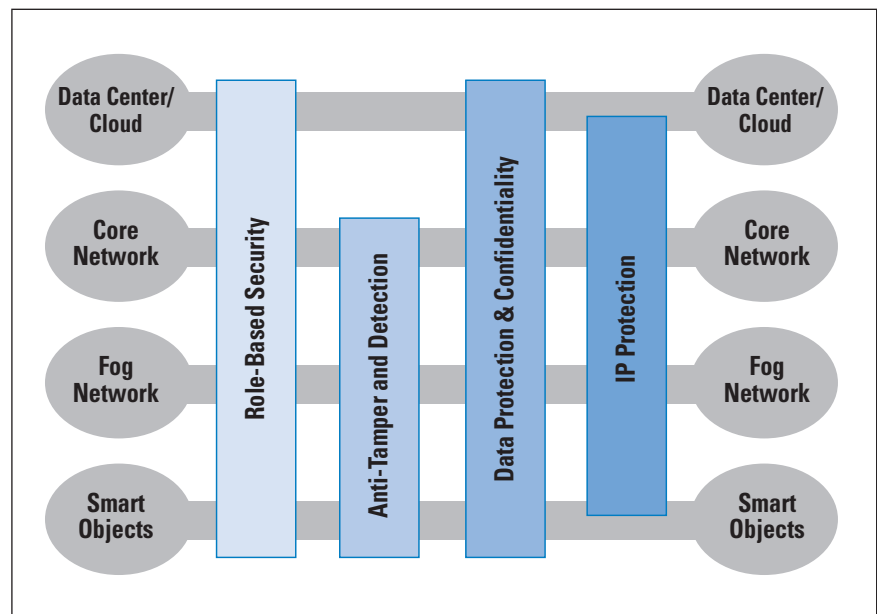
The collaboration and processes level recognizes that people must be able to communicate and collaborate to make an IoT useful. This level may involve multiple applications and exchange of data and control information across the Internet or an enterprise network.

The IWF views the IoT reference model as an industry-accepted framework aimed at standardizing the concepts and terminology associated with IoT. More importantly, the IWF model sets out the functionalities required and concerns that must be addressed before the industry can realize the value of the IoT. This model is useful both for suppliers who develop functional elements within the model and customers for developing their requirements and evaluating vendor offerings.

An IoT Security Framework

Cisco Systems, which has played a lead role in the development of the IoT World Forum Reference Model, has developed a framework for IoT security^[13] that serves as a useful complement to the World Forum IoT Reference Model. Figure 7 illustrates the security environment related to the logical structure of an IoT.

Figure 7: IoT Security Environment



The Cisco IoT model is a simplified version of the World Forum IoT Reference Model. It consists of the following levels:

- *Smart Objects/Embedded Systems*: This level consists of sensors, actuators, and other embedded systems at the edge of the network. This part of an IoT is the most vulnerable part. The devices may not be in a physically secure environment and may need to function for years. Availability is certainly of concern. Also network managers need to be concerned about the authenticity and integrity of the data generated by sensors and about protecting actuators and other smart devices from unauthorized use. Privacy and protection from eavesdropping may also be requirements.
- *Fog/Edge Network*: This level is concerned with the wired and wireless interconnection of IoT devices. In addition, a certain amount of data processing and consolidation may be done at this level. A key concern is the wide variety of network technologies and protocols that the various IoT devices use and the need to develop and enforce a uniform security policy.
- *Core Network*: The core network level provides data paths between network center platforms and the IoT devices. The security issues here are those confronted in traditional core networks. However, the vast number of endpoints to interact with and manage creates a substantial security burden.
- *Data Center/Cloud*: This level contains the application, data storage, and network management platforms. IoT does not introduce any new security issues at this level, other than the necessity of dealing with huge numbers of individual endpoints.

Within this four-level architecture, the Cisco model defines four general security capabilities that span multiple levels:

- *Role-Based Security*: *Role-Based Access Control* (RBAC) systems assign access rights to roles instead of individual users. In turn, users are assigned to different roles, either statically or dynamically, according to their responsibilities. RBAC enjoys widespread commercial use in cloud and enterprise systems and is a well-understood tool that can be used to manage access to IoT devices and the data they generate.
- *Anti-tamper and Detection*: This function is particularly important at the device and fog network levels but also extends to the core network level. All of these levels may involve components that are physically outside the area of the enterprise that is protected by physical security measures.
- *Data Protection and Confidentiality*: These functions extend to all levels of the architecture.
- *Internet Protocol Protection*: Protection of data in motion from eavesdropping and snooping is essential between all levels.

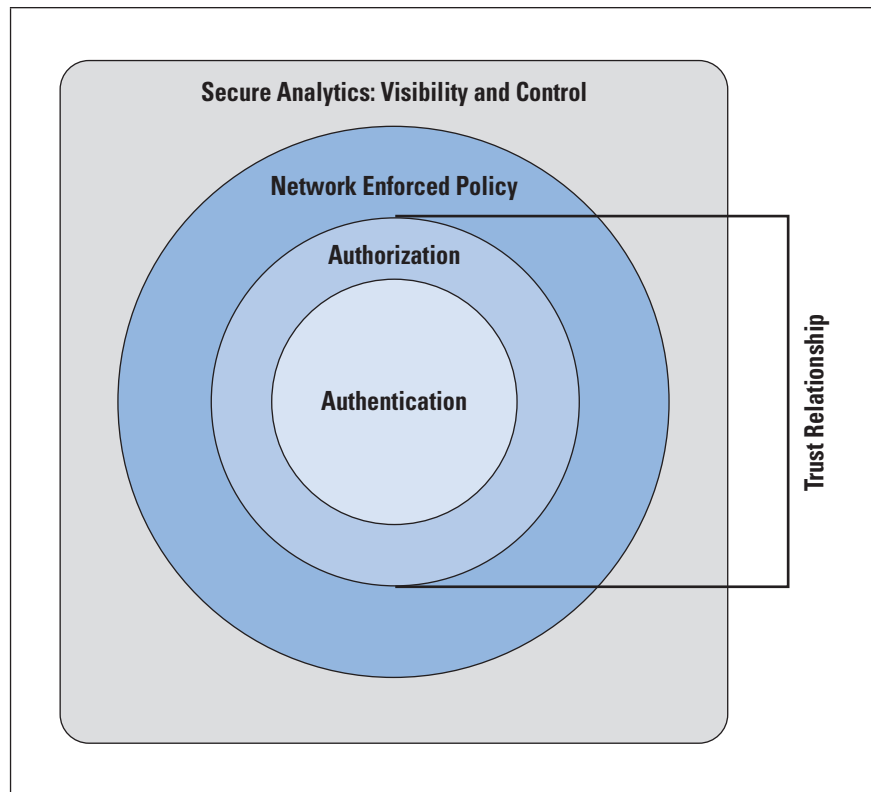
Figure 7, on page 19, maps specific security functional areas across the four layers of the IoT model. The Cisco white paper^[13] also proposes a secure IoT framework that defines the components of a security facility for an IoT that encompasses all the levels, as shown in Figure 8 on page 22. The four components follow:

- *Authentication*: This component encompasses the elements that initiate the determination of access by first identifying the IoT devices. In contrast to typical enterprise network devices, which may be identified by a human credential (for example, username and password or token), the IoT endpoints must be fingerprinted by means that do not require human interaction. Such identifiers include RFID, x.509 certificates, or the MAC address of the endpoint.
- *Authorization*: Authorization controls access of a device throughout the network fabric. This element encompasses access control. Together with the authentication layer, it establishes the necessary parameters to enable the exchange of information between devices and between devices and application platforms and enables IoT-related services to be performed.
- *Network Enforced Policy*: This component encompasses all elements that route and transport endpoint traffic securely over the infrastructure, whether control, management, or actual data traffic.
- *Secure Analytics, including Visibility and Control*: This component includes all the functions required for central management of IoT devices. It involves, firstly, visibility of IoT devices, meaning simply that central management services are securely aware of the distributed IoT device collection, including identity and attributes of each device. Building on this visibility is the ability to exert control, including configuration, patch updates, and threat countermeasures.

An important concept related to this framework is that of *trust relationship*. In this context, trust relationship refers to the ability of the two partners to an exchange to have confidence in the identity and access rights of the other. The authentication component of the trust framework provides a basic level of trust, which is expanded with the authorization component.

The Cisco white paper^[13] gives the example that a car may establish a trust relationship with another car from the same vendor. That trust relationship, however, may allow cars to exchange only their safety capabilities. When a trusted relationship is established between the same car and its dealer's network, the car may be allowed to share additional information such as its odometer reading and last maintenance record.

Figure 8: Secure IoT Framework



Conclusions

According to the McKinsey report cited earlier^[4], approximately 40 percent of the total economic value of the IoT is driven by the ability of all the physical devices to talk to each other via computers, that is, interoperability. If interoperability is limited, the IoT might be only a \$7 trillion opportunity, whereas widespread interoperability could achieve an IoT value to the global economy of over \$11 trillion by 2025. On average, 40 percent of the total value that can be unlocked requires different IoT systems to work together. Table 3, based on the McKinsey report, estimates the percent of economic value that requires interoperability between IoT systems for different sectors.

To achieve the type of interoperability needed to realize these benefits, standards need to be developed at all levels of IoT functionality, from the device layer to the application layer (Figure 4). While such standardization is still in its infancy, the architectural models described here provide a useful framework for future efforts.

Table 3: Value Added by IoT Interoperability

Setting	Value Potential Requiring Interoperability (\$ Trillion)	% of Total Value	Examples of How Interoperability Enhances Value
Factories	1.3	36	Data from different types of equipment used to improve line efficiency
Cities	0.7	43	Video, cellphone data, and vehicle sensors to monitor traffic and optimize flow
Retail	0.7	57	Payment and item-detection system linked for automatic checkout
Work sites	0.5	56	Linking worker and machinery location data to avoid accidents and exposure to chemicals
Vehicles	0.4	44	Equipment usage data for insurance underwriting, maintenance, and presales analytics
Agriculture	0.3	20	Multiple sensor systems used to improve farm management
Outside	0.3	29	Connected navigation between vehicles and between vehicles and GPS/traffic control
Home	0.1	17	Linking chore automation to security and energy system to time usage
Offices	>0.1	30	Data from different building systems and other buildings used to improve security

References

- [1] Lake, D., Rayes, A., and Morrow, M., "The Internet of Things," *The Internet Protocol Journal*, Volume 15, No. 3, September 2012.
- [2] Stankovic, J., "Research Directions for the Internet of Things," *Internet of Things Journal*, Volume 1, No. 1, 2014.
- [3] Cisco Systems, "Embracing the Internet of Everything to Capture Your Share of \$14.4 Trillion," White Paper, 2013.
http://www.cisco.com/web/about/ac79/docs/innov/IoT_Economy_Insights.pdf
- [4] McKinsey Global Institute, "The Internet of Things: Mapping the Value Beyond the Hype," June 2015.
http://www.mckinsey.com/insights/business_technology/the_internet_of_things_the_value_of_digitizing_the_physical_world
- [5] ITU-T, "Overview of the Internet of Things," Recommendation Y.2060, June 2012.
- [6] McEwen, A., and Cassimally, H., *Designing the Internet of Things*, ISBN-13: 978-1118430620, Wiley, 2013.

- [7] Beecham Research. “M2M Sector Map,” September 2011.
<http://www.beechamresearch.com/downloads.aspx?page=2>
- [8] Sutaria, R., and Raghunath, G., “Making sense of interoperability: Protocols and Standardization initiatives in IoT,” *International Conference on Recent Trends in Communication and Computer Networks – ComNet 2013*, 2013.
- [9] Ferguson, J., and Redish, A., “Wireless Communication with Implanted Medical Devices Using the Conductive Properties of the Body,” *Expert Review of Medical Devices*, Volume 6, No. 4, 2011, <http://www.expert-reviews.com>.
- [10] ITU-T, “Common Requirements and Capabilities of a Gateway for Internet of Things Applications,” Recommendation Y.2067, June 2014.
- [11] Cisco Systems, “The Internet of Things Reference Model,” White Paper, 2014. <http://www.iotwf.com/>
- [12] Vaquero, L., and Rodero-Merino, L., “Finding Your Way in the Fog: Towards a Comprehensive Definition of Fog Computing,” *ACM SIGCOMM Computer Communication Review*, October 2014.
- [13] Frahim, J., et al., “Securing the Internet of Things: A Proposed Framework,” Cisco White Paper, March 2015.
- [14] Douglas Comer, “The ZigBee IP Protocol Stack,” *The Internet Protocol Journal*, Volume 17, No. 2, December 2014.

WILLIAM STALLINGS is an independent consultant and author of numerous books on security, computer networking, and computer architecture. His latest book is *Foundations of Modern Networking: SDN, NFV, QoE, IoT, and Cloud* (Pearson, 2016). He maintains a resource site for computer science students and professionals at ComputerScienceStudent.com and is on the editorial board of *Cryptologia*. He has a Ph.D. in computer science from M.I.T. He can be reached at ws@shore.net

The RFC Series – Beyond ASCII

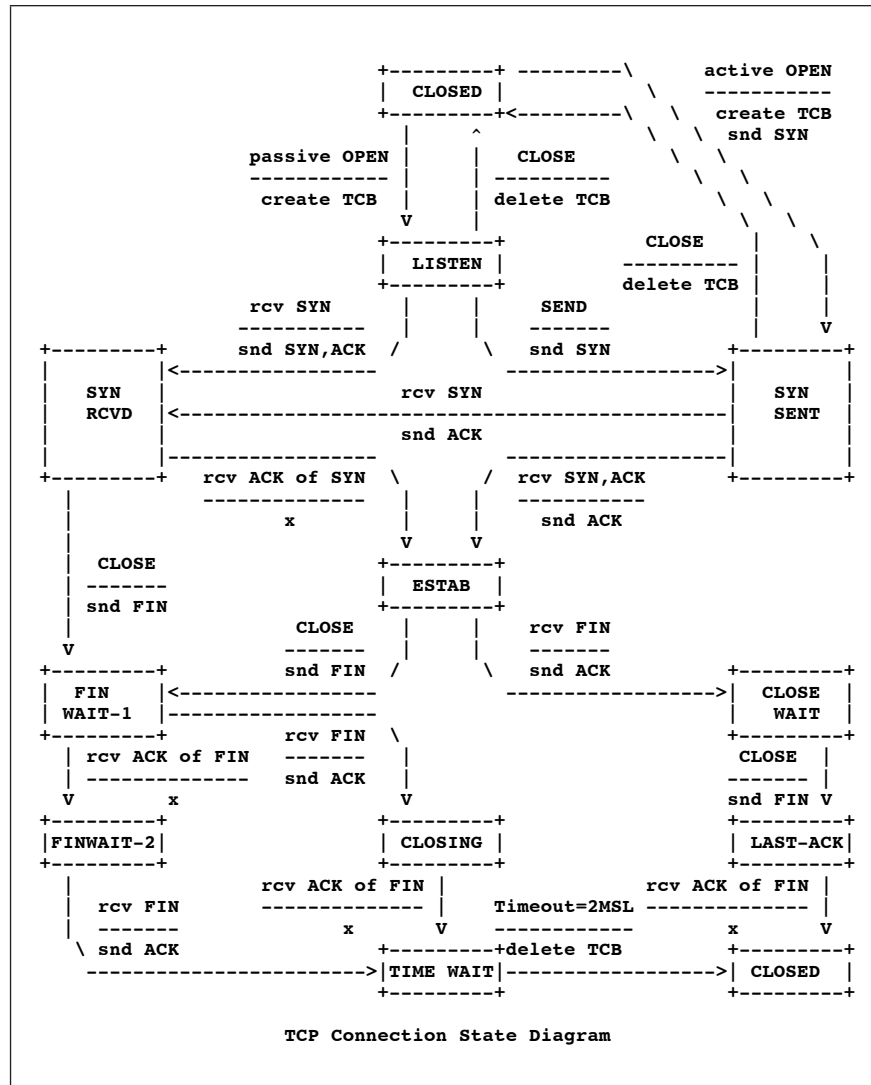
by Heather Flanagan, RFC Series Editor

The *Request for Comments* (RFC) Series began on April 7, 1969, with RFC 1^[1]. Since then, over 7,500 RFCs have been published. Most (though not all) of those RFCs have been in a plain text, *American Standard Code for Information Interchange* (ASCII)-only format. The RFC Series has always been a document series that focuses on archiving material and making sure it remains available far into the future.

The choice of plain text and ASCII-only was made for a variety of reasons. The format has stood the test of time in being readable as originally intended regardless of hardware and software changes. These documents can be read on just about any device, they take up very little bandwidth and disk space, and they provide a common experience for readers. However, while plain text and ASCII were undoubtedly the correct choice for the first few decades of the series, this format has reached the limits of its usefulness.

The Internet has evolved to the point that the restrictions imposed by the old format prevent reasonable information sharing. Networks are becoming more complicated than can be reasonably drawn via “ASCII art” (Figure 1). Internationalization brings in more characters than can be covered sensibly in ASCII. And, while people can read RFCs on just about any device, the experience is not always positive, as the format of the RFC cannot flow to match the different screen sizes. The plain text, ASCII-only format was the right choice, but it no longer meets the changing requirements of the Internet community.

In 2012, I started as RFC Series Editor. My intention was to quietly learn more about the community and its expectations for the series and the RFC Editor. That goal lasted about 3 months, at which point it became obvious that there was a decades-old demand for a change in the RFC format. The only thing that had prevented any changes in that time was the inability for the community to come to any kind of rough consensus on what the change should be. Should the format be *HyperText Markup Language* (HTML)? *Portable Document Format* (PDF)? Plain text but supporting a *UCS Transformation Format 8* (UTF-8) encoding? What about *Lamport TeX* (LaTeX)? Perhaps *Extensible Markup Language* (XML)? My job was to take the input and try to reach consensus within the community, and barring that, to make the best decision possible so we could stop having the debates.

Figure 1: Example of ASCII Art
from RFC 793^[6]

The first formal output of the discussion after I started to collect input was RFC 6949^[2]. That document captured the requirements to date, and led to the decisions described in an email to the community^[3]. Those two items in turn gave a design team, formed in July 2013 shortly after RFC 6949 was published, a baseline for determining the details of a new format.

The new RFC format will include a base, unchanging XML format using an enhanced *xml2rfc* vocabulary. From that file, the RFC Editor will render HTML, PDF/A-3, and plain text. *Scalable Vector Graphics* (SVG) diagrams in black and white will be supported, and non-ASCII characters used in a carefully prescribed manner. A list of commonly asked questions regarding the format, including a reading list that describes the details of each rendered format, profile, and general guidance, is available on the RFC Format FAQ^[4].

As an important tangent to the format work, a separate project is developing that more carefully looks at the future of the digital archive process for the series. The field of digital archiving has evolved significantly over the last decade, and the RFC Editor is poised to partner with official digital archives around the world to properly store and maintain copies of all RFCs and approved Internet Drafts as per best practice in that field. More information on the considerations that are involved with properly archiving RFCs is available in [5]. This draft is expected to be updated and moved towards publication as an RFC in 2016.

The RFC Format project has reached an important milestone, where the requirements drafts are starting their path towards publication: first, review by the *RFC Series Oversight Committee* (RSOC), then review by the *Internet Architecture Board* (IAB), and finally review by the community. Upon approval for publication, expected in early 2016, the associated Requests for Proposals will go out, and work on the necessary code base to implement the format changes can start. Coding the format tools and testing the output will be a major effort in 2016. By 2017, we will see a host of changes that make the RFC Series an easily read and still easily archived document series.

References

- [1] Steve Crocker, “Host Software,” RFC 1, April 1969.
- [2] Nevil Brownlee and Heather Flanagan, “RFC Series Format Requirements and Future Development,” RFC 6949, May 2013.
- [3] Message to RFC-Interest mailing List, “Subject: Direction of the RFC Format Development effort,”
<http://www.rfc-editor.org/pipermail/rfc-interest/2013-May/005584.html>
- [4] RFC Format Change FAQ:
<http://www.rfc-editor.org/rse/format-faq/>
- [5] Heather Flanagan, “Digital Preservation Considerations for the RFC Series,” January 2015, Internet Draft, work in progress, **draft-flanagan-rfc-preservation-03**.
- [6] Jon Postel, “Transmission Control Protocol,” RFC 793, September 1981.

HEATHER FLANAGAN is the current RFC Series Editor, a role she has filled since 2012. As an independent contractor, she also works with a variety of other organizations and efforts that build open standards and work towards developing a stronger Internet. Her portfolio includes executive oversight, project management for open source projects, copyediting, and small team meeting facilitation.
E-mail: rse@rfc-editor.org

Rob Blokzijl Obituary

Dr. Robert Blokzijl, RIPE Chair Emeritus and founding member of the *Réseaux IP Européens* (RIPE) community, died aged 72 on December 1, 2015.



Dr. Robert Blokzijl

*Photo by Olaf Kolkman
[Creative Commons NC-BY]*

The Internet community has suffered a sad loss as the man who led RIPE during its first 25 years of bottom-up, consensus-driven collaboration and decision-making died at his home in the Netherlands.

As one of the founders of RIPE in 1989 and the Chair of the RIPE community for 25 years since then, Rob Blokzijl personified all the attributes that have seen the community grow into such a positive force for bringing together those who care about the development of the Internet.

Rob's roots were in the high-energy physics community, and earlier in his career he worked at the *High Energy Physics Institute* (Nikhef) in Amsterdam and later at the *European Organization for Nuclear Research* (CERN) in Geneva. He helped to build the computer networks that were essential for that branch of science. His work in this area would inform much of his contribution to the burgeoning IP networking community in Europe in the early 1980s.

Over the past 30 years, Rob has established a global reputation as a leader and a pioneer, respected for his work with organisations such as RIPE, the *RIPE Network Coordination Centre* (RIPE NCC), the *Amsterdam Internet Exchange* (AMS-IX), the *Internet Corporation For Assigned Names and Numbers* (ICANN), Nominet, and the *North Atlantic Treaty Organization* (NATO).

In 1989, Rob was co-author of the *RIPE Terms of Reference*, which stated, "The object of RIPE is to ensure the necessary administrative and technical coordination to allow the operation and expansion of a pan-European IP network." In his role as the Chair of RIPE, his vision, expertise, and effort were essential for the tremendous growth and spread of this world-respected forum, which acted as a model for many subsequent community organisations.

Rob was also one of the key figures in creating the RIPE NCC, the body responsible for managing the IP address space in Europe, the Middle East, and parts of Central Asia and coordinating the technical community in those regions. The RIPE NCC was the first *Regional Internet Registry* (RIR) in the world, and this model has become the accepted way to organise the Internet infrastructure in a more regionally specific, responsive, and efficient way.

Rob had a particular talent for being able to engage with all elements of the Internet community, from government and experienced operators to more recent members of the RIPE community to whom he could impart his insight and wisdom on the issues of the day. He brought much common sense to otherwise complicated discussions in the Internet community, and his mantra of “keep it simple” is one that he will be remembered for.

Those who knew Rob personally will miss his sense of humour. He was a storyteller with an outstanding ability to recall and relate the events from his working life, which not only amused his listeners but also enlightened them and informed their discussions. The recent RIPE 71 Meeting was the first not to be attended by Rob due to his illness, and so it was the first where he was not to be found with a cigarette and glass of wine in hand, enjoying the company of those who typically gathered outside the venue entrance to talk about serious matters in a very non-serious way.

His contributions were often officially recognised, notably in receiving Dutch royal honours by being awarded with the title *Officer in the Order of Oranje-Nassau* in 2010. He also received the *Jonathan B. Postel Service Award* in 2015 for outstanding contributions in service to the data communications community. Since standing down as RIPE Chair in 2013, he has enjoyed the title RIPE Chair Emeritus.

To many of us in the RIPE community and beyond, Rob was a mentor, a friend, a trusted confidante and always the voice of reason. His legacy stretches from the physical networks the Internet is made of to the community he built and the wisdom he injected into that community’s make-up from the very beginning. His legacy will continue to be felt as the community continues to grow and its participants often ask themselves, “What would Rob do?”

The RIPE NCC has set up a webpage where you can leave your own tribute to Rob:

<https://labs.ripe.net/Members/mirjam/tribute-to-dr-robert-blokzij-1943-2015>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Publication of this journal is made possible by:

Supporters



Diamond Sponsors



Ruby Sponsor



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



Individual Sponsors

Lyman Chapin, Steve Corbató, Dave Crocker, Jay Etchings, Martin Hannigan, Hagen Hultzs, Dennis Jennings, Jim Johnston, Merike Kaeo, Bobby Krupczak, Richard Lamb, Tracy LaQuey Parker, Bill Manning, Andrea Montefusco, Tariq Mustafa, Mike O'Connor, Tim Pozar, George Sadowsky, Scott Seifel, Helge Skrivervik, Rob Thomas, Tom Vest, Rick Wesson.

For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

ADDRESS SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Fred Baker, Cisco Fellow
Cisco Systems, Inc.

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Steve Crocker, Chairman
Internet Corporation for Assigned Names and Numbers

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

March 2016

Volume 19, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
What's in a DNS Name?	2
QoS and QoE	14
Fragments	41
Call for Papers	42
Supporters and Sponsors	43

FROM THE EDITOR

This *Internet Engineering Task Force* (IETF) has just completed its meeting in Buenos Aires, Argentina. This meeting was the first time the IETF has met in South America, and while Buenos Aires is “far away” from many parts of the world, the technical community seems to agree that this historic meeting was well worth the journey. You can read more about this meeting on the IETF and Internet Society websites, as well as in the latest issue of the *IETF Journal*, which on this special occasion is available in both English and Spanish.

This year, the IETF will meet in Berlin, Germany, in July and in Seoul, Korea, in November. If you are involved with developing or deploying Internet protocols, I recommend that you attend an IETF meeting if you have not done so already. There are also ways to participate remotely in these meetings if you are unable to attend in person, and the IETF now has many resources for first-time attendees at its meetings. For more information about the IETF, visit <http://ietf.org>

The *Domain Name System* (DNS) is the “human face” of the Internet, allowing us to use terms such as **facebook.com** or **isoc.org** to access a service over the Internet. A small number of these names are “reserved” in the sense that they do not appear in the global DNS system. In our first article, Geoff Huston discusses the *Special-Use Domain Name* registry, which has sparked quite a bit of debate in recent months.

Quality of Service (QoS) has been discussed in several articles in this journal over the years. In our second article, William Stallings and Florence Agboma describe QoS and the more recent concept of *Quality of Experience* (QoE) as it relates to IP networks.

We would like to remind you that this journal depends on the generous support of numerous individuals and organizations. If you would like to help support IPJ, please contact us for further details. Comments, suggestions, book reviews, and articles are always welcome. Send your messages to ipj@protocoljournal.org

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

What's in a DNS Name?

by Geoff Huston, APNIC

What's the difference between `.local` and `.here`? Or between `.onion` and `.apple`? All four of these labels are capable of being represented in the Internet's *Domain Name System* (DNS) as *generic Top Level Domains* (gTLDs), but only two of these are in fact delegated names. The other two, `.local` and `.onion` not only don't exist in the delegated name space, but by virtue of a registration in the *Internet Assigned Numbers Authority* (IANA)'s *Special Use Domain Name* registry^[1], these names cannot exist in the conventional delegated domain name space.

It seems that Internet does not have a single coherent name space, but instead it has a name space that contains a number of silent and unsignalled fracture lines, and instead of being administered by a single administrative body there are numerous people who appear to want to have a hand on the tiller! Let's look at the Internet's domain name space and try to gain some insight as to how we've managed to get ourselves into this somewhat uncomfortable position.

A Very Brief History of the DNS

It is probably an impossible challenge to consider many years of development and take the outcome of many discussions, conferences, as well as countless millions of mail messages and generate a brief but complete history of the domain name system. Here I'll offer a personal interpretation of what I recall, supplemented with reference to numerous useful sources, but nevertheless it's still a somewhat subjective narrative.

A good place to start is probably RFC 920^[2], authored by Jon Postel and Joyce Reynolds, and published in October 1984. The name model of the Internet had broken away from many other contemporary "flat" or limited hierarchy naming models used in other computer networks by adopting a hierarchical name scheme that imposed no a priori limit on the depth of the hierarchy. This meant that the apex level of the name hierarchy could be limited to a number of generic category names, leaving the lower levels of the name space hierarchy to be populated by individual name instances.

This document, RFC 920, specified a division of the apex level of the name space into a small set of so-called *Top Level Domains*. These were the category-based names of `.com`, `.edu`, `.gov`, `.mil` and `.org`, the collection of two-letter country codes as administered by the *International Organization for Standardization* (ISO) and published as ISO-3166^[22], and a temporary name of `.arpa`.

By the time RFC 1034^[3] was published in 1987, there was no distinction drawn between the name space itself and the technology of resolution of these names.

The name space and the name resolution technology that operated on this name space was collectively referred to as the *Domain Name System* (DNS). At the time, the name space was a collection of top-level names overseen by the IANA. Even in those early days there was pressure to expand the set of delegated top-level domains. Initially **.net** was added, then **.int**, but these additions appeared to exacerbate the issue rather than relieve these growing pressures. As well as efforts to clarify the nature and administration of the domain name space at the time^[4], the debate over who and how the name space could be further expanded continued as a sometime vexatious topic, particularly as the set of stakeholders and interested parties began to grow. Subsequent investigation to expand the DNS name space was undertaken by the *International Ad Hoc Committee* (IAHC)^[5], sponsored by the IANA, the *Internet Architecture Board* (IAB) and the Internet Society, with membership drawn from a number of bodies including the *Telecommunication Standardization Sector of the International Telecommunication Union* (ITU-T) and the *World Intellectual Property Organization* (WIPO). This committee produced a report that advocated the limited expansion of the collection of gTLDs by adding a further seven top-level labels to the domain name space, as well as proposing some structural changes in the name registration function to delineate the roles of name registrars and name registry operators.

However, perhaps of greater interest were other activities that were underway at the same time as the committee was undertaking its investigation. In 1995 the *National Science Foundation* (NSF) had authorized a company called Network Solutions to operate the names registry for the Internet, and permitted the company to charge an annual fee to maintain a name registry entry, and to keep the proceeds from this operation. This situation caused a significant level of discontent, as there was a general perception that the registration fee was unrelated to the cost of operation of the registry and that the registry operator was exploiting a de facto monopoly position to its benefit. A number of activities emerged in alternate name systems. These alternate name systems used the same name structure, and the same name resolution tools, but used a different set of “root” name servers. These systems were so defined to sit alongside the incumbent name system, but added a number of additional top-level labels (see, for example the Wikipedia account of AlterNIC’s brief history^[5]). At issue here was the coherence of the Internet’s name system. A user whose domain name resolvers were positioned within the name space as defined by one of these alternate name systems could use a name in a communication to another user where the same name may have been defined in a different name system and resolved in an entirely different manner.

In May 2000 the IAB published RFC 2826^[7], which argued strongly for the presentation of a single root system and thereby argued strongly for a single coherent name system: “There is no getting away from the unique root of the public DNS.”

Rather than having the DNS name space grow from the “bottom up” in several uncoordinated grass roots efforts to expand the name space, and allowing each effort to fail or survive on the level of public interest and commercial uptake, the IAB was espousing a view that any such expansion of the name space was to be a top-down effort. All such new top level names were to be implemented in a coherent manner such that all such names were visible to all Internet users at the same time. Any expansion of the domain name space was intended to be a process that included all parts of the Internet, and that at all times all public DNS names were to be equally and uniformly available to all users.

However, at much the same time as this statement was made, mid-2000, the IAB was also attempting to extricate itself and the *Internet Engineering Task Force* (IETF) from the fraught debate about the accountability of the IANA, and the nature of the role of the US Government agencies that had been funding the work of the IANA. This debate also folded in the discussion of the further expansion of this domain name space. Evidently many people at the time were interested in seeing a distinct community of interest focus on the issue of the policy of the domain name space in a manner similar to the evolution of the addressing community and the emergence of the *Regional Internet Registry* model in the 1990s. In June 2000 the IAB entered into an agreement with The *Internet Corporation for Assigned Names and Numbers* (ICANN) that effectively passed over the administrative purview of the domain name space, apart from “assignments of domain names for technical uses,” to ICANN, RFC 2860^[8].

From that point onward the focal point for the debate about the expansion of the name space, and the related debate about the monopoly position of Network Solutions was essentially ICANN. Over the ensuing years ICANN made a number of decisions in the interest of addressing perceived needs that were voiced from the community of interest. The roles of the registry and the front-end registrar function were cleaved apart and competition between registrars allowed the retail price of name registrations to be subject to competitive market pressures. In addition, a number of new gTLDs were added in a relatively ponderous and deliberative process. In 2000 the gTLDs of **.aero**, **.biz**, **.coop**, **.info**, **.museum**, **.name** and **.pro** were added to the delegated name set of the root zone of the Domain Name System. Four years later a second round saw the addition of **.asia**, **.cat**, **.jobs**, **.mobi**, **.port**, **.tel**, **.travel** and **.xxx**.

This conservative approach to augmenting this root zone delegated name set changed with the so-called “new gTLD” program, that started in 2008 with the adoption by the ICANN Board of a number of policy recommendations relating to the expansion of the gTLD delegated name space, and the subsequent 2011 launch of this program.

The application window opened on January 12, 2012, and ICANN received 1,930 applications for new gTLDs. On December 17, 2012, ICANN held a prioritization draw to determine the order in which applications would be processed during initial evaluation and subsequent phases of the program. One view is that these names were effectively sold into the market, with an application fee of \$185,000 USD per name. An alternate view was that the application process now entailed significant levels of analysis of the impact on the broader environment, including considerations of competition, security, collisions, potential trademark infringement and similar subjects, and that this fee was intended to cover some portion of the costs of this investigation of the potential impact of the delegation of this particular name as a new gTLD.

Name Tensions and Collisions

Expanding the gTLD name space did not address all of the outstanding issues, and to some extent these tensions were exacerbated by the chosen mechanism for this expansion. The new names and their “owners” were defined essentially by the actions of bidding for names. Without putting too fine a point on it, the expansion of the Domain Name System was passed to a market-based mechanism that was based on foundations of a commercial model of monetization of the name space. This shift appears to have prompted other forms of use of non-delegated top-level domain names to be a little more visible.

There are a number of examples of this change in the landscape of the domain name space.

Local Names

The first of these is the use of the name space in *private* domains. Although the public name space is held together with the coordinated set of root name servers and a common convention that all public name resolvers use these root name servers to establish content, this is only a convention for the public name space. Within private environments it is quite common to see name servers that define a local name environment as a local convenience. For example, you could call the local data server in your home network **server.home**. Not only is that name convenient for the home user, it's convenient for a vendor of home equipment, who can preconfigure server equipment and use these local private names in a pre-configured mode. There are a many names that are commonly used in private environments, probably as a result of vendors in this market domain adopting particular name conventions. The names **.home**, **.hometown**, **.belkin**, **.lan**, **.dlink**, and **.local** are all popular names in locally defined private DNS domains^[9].

What happens if ICANN were to delegate a new gTLD that was the same as a name that enjoyed considerable levels of private use? The two different interpretations of the name would interact.

These days mobility is an important consideration, and a mobile end-system configured with the name of a resource in the private name space would anticipate a “no such domain” response when the system was relocated into the public space where the name was not delegated. Delegation of the name in the public DNS may cause an unanticipated response. Equally, the public space, and the services and resources accessible via these public DNS names would not be visible within the local scope where the name is defined in a private use context. Of course this situation poses some rather challenging policy issues in the name space. Does “squatting” on a name in a private use context confer any rights on tenure of the public name? Should the public name space avoid all names used in private contexts? Given the uncoordinated use of names in private contexts is any form of common regulation of the name space even possible in this context?

Non-DNS Names

The second example is the name space that is associated with non-DNS resolution mechanisms. One of these mechanisms is *multicast DNS* (mDNS), defined in RFC 6761^[10], which replaces the conventional unicast DNS query to a specific DNS resolver with a multicast group query, directed to the link local multicast address (224.0.0.251 or ff02::fb). All members of the multicast group receive the query and the holder of the queried name can identify itself in a multicast answer. All members of the group can learn the answer in this manner. In addition to the change of the resolution mechanism from unicast to link local multicast, RFC 6762^[11], requested that the IETF (not ICANN) reserve the generic top-level domain **.local** for use by mDNS, and thereby prevent ICANN from making a conventional unicast global public DNS delegation of the same top-level name. A related specification, *Link Local Multicast Name Resolution*, defined in RFC 4795^[12] using the Multicast group address of 224.0.0.252 and ff02::1:3, elected not to define an associated name space, so the mDNS approach was unique in some respects.

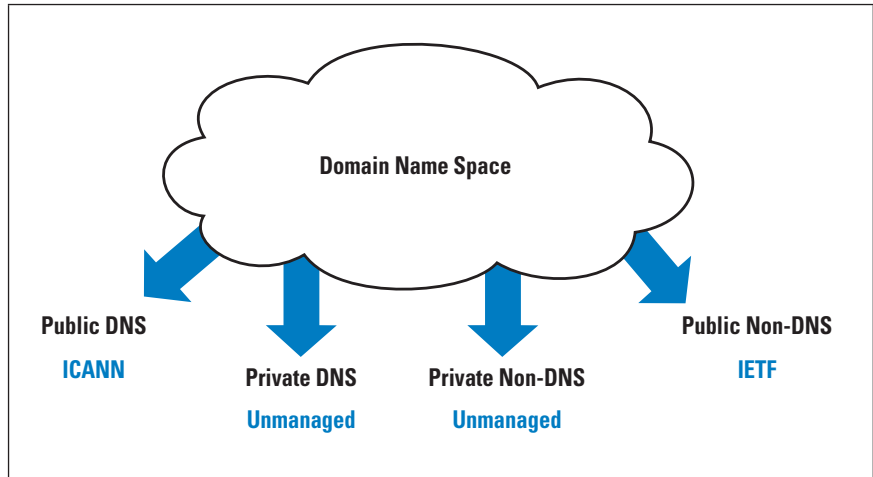
Another approach of non-DNS use of names in the domain space is the *The Onion Ring's* (TOR) use of names in the **.onion** space. Here the names within the **.onion** name space are in effect the base 32 encoded version of the public key of a defined service point, and the TOR-defined Service Directory servers are capable of performing a mapping from an encoded public key (the **.onion** name) and the desired service address. These names are not directly resolved by the DNS and connection requests for **.onion** services need to be passed into the TOR network space for resolution, RFC 7686^[13].

A third name falls into this category, and it predates the other two names by many years. The name **.localhost** refers to the local systems without further recourse to any name resolution process. It is the canonical name used to refer to oneself in the name system.

The Domain Name Space

The overall result of this process of drawing names for use out of the overall domain name space, and the entities that have some level of purview over this process is shown in Figure 1.

Figure 1: Domain Name Space Delegations



The ICANN process views the domain name space as a public good in an economic sense, and uses monetization as an intrinsic component of the name allocation function. In theory, the name space is accessible to those with an exploitation model that can recoup of expenses of acquisition of the name. In practice, the name space is accessible to those with the means to purchase a name, and there is no particular assurance that any of these names will be used in a public context. While second level names are pretty much universally accessible in `.com` or `.net`, for example, the same is probably not the case for `.google`. What was a relatively uniform common public space at the apex level of the delegated name space is now being fenced into a number of realms, many of which are private.

The publication of RFC 6761^[10] by the IETF in February 2013 essentially opened up a competing and uncoordinated channel for drawing of top-level domain names from the domain name pool. In publishing this document the IETF took what was until then a relatively static view of reserved DNS names as described in RFC 2606^[14] in 1999, and replaced it with a process that reopened up the IETF-managed name registry, using the criteria that:

“If a domain name has special properties that affect the way hardware and software implementations handle the name, that apply universally regardless of what network the implementation may be connected to, then that domain name may be a candidate for having the IETF declare it to be a Special-Use Domain Name and specify what special treatment implementations should give to that name.”
[RFC 6761]

This action is effectively unilaterally rephrasing (or “recanting”) the agreement expressed in RFC 2860 and re-defining it to mean that ICANN has purview of only those domain names that use the DNS resolution protocol, and that if the domain name uses a name resolution mechanism that does not rely on this protocol, then the name can be assigned by the IETF, via the IETF publication process. Evidently there is a set of names that are queued up to be listed using this IETF process instead of undertaking the ICANN new gTLD path^[15]. These include **.bit** (using *namecoin* resolution), **.exit** (another TOR-related name), **.gnu** and **.zkey** (using *GNU Name System* resolution), **.i2p**, **.tor** and **.carrots**.

In addition to these two parallel channels of name assignment, the private use activity continues, and names are co-opted into local use domains without any degree of effective coordination.

Clearly this story does not look good. The existence of numerous of uncoordinated activities all drawing out names from a common domain name pool is not a stable situation, nor is it in the interests of the Internet’s users. How is a user to know that names drawn from **.bit** are to be resolved using a namecoin resolution mechanism, whereas names in **.bi** or **.bid** are to be resolved using the DNS resolution protocol?

Differentiating Names?

Are there better ways to signal the resolution protocol that should be applied to a name using some additional signalling?

Should we be thinking about using a *Uniform Resource Identifier* (URI)-like syntax and using distinct schemes, such as **DNS:www.example.com** and **GNS:test.gnu**? Or using a “selector” field in a URI and using URIs of the form: **http:/namecoin/namecoin-string**?

Alternatively, we could try to push these alternate names into a single distinguished gTLD, such as **.alt**, and allow the registrars for **.alt** to register such non-DNS names in a single location in the DNS name space^[16].

We could borrow a technique used by *Internationalized Domain Names* (IDNs) and use a common prefix to denote a non-DNS name, in the same way that the character string prefix “**xn--**” denotes that the following parts of this label require pre-processing in order to produce the equivalent Unicode string. This possibility would imply that all other name forms would form part of a single name space with a single name resolution protocol, while the exception space would be clearly denoted by such a distinguished name prefix, such as, hypothetically, **.xs--gnu** for *Gnu Name System* names and **.xs--bit** names, and so on.

Behind these approaches lies a common question: What are these alternate name forms and name resolution protocols really addressing? What is the underlying issue here? If they are addressing shortfalls in the DNS, such as its lack of privacy for example, then is the appropriate answer one that includes the use of a parallel alternative name resolution protocol, or should we be looking towards the evolution of the DNS protocol to accommodate these emerging requirements? If they are addressing the ICANN position that has apparently monetized the gTLD name space and thereby blocked various other interests from accessing a gTLD name, then is the most appropriate measure for the IETF to set up of a parallel name allocation mechanism? Should the names community within ICANN undertake some deeper introspection and examine whether the gTLD program is actually catering for the full spectrum of interests in securing names for their various needs?

One Name Space?

What may be useful here is the observation that this is not a unique problem.

The radio spectrum has gone through the same process a number of times during its 100-year history, looking at the competing interests wanting access to the radio-frequency spectrum. The current spectrum allocation model contains a mix of exclusive use access arrangements. There are commercial exploitation models where actors bid for exclusive use licenses and public interest allocation models where various public sector agencies are assigned spectrum space. There are public interest and scientific use allocations, such as those used by emergency services and radio astronomers. There are also unlicensed radio spectrum allocations where there is no arrangements for exclusive access, such as are used by WiFi systems. Although a national spectrum management body is not raising revenue from these unlicensed allocations, the economic benefits of WiFi are doubtless substantial, and there is a net benefit to national economies in having this diversity of spectrum access models. The insight here is the admission that the common pool of radio spectrum space does not necessarily admit to a single exploitative model of exclusive access arrangements, and allowing a diversity of models, including that of unlicensed access, has proved to be a useful framework.

What is evident is that ICANN's gTLD process has evidently not encompassed the plurality of demand for domain names. One characterization of the outcomes of the policy for new gTLDs is that it has encouraged competitive access within a relatively narrow model of use. Access to further gTLDs within this process has many barriers, including not inconsiderable financial outlays and process overheads. The reaction has been for numerous parties to look to the IETF's management of the Special-Use Domain Name registry as an alternate means of reserving a domain name and precluding it from being used by ICANN's new gTLD program.

The rationale for entry into the IETF's Special Use Names registry—namely that the name in question uses a non-DNS resolution protocol—could be argued to be a superficial artifice that hides a more significant issue about the broad variety of the natural demand for use of names drawn from the common pool of the domain name space, and the consequent pressure for a range of means for such demands to be satisfied.

There is no doubt that the Internet's users benefit from a single coherent name space. There is considerable benefit in having the same 'name' encompass the same semantic intent and thereby 'name' the same set of services irrespective of the context, locale or time of use of that name. At the same time, the underlying technologies of name resolution, including not only the DNS resolution protocol but also other forms and means of name resolution, are subject to evolutionary pressures. It is valuable to have a means to expose these exploratory efforts in an environment of scale of use, and clearly the IETF has a role to play here. But the current mechanism of having these two bodies making uncoordinated allocations from a common name pool is not an ideal situation.

What leads to some level of unease here about the coherence of the name space is the radically different processes of evaluation of the name itself.

In the ICANN case the new gTLD process requires evaluation of the name to ensure that does not unduly infringe on existing name use, including consideration of existing brands and trademarks, designation of origin and geographic terms, issues of consumer protection through consideration of name similarity and forms of intentional passing off, potential clashes with names used by recognized international organizations, offensive terms, potential name collisions and similar environmental concerns. One could argue with the effectiveness of the process used by ICANN to evaluate these considerations, but there is some merit in the intent of ensuring that there is a process that is mindful of the larger environment of name use when considering adding a further name into this pool of use by the Internet.

The IETF's evaluation process described in RFC 6761 for admission to the Special Use Names Registry appears to admit no such similar consideration. The seven questions posed in RFC 6761 are concerned primarily with the impact of this hidden "switch" that directs applications, name resolvers, and users to understand that this name is not to be resolved by the DNS. None of these questions are concerned with the name itself, and the consequent concern is that this process could be readily abused to legitimate name squatting, and be the source of various forms of name collision. For example, the reservation of the label `.local` in the Special Use Names Registry collides with extensive conventional DNS use in local contexts^[9].

There is no record of any evaluation by the IETF of the consequences of a registration of **.local** as a reserved name for use in non-DNS contexts, with its implicit switch to a different resolution protocol of multicast DNS, that collides with pre-existing use of **.local** names in conventional local DNS contexts. Nor was there any evidence that there was consideration given to mobile users who may move in and out of environments where names in **.local** have entirely different properties and meanings, and the security issues that could result from such confusion.

However, this situation is definitely not a case of “ICANN good, IETF bad!” Far from it! But it does illustrate that there is much more to a name than might appear at the outset. The name space is indeed larger than just the DNS name resolution protocol, and this is perhaps something for ICANN to consider. At the same time names exist in a larger context of social and technical use, and this is something for the IETF to consider if it wishes to accept further reservations in the Special Use Name registry. There is also the consideration of the larger issue of whether implicit (and largely invisible) name-triggered resolution protocol switches are really in the best interests of Internet users. And for those vendors and network administrators looking for local use names to support various form of plug and play, there is the consideration of name collisions and the potential security concerns for unsuspecting users when end systems move in and out of local environments where certain name forms take on altered meanings and altered contexts of use.

If we think that a coherent and consistent name space for the Internet still has some intrinsic value, then we simply have to make some changes here to allow for a broader diversity of name use for the Internet. At the same time we must avoid stomping wilfully on each other’s toes!

References and Annotated Reading List

This is a topic that has been considered for some time, and at some considerable length, so there is no shortage of material on the topic of the name space of the Internet.

[1] Special Use Domain Names:

<http://www.iana.org/assignments/special-use-domain-names/special-use-domain-names.xhtml>

[2] Jon Postel and Joyce Reynolds, “Domain Requirements,” RFC 920, October 1984. *An early description of the DNS and its intended use.*

[3] Paul Mockapetris, “Domain Names – Concepts and Facilities,” RFC 1034, November 1987. *The original canonical specification of the DNS.*

- [4] Jon Postel, “Domain Name System Structure and Delegation,” RFC 1591, March 1994. *A description of the structure of the Domain Name space as of the early 1990’s and the description of roles and responsibilities of domain name managers.*
- [5] IAHC: <https://en.wikipedia.org/wiki/IAHC>
- [6] AlterNIC: <https://en.wikipedia.org/wiki/AlterNIC>
The Wikipedia account of the AlterNIC episode of the 1990s.
- [7] Internet Architecture Board, “IAB Technical Comment on the Unique DNS Root,” RFC 2826, May 2000. *The Internet Architecture Board’s comment on one name space and its utility for the Internet.*
- [8] Brian Carpenter, Fred Baker, and Mike Roberts, “Memorandum of Understanding Concerning the Technical Work of the Internet Assigned Numbers Authority,” RFC 2860, June 2000. *The agreement between the IETF and ICANN over names and addresses.*
- [9] Geoff Huston and George Michaelson, “On Queries to the Root,” <http://www.potaroo.net/presentations/2014-06-24-namecollide.pdf>
- [10] Stuart Cheshire and Marc Krochmal, “Multicast DNS,” RFC 6762 February 2013. *The specification of the Multicast DNS resolution protocol and the reservation of .local in the Special Use Names Registry.*
- [11] Stuart Cheshire and Marc Krochmal, “Special Use Domain Names,” RFC 6761, February 2013. *The document that re-opened the special use name registry using a rationale of non-DNS resolution technology.*
- [12] Levon Esibov, Dave Thaler, and Bernard Aboba, “Link-Local Multicast Name Resolution (LLMNR),” RFC 4795, January 2007. *A non-DNS locally scoped name resolution protocol specification.*
- [13] Jacob Appelbaum, Alec Muffet, “The “.onion” Special-Use Domain Name,” RFC 7686, October 2015. *The justification for listing .onion in the Special Use Names register.*
- [14] Donald E. Eastlake 3rd and Aliza R. Panitz, “Reserved Top Level DNS Names,” RFC 2606, June 1999.
- [15] Christian Grothoff, “Special Use Domain Names of P2P Systems,” <https://www.ietf.org/proceedings/93/slides/slides-93-dnsop-5.pdf>

- [16] Warren Kumari and Andrew Sullivan, “The ALT Special Use Top Level Domain,” Internet Draft, work in progress, September 2015, <https://tools.ietf.org/html/draft-ietf-dnsop-alt-tld-03>, *A working draft in the DNS Operations Working Group of the IETF that considers the use of .alt as a common top-level domain for Special Use contexts.*
- [17] Randall Atkinson Sally Floyd, “IAB Concerns and Recommendations Regarding Internet Research and Evolution,” RFC 3869, August 2004. *Section 3.2 considers research topics in the general area of names research.*
- [18] DNSOP Working Group – Proceedings of WG Meeting, IETF 93, July 2015, <https://www.ietf.org/proceedings/93/slides/slides-93-dnsop-5.pdf>, *A list of other names that are candidates for listing in the Special Use Names Registry.*
- [19] NSRG, “What’s In A Name: Thoughts from the NSRG,” abandoned work, September 2003, <https://tools.ietf.org/html/draft-irtf-nsrg-report-10>, *A report from the Name-space Research Group (NSRG) of the Internet Research Task Force (IRTF). It appears that the report failed to achieve consensus within the group, and was never published as an RFC.*
- [20] Joe Abley, Peter Koch, and Alain Durand, “Problem Statement for the Reservation of Top-Level Domains in the Special-Use Domain Names Registry,” <https://tools.ietf.org/html/draft-adpkja-dnsop-special-names-problem-00>, *One of those rare cases where the document’s title says it all.*
- [21] ICANN, “New gTLD program,” <http://newgtlds.icann.org/en/>, *A description of ICANN’s gTLD expansion program.*
- [22] ISO 3166 Country Codes: http://www.iso.org/iso/country_codes, *This specification is the product of TC 46, a committee of the International Organization for Standardization. This body was formed in 1947 as the outcome of collaboration between ISA, the International Federation of National Standardizing Associations, and the UNSCC, the United Nations Standards Coordinating Committee. Strictly speaking ISO does not define countries nor their standard names—that is the product of the United Nations Statistics Divisions. However ISO manages the assignment of 2 and 3 letter codes to countries, recommended as a “general purpose code.”*

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic

Quality of Service and Quality of Experience: Network Design Implications

by William Stallings, Independent Consultant
Florence Agboma, BSkyB

The Internet and enterprise IP-based networks continue to see rapid growth in the volume and variety of data traffic. Cloud computing, big data, the pervasive use of mobile devices on enterprise networks, and the increasing use of video and image traffic all contribute to the increasing difficulty in maintaining satisfactory network performance. Two key tools in measuring the network performance that an enterprise desires to achieve are *Quality of Service* (QoS) and *Quality of Experience* (QoE). QoS and QoE enable the network manager to determine if the network is meeting user needs and to diagnose problem areas that require adjustment to network management and network traffic control. This article provides an overview of QoS and QoE concepts, the relationship between the two, and the design implications of a QoE/QoS architecture.

Background and Overview

Historically, the Internet and other IP-based networks provided a *best-effort* delivery service. The term “best effort” refers to the connectionless, datagram nature of the interconnected set of networks. With best effort, a packet may be lost, duplicated, delayed, or delivered out of order, and the network does not inform sender or receiver. Traditionally, a best-effort network attempts to allocate its resources with equal availability and priority to all traffic flows, with no regard for application priorities, traffic patterns and load, or customer requirements. To protect the network from congestion collapse and to guarantee that some flows do not crowd out other flows, congestion-control mechanisms were introduced; they tended to throttle traffic that consumed excessive resources. As the intensity and variety of traffic increased, various QoS mechanisms were developed, including *Integrated Services Architecture* (ISA) and *Differentiated Services* (DiffServ) (for example, see [1]). *Service-Level Agreements* (SLAs) were also developed so that the service provided to various customers was tunable and somewhat predictable. These mechanisms and services serve two purposes: (1) allocate network resources efficiently so as to maximize effective capacity; and (2) enable networks to offer customers different levels of QoS on the basis of their requirements.

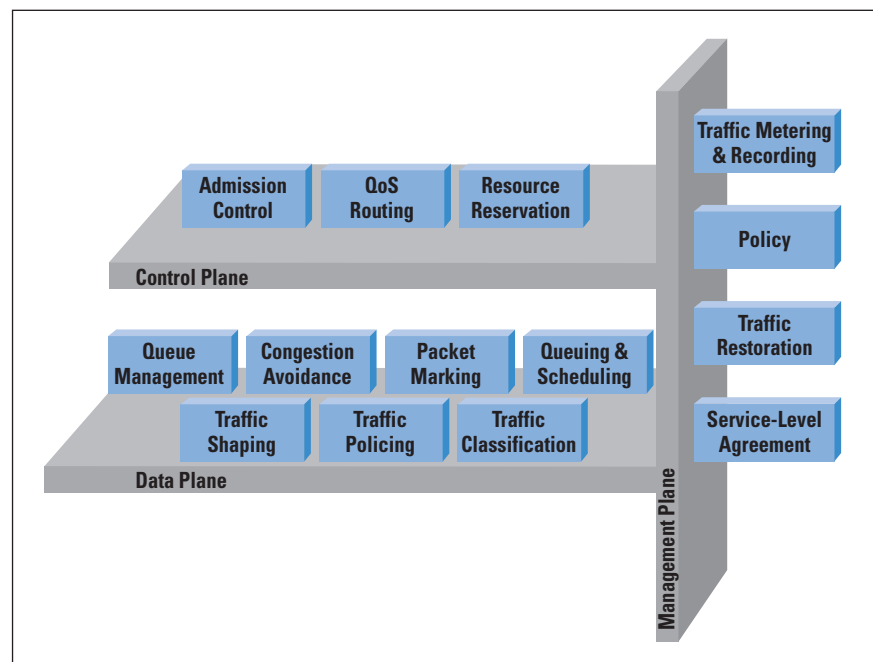
QoS is an important but increasingly insufficient tool for providing network services for many of today’s high-volume applications. To meet the needs of such applications, QoS has recently been augmented with the concept of QoE, which is a subjective measure of performance as reported by the user. Unlike QoS, which can be precisely measured, QoE relies on human opinion. QoE is important particularly when dealing with multimedia applications and multimedia content delivery.

Because QoE extends the concept of QoS to more fully tailor network services and performance to customer and user needs, it is garnering increasing attention by network protocol and system designers. The management of QoE has become a crucial concept in the deployment of future successful applications, services, and products. The greatest challenges in providing QoE are developing effective methods for converting QoE features to quantitative measures and translating QoE measures to QoS measures. Whereas now it is easy to measure, monitor, and control QoS at both the networking and application layers, and at both the end system and network sides, its management is still quite intricate.

QoS Architectural Framework

The *Telecommunication Standardization Sector of the International Telecommunication Union* (ITU-T) Recommendation Y.1291 provides an overall architectural framework that relates the various elements that go into QoS provision^[2]. Figure 1 shows the relationship between these elements, which are organized into three planes: data, control, and management. This architectural framework is an excellent overview of QoS functions and their relationships and provides a useful basis for summarizing QoS.

Figure 1: Architectural Framework for QoS Support



The *Data Plane* includes those mechanisms that operate directly on flows of data. The following discussion briefly describes each mechanism in turn.

Queue Management algorithms manage the length of packet queues by dropping packets when necessary or appropriate. Active management of queues is concerned primarily with congestion avoidance. In the early days of the Internet, the queue management discipline was to drop any incoming packets when the queue was full; it was referred to as the *tail-drop* technique. This technique has many drawbacks, including^[3]:

- There is no reaction to congestion until it is necessary to drop packets, whereas a more aggressive congestion-avoidance technique would likely improve overall network performance.
- Queues tend to be close to full, causing an increase in packet delay through a network and possibly resulting in a large batch of dropped packets for bursty traffic, necessitating many packet retransmissions.
- Tail drop may allow a single connection or a few flows to monopolize queue space, preventing other connections from getting room in the queue.

One noteworthy example of queue management is *Random Early Detection* (RED)^[4]. RED drops incoming packets probabilistically based on an estimated average queue size. The probability for dropping increases as the estimated average queue size grows.

Queuing and Scheduling Algorithms, also referred to as “queuing discipline algorithms,” determine which packet to send next; they are used primarily to manage the allocation of transmission capacity among flows.

Congestion Avoidance deals with means for keeping the load of the network sufficiently under its capacity such that the network can operate at an acceptable performance level. The specific objectives are to avoid significant queuing delays and, especially, to avoid congestion collapse.

Packet Marking encompasses two distinct functions. First, packets may be marked by network edge nodes to indicate some form of QoS that the packet should receive. An example is the *DiffServ* (DS) field in the IPv4 and IPv6 packets and the *Traffic Class* field in *Multiprotocol Label Switching* (MPLS) labels^[5]. An edge node can set the values in these fields to indicate a desired QoS. Such markings may be used by intermediate nodes to provide differential treatment to incoming packets. Second, packet marking can be used to mark packets as nonconforming; such packets may be dropped later if congestion occurs.

Traffic Classification refers to the assignment of packets to a traffic class at the edge of the network. Typically, the classification entity looks at multiple fields of a packet, such as source and destination address, application payload, and QoS markings, and determines the aggregate to which the packet belongs. This classification provides network elements a method to weigh the relative importance of one packet over another in a different class. All traffic assigned to a particular flow or other aggregate can be treated similarly. The flow label in the IPv6 header can be used for traffic classification.

Traffic Policing determines whether the traffic being presented, is on a hop-by-hop basis compliant, with prenegotiated policies or contracts. Nonconforming packets may be dropped, delayed, or labeled as nonconforming. As an example, ITU-T Recommendation Y.1221^[6] recommends the use of token bucket to characterize traffic for purposes of traffic policing.

Traffic Shaping controls the rate and volume of traffic entering and transiting the network on a per-flow basis. The entity responsible for traffic shaping buffers nonconforming packets until it brings the respective aggregate in compliance with the traffic. The resulted traffic thus is not as bursty as the original and is more predictable. For example, Y.1221 recommends the use of leaky bucket and/or token bucket for traffic shaping.

The *Control Plane* is concerned with creating and managing the pathways through which user data flows. *Admission Control* determines what user traffic may enter the network. This decision may be in part determined by the QoS requirements of a data flow compared to the current resource commitment within the network. But beyond balancing QoS requests with available capacity to determine whether to accept a request, there are other considerations for admission control. Specifically, network managers and service providers must be able to monitor, control, and enforce use of network resources and services based on policies derived from criteria such as the identity of users and applications, traffic/bandwidth requirements, security considerations, and time of day or week. RFC 2753^[7] discusses such policy-related issues.

QoS Routing determines a network path that is likely to accommodate the requested QoS of a flow. This function contrasts with the philosophy of traditional routing protocols, which generally are looking for a least-cost path through the network. RFC 2386^[8] provides an overview of the issues involved in QoS routing, which is an area of ongoing study.

Resource Reservation reserves network resources on demand for delivering desired network performance to a requesting flow. The resource-reservation mechanism that has been implemented for the Internet is the *Resource Reservation Protocol* (RSVP)^[9].

The *Management Plane* contains mechanisms that affect both control- and data-plane mechanisms. It deals with the operation, administration, and management aspects of the network. A *Service-Level Agreement* (SLA) is the agreement between a customer and a provider of a service that specifies the level of availability, serviceability, performance, operation, or other attributes of the service. SLAs are discussed subsequently.

Traffic Metering and Recording concerns monitoring the dynamic properties of a traffic stream using performance metrics such as data rate and packet-loss rate. It involves observing traffic characteristics at a given network point and collecting and storing the traffic information for analysis and further action. Depending on the conformance level, a meter can invoke necessary treatment (for example, dropping or shaping) for the packet stream. A subsequent section discusses the types of metrics that are used in this function.

Traffic Restoration refers to the network response to failures. It encompasses numerous protocol layers and techniques.

Policy refers to a set of rules for administering, managing, and controlling access to network resources. The rules can be specific to the needs of the service provider or reflect the agreement between the customer and service provider, which may include reliability and availability requirements over a period of time and other QoS requirements.

Service-Level Agreements

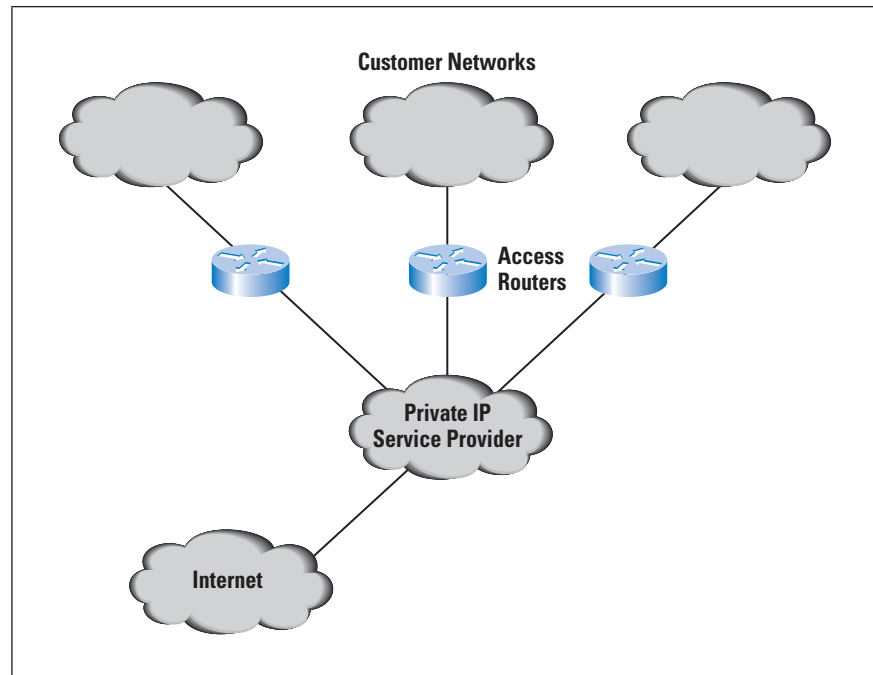
An SLA is a contract between a network provider and a customer that defines specific aspects of the service that is to be provided. The definition is formal and typically defines quantitative thresholds that must be met. An SLA typically includes the following information:

- *A description of the nature of service to be provided:* A basic service would be IP-based network connectivity of enterprise locations plus access to the Internet. The service may include additional functions such as web hosting, maintenance of domain name servers, and operation and maintenance tasks.
- *The expected performance level of the service:* The SLA defines numerous metrics, such as delay, reliability, and availability, with numerical thresholds.
- *The process for monitoring and reporting the service level:* This function describes how performance levels are measured and reported.

The types of service parameters included in an SLA for an IP network are similar to those provided for *Frame Relay* and *Asynchronous Transfer Mode* (ATM) networks. A key difference is that, because of the unreliable datagram nature of an IP network, it is more difficult to realize tightly defined constraints on performance, compared to the connection-oriented Frame Relay and ATM networks.

Figure 2 shows a typical configuration that lends itself to an SLA. In this case, a network service provider maintains an IP-based network. A customer has many private networks (for example, LANs) at various sites. Customer networks are connected to the provider via access routers at the access points. The SLA dictates service and performance levels for traffic between access routers across the provider network. In addition, the provider network links to the Internet and thus provides Internet access for the enterprise.

Figure 2: Typical Framework for Service-Level Agreement



For example, the standard SLA provided by Cogent Communications for its backbone networks includes the following items:

- *Availability*: 100% availability
- *Latency (delay)*: Monthly average Network Latency for packets carried over the COGENT Network between Backbone Hubs for the following regions: Intra-North America: ≤ 45 ms; Intra-Europe: ≤ 35 ms; New York to London: ≤ 85 ms; Los Angeles to Tokyo: ≤ 120 ms.

Latency is defined as the average time taken for an IP packet to make a round trip between Backbone Hubs within a region. COGENT monitors aggregate latency within the COGENT Network by monitoring round-trip times between a sample of Backbone Hubs on an ongoing basis.

- *Network Packet Delivery (reliability)*: Average monthly Packet Loss no greater than 0.1% (or successful delivery of 99.9% of packets). Packet Loss is defined as the percentage of packets that are dropped between Backbone Hubs on the COGENT Network.

An SLA can be defined for the overall network service. In addition, SLAs can be defined for specific end-to-end services available across the carrier's network, such as a virtual private network, or differentiated services.

IP Performance Metrics

The *IP Performance Metrics Working Group* (IPPM) is chartered by the *Internet Engineering Task Force* (IETF) to develop standard metrics that relate to the quality, performance, and reliability of Internet data delivery. Two trends dictate the need for such a standardized measurement scheme:

- The Internet has grown and continues to grow at a dramatic rate. Its topology is increasingly complex. As its capacity has grown, the load on the Internet has grown at an even faster rate. IP-based enterprise networks have exhibited similar growth in complexity, capacity, and load. The sheer scale of these networks makes it difficult to determine quality, performance, and reliability characteristics.
- The Internet serves a large and growing number of commercial and personal users across an expanding spectrum of applications. Similarly, private networks are growing in terms of user base and range of applications. Some of these applications are sensitive to particular QoS parameters, leading users to require accurate and understandable performance metrics.

A standardized and effective set of metrics enables users and service providers to have an accurate common understanding of the performance of the Internet and private internets. Measurement data is useful for a variety of purposes, including:

- Supporting capacity planning and troubleshooting of large complex internets
- Encouraging competition by providing uniform comparison metrics across service providers
- Supporting Internet research in such areas as protocol design, congestion control, and quality of service
- Verification of service-level agreements

The metrics are defined in three stages:

- *Singleton metric*: The most elementary, or atomic, quantity that can be measured for a given performance metric. For example, for a delay metric, a singleton metric is the delay experienced by a single packet.

- *Sample metric*: A collection of singleton measurements taken during a given time period. For example, for a delay metric, a sample metric is the set of delay values for all of the measurements taken during a one-hour period.
- *Statistical metric*: A value derived from a given sample metric by computing some statistic of the values defined by the singleton metric on the sample. For example, the mean of all the one-way delay values on a sample might be defined as a statistical metric.

The measurement technique can be either *active* or *passive*.

Active techniques require injecting packets into the network for the sole purpose of measurement. This approach has several drawbacks. The load on the network is increased, in turn possibly affecting the desired result. For example, on a heavily loaded network, the injection of measurement packets can increase network delay, so that the measured delay is greater than it would be without the measurement traffic. In addition, an active measurement policy can be abused for denial-of-service attacks disguised as legitimate measurement activity.

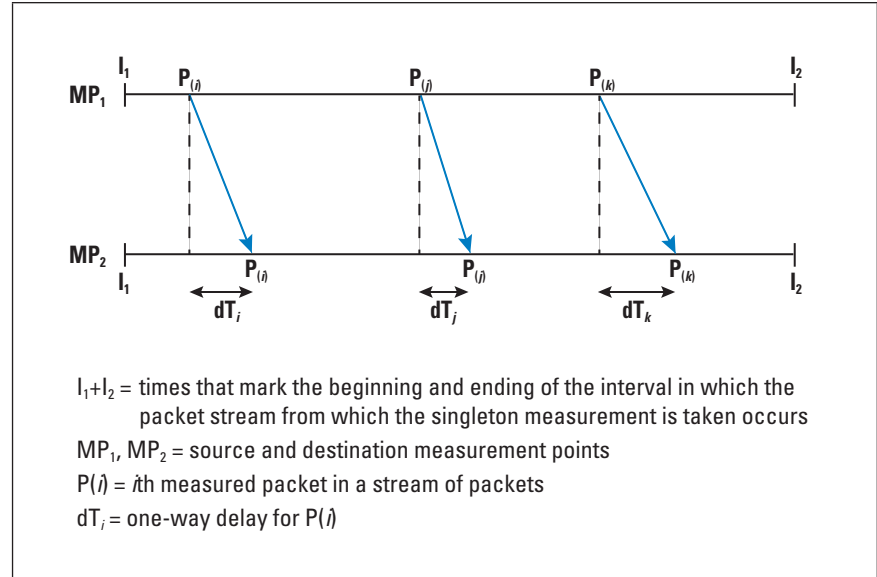
Passive techniques observe and extract metrics from existing traffic. This approach can expose the contents of Internet traffic to unintended recipients, creating security and privacy concerns. So far, the metrics defined by the IPPM working group are all active.

For the sample metrics, the simplest technique is to take measurements at fixed time intervals, known as periodic sampling. There are several problems with this approach. First, if the traffic on the network exhibits periodic behavior, with a period that is an integer multiple of the sampling period (or vice versa), correlation effects may result in inaccurate values.

Also, the act of measurement can perturb what is being measured (for example, injecting measurement traffic into a network alters the congestion level of the network), and repeated periodic perturbations can drive a network into a state of synchronization (for example, [10]), greatly magnifying what might individually be minor effects. Accordingly, RFC 2330^[11] recommends Poisson sampling. This method uses a Poisson distribution to generate random time intervals with the desired mean value.

Figure 3 illustrates the packet-delay variation metric. This metric is used to measure jitter, or variability, in the delay of packets traversing the network. The singleton metric is defined by selecting two packet measurements and measuring the difference in the two delays. The statistical measures use the absolute values of the delays.

Figure 3: Model for Defining Packet-Delay Variation



QoS for Streaming Video

It is worthwhile to comment on the relationship between QoS concerns and streaming video, which perhaps accounts for the greatest volume of traffic on the Internet. First, consider the network transmission requirements for real-time video, such as video teleconferencing. For such applications QoS metrics such as latency and jitter are important, and they impose significant requirements on the networks through which the real-time traffic passes.

For streaming video, in contrast, such QoS requirements can be relaxed by equipping the application to handle a wider range of responses. A motivation for this approach is that QoS is not ubiquitously deployed in networks and thus is not available for all possible streams between content deliverers and content consumers. In the case of streaming video in particular, applications have come to assume that the receiving device is equipped with generous levels of memory, and the application uses this memory as a playback buffer. This system allows the application to use TCP for reliable delivery. The TCP connection can be exposed to far higher levels of network jitter and even packet loss than would be tolerable by a real-time packet stream. Thus, to a significant extent, this approach to video streaming avoids the need to use QoS as a strict precondition for streaming video over the Internet.

From QoS to QoE

The literature contains numerous different, though similar, definitions of QoE. To provide a common working definition, the EU-sponsored *European Network on Quality of Experience in Multimedia Systems and Services* (see [12]) has developed a definition of QoE that reflects broad industry and academic consensus:

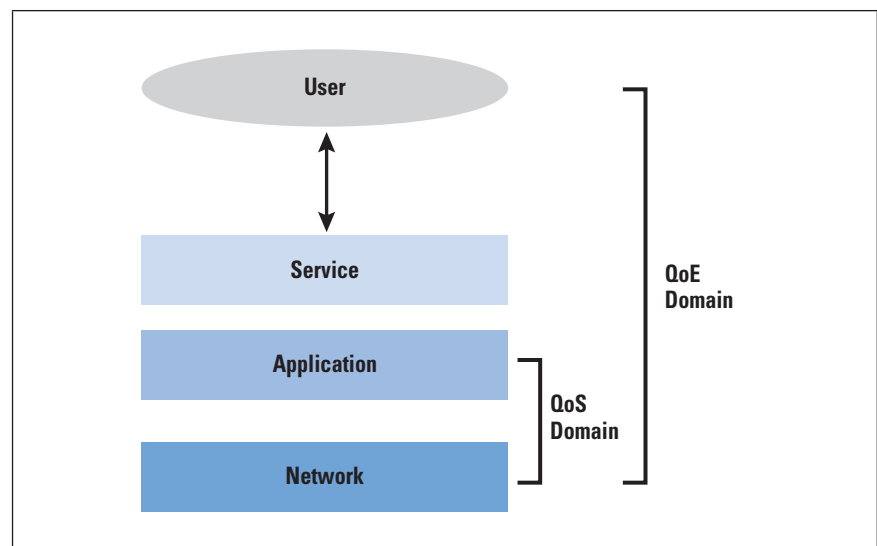
“Quality of Experience (QoE) is the degree of delight or annoyance of the user of an application or service. It results from the fulfillment of his or her expectations with respect to the utility and/or enjoyment of the application or service in the light of the user’s personality and current state.”

The QoE methods to measure this degree of delight or annoyance are discussed later in this article.

The Layers of QoE

For any type of service (for example, IPTV), multiple QoS parameters contribute to the overall user’s perception of quality. The concept of QoE in addition to QoS mechanisms has been proposed in [13], [14], [15], and [16] as a QoE/QoS layered approach by which the requirements of the users drive network-dimensioning strategies. The QoE/QoS layered approach does not substitute the QoS aspect of the network, but instead, user and service-level perspectives are both complementary as shown in Figure 4.

Figure 4: QoE/QoS Layered Model



The levels from the layered approach follow:

- *User*: The user interacts with the service. It is the user’s degree of delight or annoyance from using the service that is to be measured.
- *Service*: This level is the virtual level where the user’s experience of the overall performance of the service can be measured. It is the interface where the user interacts with the service (for example, the visual display to the user).
- *Application-level QoS (AQoS)*: This level deals with the control of application-specific parameters such as content resolution, bitrate, frame rate, colour depth, codec type, layering strategy, sampling rate, and number of channels.

- *Network-level QoS (NQoS)*: This level is concerned with the low-level network parameters such as service coverage, bandwidth, delay, throughput, and packet loss.

Background on the Layers of QoE

Being linked to human perception, QoE is hard to describe quantitatively, and it varies from person to person. The complexities of QoE at the user level stem from the differences between individual user characteristics, of which some might be of a variant nature (that is, changing often with time), whilst others are of a relatively stable nature. Examples could include gender, age, attitudes, prior experience, expectations, socio-economic status, cultural background, educational level, etc. Thus, it becomes a challenge to derive unified QoE metrics for all users and their contexts. Reiter et al.^[17] provide more details on the human factors that may influence QoE. The current practice in any QoE measurement is to identify and control the relatively stable characteristics of a user in a way that is satisfactory to at least a large proportion of the potential user group.

Multiple layers might impact the user QoE for any application. The service level provides a virtual level where the user's tolerance thresholds for a particular service could be measured. As an illustration, the QoE measures from the user perspective for streaming applications could be a) start-up time, b) audiovisual quality, c) channel-change delay, and d) buffering interruptions, whilst the QoE measures for web-browsing applications could be waiting times.

Quite often, the network capacity will dictate the bandwidth that may be used for transmission. At the application-level QoS, there might be a trade-off between quality and size. For audio, a higher sampling rate, for example, 96 kHz, might allow for more information to be perceived compared with 48 kHz, but at the expense of a bigger file size. Traditional telephony speech might be limited to 8 kHz because of the bandwidth capacity. For video, a high resolution might require more bandwidth than low resolution. There are huge varieties of device screens in all kinds of sizes, featuring varied aspect ratios. The one commonality in this array of equipment is that they are all capable of rescaling the video, as is the case for a native player going to the full-screen mode. For a given bitrate, there might be a trade-off between lower resolutions in pixels (images being slightly blurred) with fewer artifacts versus higher resolutions that provide a sharper image but possibly with more artifacts. Most compression standards might use a block-based and motion-compensation coding scheme, and as a result additional compression artifacts are added to the decoded video.

Network-level QoS parameters could impact QoE in a variety of ways.

The network delay could impact QoE, especially for interactive services. For instance, the interactive nature of web browsing that requires multiple retrieval events within a certain window of time might be affected by delay variations of the network. *Voice over IP* (VoIP) services might have the stringent response-time demands, whereas email services might tolerate much longer delays. The different distribution methods of streaming video over the network might affect QoE in different ways. HTTP-based adaptive streaming, which uses TCP, might react to bandwidth constraints and CPU capacity in two ways: either the streaming switches between streaming the different bitrate encodings depending on available resources, or a frame freeze (rebuffering) due to incoming packet starvation in the player buffer. The continuous bitrate switches and rebuffering affect QoE badly. The other distribution method, *User Datagram Protocol* (UDP), might use multicast to replicate the streams throughout the network. Quite often, a resilient coding scheme and a flow-control mechanism might be implemented to maintain the viewing experience despite the effects of poor network conditions.

For reasons that should now be apparent, the background on the layers of QoE suggests that the effect of QoE could be an attribute of only the application layer or a combination of both the application and network layers. Although the trade-offs between quality and network capacity may begin with application-level QoS due to network capacity considerations, an understanding of the user requirements at the service level (that is, in terms of QoE measures) would enable a better choice of application-level QoS parameters to be mapped onto the network-level QoS parameters. A scenario that aims at controlling QoE using QoS parameters as actuators is discussed later in this article. Taking the QoE/QoS approach as a whole entity rather than single entities might aid in providing better QoEs, and potentially could lead to a better-managed delivery infrastructure.

Key Factors That Determine QoE

The nature of QoE, which comprises many layers of interaction between the enabling elements of service delivery and the human user, makes measuring and improving user experiences a challenging task. To understand QoE we must account for both technical and nontechnical factors. Many factors contribute to producing a good QoE. Moller et al.^[18] provide useful perspectives on factors that influence QoE. Here, we discuss the following key factors that might influence QoE:

User Demographics: The context of demographics herein refers to the relatively stable characteristics of a user that might indirectly influence perception, and intimately affects other technical factors to determine QoE. The findings from [19] suggest that user demographics may influence QoE. In the context of adopting HD voice telephony, the different user groups had significantly different quality ratings.

The grouping of users was based on demographic characteristics such as their attitudes towards adoption of new technologies, socio-demographic information, socio-economic status, and prior knowledge. Cultural background is another user demographic factor that might also influence perception due to cultural attitude to quality^[20].

Type of Device: Different device types possess varying characteristics that may affect QoE. An application designed to run on more than one device type, for example on a connected TV device such as Roku and on an iOS device such as an iPhone, may not deliver the same QoE on every device.

Content: The content being distributed via Internet delivery can range from interactive content specifically curated according to personal interests to content that is produced for linear TV transmission. The different characteristics of the video might require different system properties in terms of the quality being produced. According to [21], people tend to watch *Video on-Demand* (VoD) content with a higher level of engagement than its competing alternative, linear TV. This trend may be because users make an active decision to watch a particular VoD content, and as a result, give their full attention to it. One could infer that for VoD, users might be less tolerant of any quality degradations because of their high level of engagement.

Connection Type: The type of connection used to access the service influences users' expectations and their QoEs. Users have been found to have lower expectations when using 3G in contrast to a wireline connection when in fact both connection types were identical in terms of their technical conditions^[22]. Agboma et al.^[23] found users' expectations to be considerably lowered and more tolerant to visual impairments on small devices. Conventional QoS management practices cannot account for these psychological factors.

Media (audiovisual) Quality: This factor may be observed as a significant one affecting QoE, because it is the part of a service that users notice most. The integration of audio and video quality appears to be content-dependent. For less-complex scenes (for example, head and shoulder content), audio quality is slightly more important than video quality. In contrast, for high-motion content, video quality has been found to be significantly more important than audio quality^[24]. Other studies^[25] suggest that the optimum audio/video bitrate allocation depends on scene complexity. For instance, visually complex scenes would benefit from the allocation of higher bitrates, with relatively more bits allocated towards audio, because high-audio bitrates seem to produce the best overall audiovisual quality.

Network: Content delivery via the open Internet is highly susceptible to the effects of delays, jitter, packet loss, and available bandwidth. For users, delay and jitter cause frame freeze and the lack of lip synchronization between what is heard (audio) and what is seen (video). Content delivery is guaranteed using a TCP/IP delivery mechanism. However, under bad network conditions, the frequency of rebuffering, and the implementation of a video player heuristics, might affect QoE. Rebuffering interruptions in IP video playback is considered the worst degradation on user QoE and should be avoided at the cost of startup delay^[26]. On the same note, QoE for a given startup delay strongly depends on the concrete application context and user expectations^[27]. In spite of the different QoE factors that are concerned with the network, reliability and a strong wireless signal is crucial for consuming TV-like services anytime, anywhere, and from any device.

Usability: Another QoE factor is the amount of effort that is required to navigate through the service. The design should render good quality without a great deal of technical input from the user before or after service consumption.

Cost: The long-established practice of judging quality by price implies that expectations are price-dependent. If the tariff for a certain service quality is high, users may be highly sensitive to any quality degradations. A scenario of QoE-based charging is demonstrated in [28] and [29] to analyze a situation where price is used to reflect the value of a service, and at the same time, part of the user-context factor. While [28] offers a trade-off between user expectations and price, its deployment may yield unexpected complexities.

QoE Measurement Methods

QoE measurement techniques evolved through the adaptation and application of psychophysics methods during the early stages of television systems (See [30] for more details). Here, we introduce three QoE measurement methods: subjective assessment methods, objective assessment methods, and end-user device analytics as an alternative to measure QoE. Hereafter, streaming video will be the focus of discussion.

Subjective Assessment

In *Subjective Assessment* of QoE, experiments are carefully designed to a high level of control (such as in a controlled laboratory, field tests, or crowdsourcing environments) so that the validity and reliability of the results can be trusted. It might be useful to consult expert advice during the initial design of the subjective experiment, because the topics of experimental design, experimental execution, and statistical analysis are complex. The different phases discussed in the following paragraphs are an abstract. A methodology to obtain subjective QoE data might consist of the following phases:

Characterize the Service: The task at this stage is to choose the QoE measures that affect the users' experience the most. As an example, consider a multimedia conferencing service; previous studies have shown that the quality of the voice takes precedence over the quality of video^[31]. Also, the video quality required for such applications does not demand a very high frame rate, provided that audio-to-video synchronization is maintained. Thus, the resolution of individual frames can be considerably lower than the case of other video streaming services, especially when the size of the screen is small (such as that of a mobile phone). Therefore, in multimedia conferencing, the QoE measures might be prioritized as voice quality, audio-video synchronization, and image quality.

Design and Define Test Matrix: Once the service has been characterized, the QoS factors that affect the QoE measures can be identified. For instance, the video quality in streaming services might be directly affected by network parameters such as bandwidth, packet loss, and encoding parameters such as frame rate, resolution, and codec. The capability of the rendering device will also play a significant role in terms of screen size and processing power. However, testing such a large combination of parameters may not be feasible. This draft matrix could be reduced to more achievable realistic test conditions by eliminating the combinations that have similar effects on QoE.

Specify Test Equipment and Materials: Subjective tests should be designed to specify test equipment that will allow controlled enforcement of the test matrix. For instance, to assess the correlation between NQoS parameters and the perceived QoE in a streaming application, at least a client device and a streaming server separated by an emulated network are needed. If the objective is to evaluate how different device capabilities impact QoE, then a video content is chosen to produce formats that can run in each of the client devices under scrutiny.

Identify Sample Population: A representative sample population is identified, possibly covering different classes of users categorized by the user demographics that are of interest to the experimenter. Depending on the target environment for the subjective test, at least 24 test subjects (that is, participants) have been suggested for a controlled environment (for example, a laboratory) and at least 35 test subjects for a public environment. Fewer subjects may be used for pilot studies to indicate trending^[32]. The use of *crowdsourcing* in the context of subjective assessment is still nascent, but it has the potential to further increase the size of the sample population and could reduce the completion time of the subjective test.

Subjective Methods: The recommendations include several subjective assessment methodologies, for example [33], [34], and [35]. A recent recommendation^[32] addresses the context of Internet video and distribution of quality television in any environment. These recommendations provide guidelines on topics such as the different subjective test designs, rating scales, and experiment durations.

There has also been some interest in developing alternative subjective methodologies for time-varying system characteristics, see [36] and [37]. Typically, each test subject is presented with the test conditions under scrutiny along with a set of rating scales that allows the correlation of the users' responses with the actual QoS test conditions being tested. There are several rating scales, depending on the design of the experiment. Other scale methods can be found in [32], and acceptance methods in [38].

Analysis of Results: When the test subjects have rated all QoS test conditions, a post-screening process might be applied to the data to remove any erroneous data from a test subject who appears to have voted randomly. Depending on the design of the experiment, a variety of statistical approaches could be used to analyze results. For the sake of brevity, statistical analysis of the results is outside the scope of this article. The most common quantification method is the *Mean Opinion Score* (MOS), which is the average of the opinions collected for a particular QoS test condition. The results from subjective assessment experiments are used to quantify QoE, and to model the impacts of QoS factors. Other authors^[39] have gone beyond deriving QoS MOS functions to QoE management. The rationale and limitations of MOS are discussed later in this article.

Subjective experiments require significant planning and design so as to produce reliable subjective MOS ratings. However, they are time-consuming and expensive to carry out and are not feasible for real-time in-service monitoring. Therefore, the use of objective models is often desirable.

Objective Assessment

In *Objective Assessment* of QoE, computational algorithms provide estimates of audio, video, and audiovisual quality as perceived by the user. Each objective model targets a specific service type. For example, ITU-T P.1201 predicts the impact of IP network impairments for IPTV applications and multimedia streaming applications^[40]. ITU-T J.341 targets video-quality measurement in *High-Definition Television* (HDTV) noninteractive applications^[41]. Other proprietary objective models do exist. For a given scope of quality features, the goal of any objective model is to find the optimum fit that strongly correlates with data obtained from subjective experiments. The following phases presented here should not be considered exhaustive, but they are meant to illustrate a process of obtaining objective QoE data. A methodology to obtain objective QoE data might consist of the following phases:

Database of Subjective Data: A starting point might be the collection of a group of subjective datasets as this list could serve as benchmark for training and verifying the performance of the objective model. A typical example of one of these datasets might be the subjective QoE data generated from well-established subjective testing procedures, as discussed earlier. The selection of the subjective datasets should typically reflect the use cases of the objective model.

Preparation of Objective Data: The data preparation for the objective model might typically include a combination of the same QoS test conditions as found in the subjective datasets, as well as other complex QoS conditions. A variety of pre-processing procedures might be applied to the video data prior to training and refinement of the algorithm.

Objective Methods: Various algorithms in the literature could provide estimates of audio, video, and audiovisual quality as perceived by the user. Some algorithms might be developed, and test subjects trained, for a specific perceived quality artifact, whilst other algorithms and test subjects might be trained for a wider scope of perceived quality artifacts. Examples of the perceived artifacts might include blurring, blockiness, unnatural motion, pausing, skipping, rebuffering, and imperfect error concealment after transmission errors. See [42] for a good overview, and for the development of objective video-quality prediction.

Verification of Results: After the objective algorithm has processed all QoS test conditions, the predicted values might benefit from a post-screening process to remove any outliers, the same concept as applied to the subjective datasets. The predicted values from the objective algorithm might be in a different dimension compared to the subjective QoE datasets. The predicted values might be transformed to the same scale as obtained in the subjective experiments (for example, into the MOSs) to allow for a linear comparison, and so that the optimum fit between the predicted QoE values and subjective QoE data can be obtained. This transformation might be an integral module of the objective model. The statistical analysis that might be applied to calibrate the scale of an objective model is outside the scope of this article.

Validation of Objective Model: The objective data analysis might be evaluated with respect to its prediction accuracy, consistency, and linearity by using different subjective datasets. It may also be worth noting that the performance of the model might depend on the training datasets and the verification procedures. See [43] for more details on calibrating and validating objective models. The *Video Quality Experts Group* (VQEG) validates the performance of objective perceptual models that result in ITU recommendations and standards for objective quality models for both television and multimedia applications^[44]. The practical deployment of such objective models is discussed in the background of QoE measurement methods.

End-User Device Analytics

End-user device analytics could provide an alternative to measuring QoE as experienced by the end user. Real-time data such as the connection time, bytes sent, and average playback rate are collected by the video player application for each video viewing session and fed back to a server module where the data is pre-aggregated and then turned into actionable QoE measures.

Some of the metrics reported for per-user and aggregate viewing sessions include startup delay, rebuffering delays, average bitrates, and frequency of bitrates switches. Operators might be more inclined to associate viewers' engagement to QoE because better QoEs might make viewers less likely to abandon a viewing session, leading to increased monetization of assets and low subscriber turnover.

The definition of viewer engagement may have different meanings for different operators and contexts. First, operators might like to know which viewer engagement metrics affect QoE the most in order to guide the design of the delivery infrastructures. Secondly, they might also like to quickly identify and resolve service outages and other quality issues. A minute of encoder glitch could replicate throughout the *Internet Service Providers* (ISPs) and the various delivery infrastructures, and affect all their customers. Operators might like to know the scale of this impact, and how it affects users' engagement. The cost of getting the viewer experience wrong often makes the news headlines^[45]. Finally, they would like to understand their customers' demographics (connection methods, type of device, and bitrates of the consumed asset) within a demographic region so that they can fully monetize their assets, and use their other resources strategically.

QoE enthusiasts advocate QoE measurement to be a multidisciplinary approach that seeks to explain its findings, building on general laws of perception, sociology, and user psychology^[39]. If the end-user device analytics is used as a means of QoE measurement, many unexplained variables may not be accounted for (for example, why a user exits a service). For instance, the viewer's tolerance thresholds of QoS parameters for live and VoD might have different QoE patterns, and another dimension of complexity to include might be premium vs. free content. Also, a lack of interest in watching the content, not necessarily an effect of poor QoE, might make a user exit a service. Viewer engagement metrics is measured objectively bypassing subjective studies and surveys.

Some attempts at addressing these unexplained variables can be found in the literature. For example, [46] used the fraction of video viewed as a measure of engagement because this factor can be measured objectively. The data that appeared to belong to early quitters were systematically removed from their analysis to provide a clearer understanding of how the QoE measures affected viewer engagement. A slightly different approach to measuring viewer engagement can be found in [47]. However, with a huge database of aggregated data in the range of millions of viewing session logs, some profile classes of how the QoE measures affect viewer engagement might emerge.

Background on QoE Measurement Methods

The MOS appears to be the *de facto* standard metric for QoE. The possible reasons could be its long-term establishment in telephony networks, perhaps its widespread acceptance on the merits that it can be easily understood, and a metric for benchmarks. There are different types of MOS values, and different test methodologies to produce them, (see [32] for more details). Figure 5 shows the five-point absolute category rating MOS scale that is commonly used.

Figure 5: Five-Point Rating Scale

Score	Label
5	Excellent
4	Good
3	Fair
2	Poor
1	Bad

The MOS value is the average opinion for a given QoS tests condition, not necessarily for the individual users because different users have different opinions. Additional information such as statistical uncertainty in terms of confidence intervals is usually encouraged. The MOS is considered to be characteristic of only the experiment and the group of test subjects from which it was derived.

While there is a reference methodology to produce MOS, it has to be interpreted within context. First, the MOS value obtained for a particular QoS test condition in a subjective experiment may depend on the range of the QoS test conditions used in the experiment. This dependence might be due to test subjects who re-calibrate their use of the rating scale to the conditions in the experiment. An appropriately designed experiment that has a practice period at the start of the experiment, and the test conditions include the best and worst conditions, minimizes the effects of the aforementioned behavior.

Secondly, direct comparisons of MOS scores obtained from separate experiments are generally not meaningful. They are meaningful only if the experiments have been specially designed to enable such comparisons. Data from such specially configured experiments must be studied and shown that their MOS comparisons are statistically valid. Biases in the rating scale interpretation might exist such as differences in interpretation and use of rating scales across cultures; test subject profile, for example, age and technology exposure; test environment; and the presentation order of the test conditions^[48].

Thirdly, it is possible that different objective models that have been trained and optimized using different subjective contexts will predict non-identical MOS values for the same QoS conditions. Objective models are usually developed and optimized for a specific scope of quality features. As a consequence, comparisons between MOS predictions and thresholds can be reliably made only if the thresholds are chosen in the context of the MOS model. See [46] for MOS interpretation and reporting.

The ITU standardization activities^[49] classify objective quality assessment methods into five main categories: media-layer models, parametric packet-layer models, parametric planning models, bitstream-layer models, and hybrid models. Depending on the input parameters to the model and the specific service type (for example, speech, audio, video, and multimedia), each category is aimed at predicting QoE.

The practical deployment of such objective models might benefit in-service monitoring positioned at any one node, or at all nodes, within the content-delivery ecosystem. The node(s) could be at the headend incoming feeds, distribution networks, and at locations of endpoints that are customer equipment. For example, the media-layer model (also known as the perceptual model) might be best suited for quality assurance at the headend, and for benchmark purposes. The media-layer model uses both the source video and the impaired video to predict QoE. The parametric-packet layer model uses only packet-header information such as IP packets, for example, UDP, and extracts the required information to predict QoE. It does not consider encoding distortions, but owing to its light-weight implementation it might be appropriate for in-service monitoring distribution networks. On the other hand, the bitstream-layer models might be appropriate for locations of endpoints that are customer equipment, because they analyze up to the bitstream level. This model in turn does not consider the decoded (impaired) output. The perceptual hybrid model, which combines the media layer and bitstream layer, might use both the bitstream information and the decoded output to predict QoE. While objective assessment seems to offer real-time QoE measurements, some categories might not meet industry demands. Hence, end-user device analytics as a method to QoE measurement appears to be an alternative approach.

Currently, there is the lack of a reference methodology for end-user device analytics as a method of QoE measurement, analogous to MOSs found in subjective assessments^[32] and objective assessments^[42]. A limiting factor to this development might be the restricted rights governing service providers, or the likes, on the usage of such databases. This situation makes it challenging for researchers, service providers, and delivery infrastructures to focus their efforts on developing better delivery infrastructures.

Subjective experiments could still be the most accurate way to measure perceived QoE, and the only way to obtain reliable ground-truth data used in benchmarking objective QoE models.

Linking QoS to QoE

A first glance indicates a considerable mismatch between the concepts of QoS and QoE. This mismatch can be seen in the ITU-T definitions of the two terms. QoS is defined in ITU-T E.800^[50] as the “totality of characteristics of a telecommunications service that bear on its ability to satisfy stated and implied needs of the user of the service.” These characteristics are objective measures consisting of quantitative variables and attributes that may be present or absent. Thus, QoS characteristics are objective and can be objectively measured. By contrast, ITU-T P.10^[51] defines QoE as “the overall acceptability of an application or service, as perceived subjectively by the end-user.” Thus, QoE is called out specifically as consisting of subjective measures.

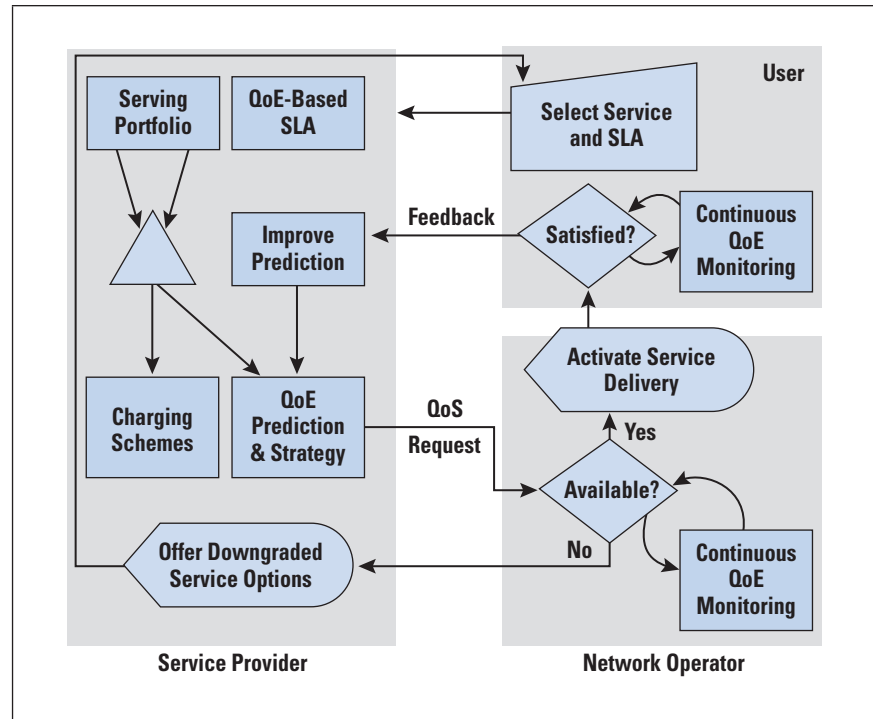
Both network service providers and customers have become accustomed to developing SLAs based on QoS measures. Recent years have seen a growing awareness on the part of both providers and customers that QoE is the more important concept, and that ways must be found to tie network performance parameters, as committed in an SLA, to the desired user QoE. For a typical network and service environment for a particular customer, numerous QoS metrics will impact overall QoE. The focus of ongoing research and product development is on developing reliable techniques, acceptable to both service provider and customer, for correlating QoS performance metrics with QoE as measured by MOS. ITU-T G.1080 identifies two ways in which such correlations can be exploited:

- Given a QoS measurement, predict the expected QoE for a user, with appropriate assumptions.
- Given a target QoE for a user, deduce the required network service performance, with appropriate assumptions.

Service providers can take the first approach and provide a range of QoS offerings with an outline of the QoE that their customers might reasonably expect. Customers can take the second approach by defining the required QoE and then determining what level of service will meet their needs.

Figure 6 illustrates a scenario for the second approach, where the user can make a selection from a range of services, including the required *level of service* (SLA). By contrast to the purely QoS-based management, the SLA here is not expressed in terms of raw network parameters. Instead, the user indicates a QoE target; it is the service provider that maps this QoE target together with the type of service selected, onto QoS demands.

Figure 6: User-Centric Service Delivery



For instance, in the case of multimedia streaming service, the user may simply choose between two QoE levels (high or low). The service provider selects the appropriate quality-prediction model and management strategy (for example, minimize network resource consumption) and forwards a QoS request to the operator. It is possible that the network cannot sustain the required level of QoS, making it impossible to deliver the requested QoE. This situation leads to a signal back to the user, prompting a reduced set of services/QoE values.

Assuming that the network can support the service, delivery can be activated. During service operation, two monitoring and control loops run concurrently: one at network level and the other at service level. The latter allows the user to switch to a different level of QoE (for example, to get a cheaper service or to request higher quality). If the user generates no explicit feedback, it means that the user is satisfied, confirming that the quality-prediction model is working. In this way, the quality-prediction model continues to be redefined during service delivery, allowing it to evolve as user needs and devices change over time.

There are so many varied components with this approach, which extends from the management complexity of a *Content-Delivery Network* (CDN) service based on QoE SLAs to service billing. However, focusing on integrating QoE as part of the management methodology during the design and development of services ensures a user-centric perspective, and helps to move beyond the current network design principles.

References

- [1] Geoff Huston, “QoS—Fact or Fiction?” *The Internet Protocol Journal*, Volume 3, No. 1, March 2000.
- [2] ITU-T, “An Architectural Framework for Support of Quality of Service in Packet Networks,” Recommendation Y.1291, May 2004.
- [3] Sally Floyd and Van Jacobson, “Random Early Detection Gateways for Congestion Avoidance,” *IEEE/ACM Transactions on Networking*, August 1993.
- [4] Jon Crowcroft, Bob Braden, Steve Deering, Sally Floyd, Bruce Davie, Van Jacobson, and Deborah Estrin, “Recommendations on Queue Management and Congestion Avoidance in the Internet,” RFC 2309, April 1998.
- [5] William Stallings, “MPLS,” *The Internet Protocol Journal*, Volume 4, No. 3, September 2001.
- [6] ITU-T, “Traffic Control and Congestion Control in IP-based Networks,” Recommendation Y.1221, June 2010.
- [7] Dimitrios Pendarakis, Raj Yavatkar, and Roch Guerin, “A Framework for Policy-based Admission Control,” RFC 2753, January 2000.
- [8] Eric S. Crawley, Raj Nair, Bala Rajagopalan, and Hal Sandick, “A Framework for QoS-based Routing in the Internet,” RFC 2386, August 1998.
- [9] Lixia Zhang, Steve Berson, Shai Herzog, and Sugih Jamin, “Resource ReSerVation Protocol (RSVP) – Version 1 Functional Specification,” RFC 2205, September 1997.
- [10] Sally Floyd, “TCP and Explicit Congestion Notification,” *ACM Computer Communication Review*, October 1994.
- [11] Vern Paxson, Guy Almes, Jamshid Mahdavi, and Matt Mathis, “Framework for IP Performance Metrics,” RFC 2330, May 1998.
- [12] Sebastian Möller, Patrick Le Callet, and Andrew Perkis, “Qualinet White Paper on Definitions on Quality of Experience – Output Version of the Dagstuhl Seminar 12181,” *European Network on Quality of Experience in Multimedia Systems and Services*, COST Action IC 1003, 2012.
- [13] Florence Agboma and Antonio Liotta, “QoE for Mobile TV Services,” In *Multimedia Transcoding in Mobile and Wireless Networks*, IGI Global, 2009.

- [14] Mario Siller, “An Agent-Based Platform to Map Quality of Service to Experience in Active and Conventional Networks,” University of Essex, 2006.
- [15] Milos Ljubojević et al., “Influence of Resolution and Frame Rate on the Linear in-Stream Video Ad QoE,” *Journal of Information Technology and Applications*, Volume 4, Issue 1, June 2014.
- [16] Jingjing Zhang and Nirwan Ansari, “On Assuring End-to-End QoE in Next Generation Networks: Challenges and a Possible Solution,” *IEEE Communications Magazine*, Volume 49, Issue 7, July 2011.
- [17] Ulrich Reiter et al., “Factors Influencing Quality of Experience,” In *Quality of Experience: Advanced Concepts, Applications and Methods*, 2014.
- [18] Sebastian Möller, M. Waltermann, and M. Garcia, “Features of Quality of Experience,” In *Quality of Experience: Advanced Concepts, Applications and Methods*, 2014.
- [19] Miguel Ríos Quintero and Alexander Raake, “Is taking into account the subjects’ degree of knowledge and expertise enough when rating quality?” In *Fourth International Workshop on Quality of Multimedia Experience (QoMEX)*, July 2012.
- [20] Telchemy, “Voice Quality Measurement,” In Voice over IP Performance Management, 2014. <http://www.telchemy.com/appnotes/TelchemyVoiceQualityMeasurement.pdf>
- [21] Ofcom, *On-Demand Services: Understanding Consumer Choices*, October 2012.
- [22] Andreas Sackl et al., “Wireless Vs. Wireline Shootout: How User Expectations Influence Quality of Experience,” In *Fourth International Workshop on Quality of Multimedia Experience (QoMEX)*, July 2012.
- [23] Florence Agboma and Antonio Liotta, “Addressing User Expectations in Mobile Content Delivery,” *Mobile Information Systems*, Volume 3, Issue 3–4, December 2007.
- [24] David S. Hands, “A Basic Multimedia Quality Model,” *IEEE Transactions on Multimedia*, Volume 6, No. 6, December 2004.
- [25] Stefan Winkler and Christof Faller, “Maximizing Audiovisual Quality at Low Bitrates,” *Proceedings of International Workshop on Video Processing and Quality Metrics for Consumer Electronics*, January 2005.

- [26] Tobias Hossfeld et al., “Initial Delay Vs. Interruptions: Between the Devil and the Deep Blue Sea,” In *Fourth International Workshop on Quality of Multimedia Experience (QoMEX)*, July 2012.
- [27] Sebastian Egger et al., “Waiting Times in Quality of Experience for Web Based Services,” In *Fourth International Workshop on Quality of Multimedia Experience (QoMEX)*, July 2012.
- [28] Peter Reichl, Patrick Maillé, Patrick Zwickl, and Andreas Sackl, “A Fixed-Point Model for QoE-Based Charging,” In *Proceedings of the 2013 ACM SIGCOMM Workshop on Future Human-Centric Multimedia Networking*, 2013.
- [29] Andrew Perkis, Peter Reichl, and Sergio Beker, “Business Perspectives on Quality of Experience,” In *Quality of Experience*, 97–108, 2014.
- [30] Florence Agboma, “Quality of Experience Management in Mobile Content Delivery Systems,” University of Essex, Ph.D. Thesis, 2009.
- [31] Martina Angela Sasse et al., “Remote Seminars through Multimedia Conferencing: Experiences from the Mice Project,” In *INET ’94*, 1994.
- [32] ITU-T, “Methods for the Subjective Assessment of Video Quality, Audio Quality and Audiovisual Quality of Internet Video and Distribution Quality Television in Any Environment,” Recommendation P.913, 2014.
- [33] ITU-R, “Methodology for the Subjective Assessment of the Quality of Television Pictures,” Recommendation BT.500, 2012.
- [34] ITU-T, “Subjective Video Quality Assessment Methods for Multimedia Applications,” Recommendation P.910, 2008.
- [35] ITU-R, “General Methods for the Subjective Assessment of Sound Quality,” Recommendation BS.1284-1, 2003.
- [36] Jesus Gutierrez et al., “Subjective Evaluation of Transmission Errors in IPTV and 3DTV,” In *IEEE Visual Communications and Image Processing (VCIP)*, 2011.
- [37] Nicolas Staelens et al., “Assessing Quality of Experience of IPTV and Video on Demand Services in Real-Life Environments,” *IEEE Transactions on Broadcasting*, Volume 56, Issue 4, December 2010.

- [38] Hendrik Knoche, John D. McCarthy, and Martina Angela Sasse, “Can Small Be Beautiful?: Assessing Image Resolution Requirements for Mobile TV,” In *ACM International Conference on Multimedia*, 2005.
- [39] Tobias Hossfeld et al., “Quality of Experience Management for YouTube: Clouds, Fog and the Aquareyoum,” *PIK – Praxis der Informationsverarbeitung und Kommunikation*, Volume 35, No. 3 2012.
- [40] ITU-R, “Parametric Non-Intrusive Assessment of Audiovisual Media Streaming Quality,” Recommendation P.1201, 2013.
- [41] ITU-T, “Objective Perceptual Multimedia Video Quality Measurement of HDTV for Digital Cable Television in the Presence of a Full Reference,” Recommendation J.341, 2011.
- [42] Marcus Barkowsky et al., “Hybrid video quality prediction: reviewing video quality measurement for widening application scope,” In *Multimedia Tools and Applications*, Volume 74, No. 2, Springer, 2015.
- [43] ITU-T, “Methods, Metrics and Procedures for Statistical Evaluation, Qualification and Comparison of Objective Quality Prediction Models,” Recommendation P.1401, 2012.
- [44] Kjell Brunnström, David Hands, Filippo Speranza, and Arthur Webster, “VQEG Validation and ITU Standardization of Objective Perceptual Video Quality Metrics [Standards in a Nutshell],” *IEEE Signal Processing Magazine*, Volume 26, No. 3, May 2009.
- [45] BBC News, “Xbox and PlayStation resuming service after attack,” <http://www.bbc.co.uk/news/uk-30602609>
- [46] Athula Balachandran, Vyas Sekar, Aditya Akella, Srinivasan Seshan, Ion Stoica, and Hui Zhang “Developing a Predictive Model of Quality of Experience for Internet Video,” In *ACM SIGCOMM Computer Communication Review*, Volume 43, No. 4, October 2013.
- [47] Florin Dobrian, Asad Awan, Dilip Joseph, Aditya Ganjam, Jibin Zhan, Vyas Sekar, Ion Stoica, and Hui Zhang, “Understanding the Impact of Video Quality on User Engagement,” *ACM SIGCOMM Computer Communication Review*, Volume 41, No. 4, August 2011.

- [48] ITU-T, “Mean Opinion Score Interpretation and Reporting,” Recommendation P.800.2, 2013.
- [49] Akira Takahashi et al., “Standardization activities in the ITU for a QoE assessment of IPTV,” *IEEE Communications Magazine*, Volume 46, No. 2, February 2008.
- [50] ITU-T, “Definitions of terms related to quality of service,” Recommendation E.800, September 2008.
- [51] ITU-T, “Vocabulary for performance and quality of service,” Recommendation P.10, July 2006.

WILLIAM STALLINGS is an independent consultant and author of numerous books on security, computer networking, and computer architecture. His latest book is *Foundations of Modern Networking: SDN, NFV, QoE, IoT, and Cloud* (Pearson, 2016). He maintains a resource site for computer science students and professionals at ComputerScienceStudent.com and is on the editorial board of *Cryptologia*. He has a Ph.D. in computer science from M.I.T. He can be reached at ws@shore.net

FLORENCE AGBOMA holds an M.Sc. and a Ph.D. in Electronic Systems Engineering from the University of Essex, United kingdom, in 2005 and 2009, respectively. Her Ph.D. research focused on Quality of Experience for Mobile Content Delivery Systems. She currently works as a Technology Analyst at British Sky Broadcasting (BSkyB), London. Her interests include subjective video quality assessment, psychophysical methods, and quality of experience management across devices and platforms. E-mail: fagbom@googlemail.com

IANA Stewardship Transition Moves to Final Phase

An historic proposal for the global community to assume stewardship of the *Internet Assigned Numbers Authority* (IANA) functions, produced after nearly two years of work by the global Internet community, has been delivered to the U.S. Government for its consideration. The proposal would remove U.S. Government oversight over a set of fundamental Internet administrative functions, including management of the global pool of Internet number resources (IPv4 and IPv6 addresses and *Autonomous System Numbers*), and replace it with a set of arrangements for community-based oversight.

The proposal, developed by the *IANA Stewardship Transition Coordination Group* (ICG), is based on input from three operational communities, including the *Internet Number Community* (those with an interest in the global management of Internet number resources). The contributions of the Internet Number Community were coordinated via a *Consolidated RIR IANA Stewardship Proposal* (CRISP) Team made up of community members drawn from each of the five RIR regions.

While the ICG published the final draft of its proposal in October 2015, elements of the proposal relied upon the adoption of a set of recommendations regarding the accountability of ICANN to its community. These recommendations were developed separately by a *Cross Community Working Group on Enhancing ICANN Accountability* (CCWG) and were adopted by the ICANN Board at its meeting in March 2016 in Marrakech, Morocco. The Board was at that point able to pass on both the ICG and CCWG documents to the *National Telecommunication and Information Administration* (NTIA), an agency of the U.S. Government.

The U.S. Government will now review the proposal to ensure that it meets the criteria set out by the NTIA when they first announced their intention, in March 2014, to pass stewardship of the IANA functions to the global community. If approved, the *Regional Internet Registries* (RIRs) and ICANN will continue their work towards implementation of the proposal, which will be completed prior to the expiration of ICANN's current contract with NTIA in September 2016.

For more information, see:

<http://www.ianacg.org/>

<https://www.nro.net/nro-and-internet-governance/iana-oversight/about-the-proposal>

<http://www.ianacg.org/icg-files/documents/IANA-transition-proposal-final.pdf>

<https://community.icann.org/display/acctcrosscomm/CCWG+on+Enhancing+ICANN+Accountability>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Publication of this journal is made possible by:

Supporters



Diamond Sponsors



Ruby Sponsor



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



Individual Sponsors

Lyman Chapin, Steve Corbató, Dave Crocker, Jay Etchings, Martin Hannigan, Hagen Hultzs, Dennis Jennings, Jim Johnston, Merike Kaeo, Bobby Krupczak, Richard Lamb, Tracy LaQuey Parker, Bill Manning, Andrea Montefusco, Tariq Mustafa, Mike O'Connor, Tim Pozar, George Sadowsky, Scott Seifel, Helge Skrivervik, Rob Thomas, Tom Vest, Rick Wesson.

For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

ADDRESS SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Fred Baker, Cisco Fellow
Cisco Systems, Inc.

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Steve Crocker, Chairman
Internet Corporation for Assigned Names and Numbers

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

June 2016

Volume 19, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Fragmentation	2
Resource Discovery in IoT....	13
The IANA Transition.....	26
Fragments	29
Call for Papers.....	30
Supporters and Sponsors	31

FROM THE EDITOR

A major design feature of the *Internet Protocol* (IP) is its ability to run over a variety of underlying network technologies. If you look through the *Request For Comments* (RFC) document series, you will find numerous specifications of the form “IP over xxx,” where “xxx” is anything from Ethernet to X.25, Frame Relay, Bluetooth, WiFi, and even “Avian Carriers” (pigeons), the latter being one of the more famous April Fools RFCs. Because each of these technologies has different capabilities in terms of how much data can be carried in a “packet” or datagram, IP employs the concept of *fragmentation* and *reassembly* in cases where the originating datagram is larger than what the underlying network medium can support. In our first article, Geoff Huston explains fragmentation and reassembly for both IPv4 and IPv6. Special thanks go to Mansour Ganji of Vodafone New Zealand for suggesting this topic.

We’ve covered various aspects of the *Internet of Things* (IoT) in previous editions of this journal. This time, Akbar Rahman and Chonggang Wang discuss ongoing work within the *Internet Engineering Task Force* (IETF) and elsewhere to develop Resource Discovery mechanisms for IoT devices.

The long-awaited proposal to transition the *Internet Assigned Numbers Authority* (IANA) Stewardship Functions from the U.S. Government to a new entity was finally submitted in early March of this year. Vint Cerf explains the history and background of this process. At the end of his article you will find pointers to further information about this important Internet milestone.

As always, we welcome your feedback, suggestions, book reviews, articles, and sponsorship support. You can contact us by e-mail to ipj@protocoljournal.org and visit our website for subscription information, back issues, author guidelines, sponsor information, and much more.

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

Fragmentation

by Geoff Huston, APNIC

One of the more difficult design exercises in packet-switched network architectures is that of the design of packet fragmentation. In this article I will examine *Internet Protocol* (IP) *packet fragmentation* in detail and look at the design choices made by IP Version 4, and then compare that with the design choices made by IP Version 6.

Packet-switched networks dispensed with a constant time base, in turn allowing individual packets to be sized according to the needs of the application as well as the needs of the network. Smaller packets have a higher ratio of packet header to payload, and are consequently less efficient in data carriage. On the other hand, within a packet-switching system the smaller packet can be dispatched faster, reducing the level of *head-of-line blocking* in the internal queues within a packet switch and potentially reducing network-imposed jitter as a result. Larger packets allow larger data payloads, in turn allowing greater carriage efficiency. Larger payloads per packet also allows a higher internal switch capacity when measured in terms of data throughput. But larger packets take longer to be dispatched, potentially causing increased jitter.

Various network designs have adopted various parameters for packet size. Ethernet, standardized in the mid-1970s, adopted a variable packet size, with supported packet sizes of between 64 and 1,500 octets. *Fiber Distributed Data Interface* (FDDI), a fibre ring local network, used a packet size of up to 4,478 octets. Frame Relay used a variable packet size of between 46 and 4,470 octets. The choice of variable-sized packets allows applications to refine their behaviour. Jitter and delay-sensitive applications, such as digitised voice, may prefer to use a stream of smaller packets to attempt to minimise jitter, while reliable bulk data transfer may choose a larger packet size to increase carriage efficiency. The nature of the medium may also have a bearing on this choice. If there is a high *Bit Error Rate* (BER) probability, then reducing the packet size minimises the impact of sporadic errors within the data stream, possibly increasing throughput.

IPv4 and Packet Fragmentation

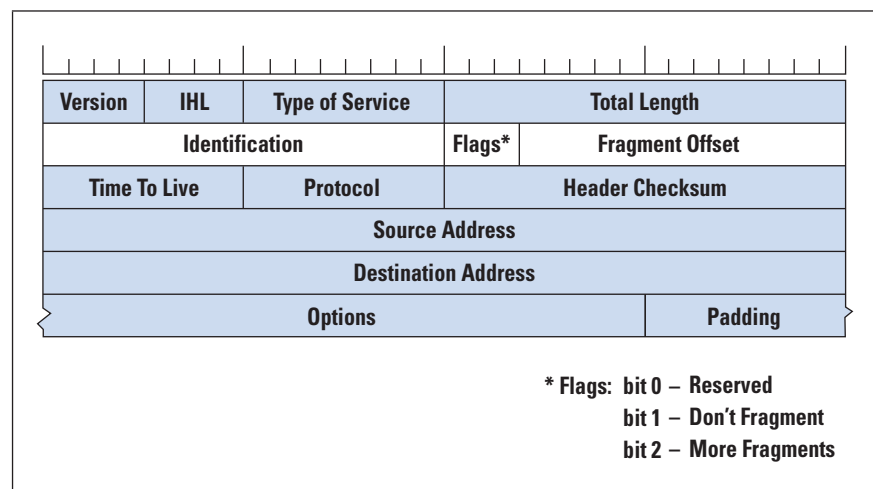
In designing a network protocol that is intended to operate over a wide variety of substrate carriage networks, the designers of IP could not rely on a single packet size for all transmissions. Instead the IP designers of the day provided a packet-length field in the IP Version 4 header^[1]. This field was a 16-bit octet count, allowing for an IP packet to be anywhere from the minimum size of 20 octets (corresponding to an IP header without any payload) to a maximum of 65,535 octets. So IP itself supports a variable size packet format. But which packet size should an implementation use?

The tempting answer is to use the maximum size permitted by the network interface of the local device, with the caveat that an application may nominate the explicit use of smaller-sized packets. But there is a complication here. The Internet was designed as an “inter-network” network system, allowing an IP packet to undertake an end-to-end journey from source to destination across numerous different networks. For example, consider a host connected to a FDDI network, which is connected to an Ethernet network. The FDDI-connected host may elect to send a 4,478-octet packet, which will fit into a FDDI network, but the packet switch that attempts to pass the packet into the Ethernet network will be unable to do so because it is too large.

The solution adopted by IPv4 was *forward fragmentation*. The basic approach is that any IP router that is unable to forward an IP packet into the next network because the packet is too large for this network may split the packet into a set of smaller IP fragments, and forward each of these fragments. The fragments continue along the network path as autonomous packets, and the addressed destination host is responsible for reassembling these fragments back into the original IP packet.

The behaviour is managed by a 32-bit field in the IPv4 header, which is subdivided into three sub-fields (Figure 1).

Figure 1: IPv4 Packet Header
Fragmentation Fields



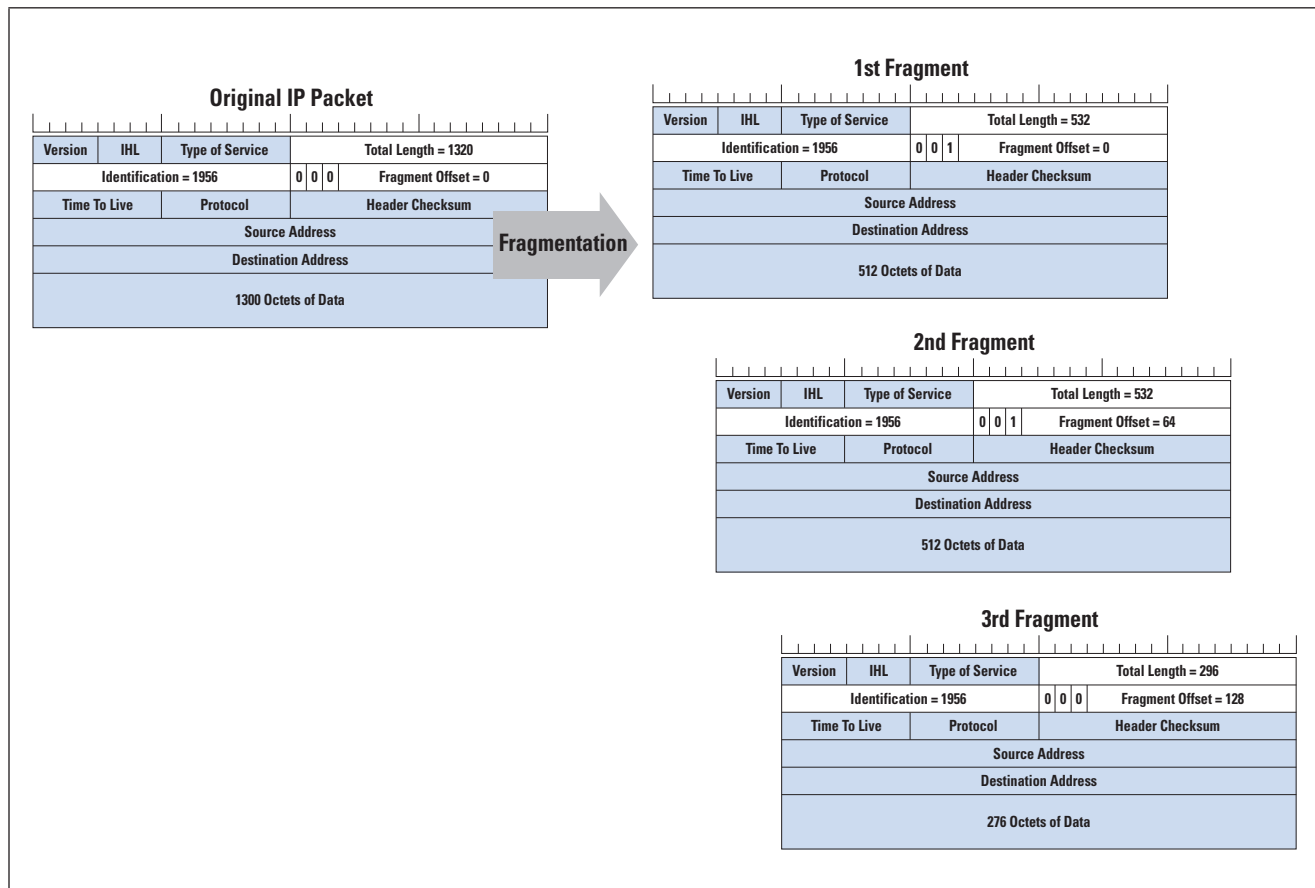
The first sub-field is a 16-bit packet identifier which allows fragments that share a common packet-identifier value to be identified as fragments of the same original packet.

The second sub-field is a 3-bit vector of flags. The first bit is unused. The second is the *Don't Fragment* flag. If this flag is set the packet cannot be fragmented, and must be discarded when it cannot be forwarded. The third bit is the *More Fragments* field, and is set for all fragments except the final fragment.

The third sub-field is the fragmentation offset value that is the offset of this fragment from the start of the IP payload of the original packet, measured in *octawords* (64-bit units).

For example, a router attempting to pass a 1320-octet IP packet into a network whose maximum packet size is 532 octets would need to split the IP packet into three parts. The first packet would have a fragmentation offset of 0 and the *More Fragments* bit set. The total length would be 532 octets, and the IP payload would be 512 octets, making a total of 532 octets for the packet. The second packet would have a fragmentation offset value of 64, the *More Fragments* bit set, total length of 532, and an IP payload of 512 octets, making a total of 532 octets for the packet. The third packet would have a fragmentation offset value of 128, the *More Fragments* bit clear, total length of 296, and an IP payload of 276 octets, making a total of 296 octets for the packet (Figure 2).

Figure 2: Example of IPv4 Packet Fragmentation



The advantage of this approach is that as long as it is permissible to fragment the IP packet, all packet flows are “forward,” meaning that the sending host is unaware that packet fragmentation is occurring, and all the IP fragment packets continue to head towards the original destination, where they are reassembled.

Another advantage is that while the router performing the fragmentation has to expend resources to generate the packet fragments, the ensuing routers on the path to the destination have no additional processing overhead, assuming that they do not need to further fragment these IP fragments. Fragments can be delivered in any order, so the fragments may be passed along parallel paths to the destination.

To complete the IPv4 story we must describe the IPv4 behaviour when the *Don't Fragment* bit is set. The router that is attempting to fragment such a packet is forced to discard it. Under these circumstances the router is expected to generate an *Internet Control Message Protocol* (ICMP) “Unreachable” error (type 3, code 4), and in later versions of the IP specification it was expected to add the *Maximum Transmission Unit* (MTU) of the next-hop network into the ICMP packet. The original sender would react to receiving such an ICMP message by changing its local maximum packet size associated with that particular destination address, and thus it would “learn” a viable packet size for the path between the source and destination.

Evaluating IPv4 Fragmentation

A case has been made that the IP approach to fragmentation contributed to its success. This design allowed transport protocols to operate without consideration of the exact nature of the underlying transmission networks, and avoid additional protocol overhead in negotiating an optimal packet size for each transaction. Large *User Datagram Protocol* (UDP) packets could be transmitted and fragmented in real time as required without requiring any form of end-to-end network path packet size discovery. This approach allowed IP to be used on a wide variety of substrate networks without requiring extensive tailoring.

But it wasn't all good news.

Cracks in the IP fragmentation story were described in a 1987 paper by Kent and Mogul, “Fragmentation Considered Harmful.”^[2]

TCP has always attempted to avoid IP fragmentation. The initial opening handshake of *Transmission Control Protocol* (TCP) exchanges the local and remote *Maximum Segment Size* (MSS), and the sender will not send a TCP segment larger than that notified by the remote end at the start of the TCP session. The reason that TCP attempted to avoid fragmentation was that fragmentation was inefficient under conditions of packet loss in a TCP environment. Lost fragments can be repaired only by resending the entire packet, including resending all those fragments that were successfully transmitted in the first place. TCP can perform a data repair more efficiently if it limits its packet size to one that does not entail packet fragmentation.

This form of fragmentation also posed vulnerabilities for hosts. For example, an attacker could send a stream of fragments with a close to maximally sized fragment offset value, and random packet identifier values.

If the receiving host believed that the fragments represented genuine incoming packets, then a credulous implementation might generate a reassembly buffer for each received fragment that may represent a memory buffer starvation attack. It is also possible, either through malicious attack or by poor network operation, that fragments may overlap or overrun, and the task of reassembly requires care and attention in implementation of fragment reassembly.

Lost fragments represent a slightly more involved problem than lost packets. The receiver has a packet reassembly timer upon the receipt of the first fragment, and will continue to hold this reassembly state for the reassembly time. The reassembly timer is a factor in the maximal count of packets in flight, because the packet identifier cannot be recycled within a period defined by the sender-receiver path delay plus the reassembly timer of the receiver. For higher-delay high-capacity network paths, this limit of 65,535 packets in flight can be a potential performance bottleneck^[3].

Fragmentation also consumes router processing time, forcing the processing of oversized packets from a highly optimised fast path into a processor queue.

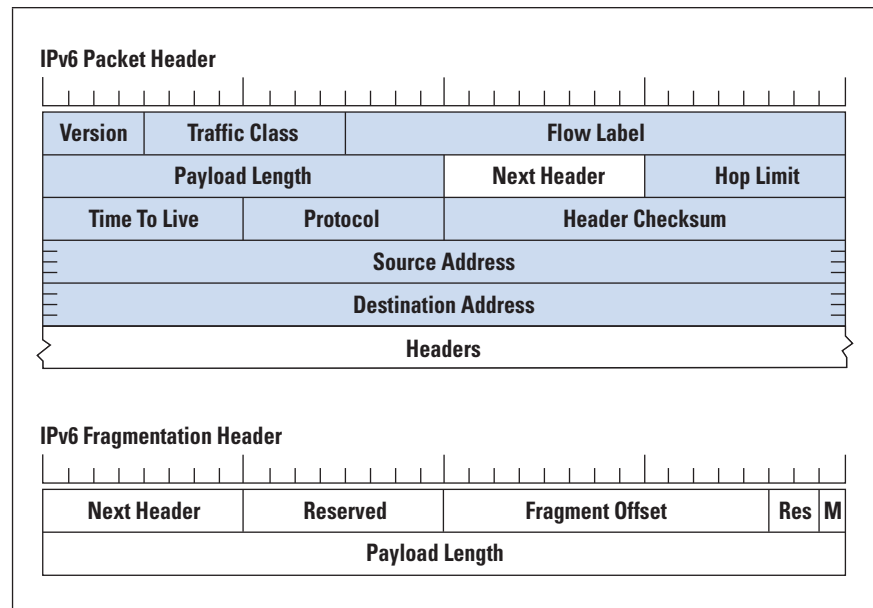
And then there is the “middleware problem.” Filters and firewalls perform their function by applying a set of policy rules to the packet stream. But these rules typically require the presence of the transport layer header. How can a firewall handle a fragment? One option is to pass all trailing fragments through without inspection, but this process exposes the internal systems to potential attack[4]. Another option is to have the firewall rebuild the original packet, apply the filter rules, and then refragment the packet and forward it on if the packet is accepted by the filter rules. However, this process now exposes the firewall to various forms of memory starvation attack. *Network Address Translators* (NATs)[5] that use the transport-level port addresses as part of the NAT binding table have a similar problem with trailing fragments. The conservative approach is for the NAT to reassemble the IP packet at the NAT, apply the NAT transform, and then pass the packet onward, fragmenting as required.

IPv6 and Fragmentation

When it came time to think about the design of what was to become IPv6, the *forward fragmentation* approach was considered to be a liability, and while it was not possible to completely discard IP packet fragmentation in IPv6, there was a strong desire to redefine its behaviour.

The essential change between IPv4 and IPv6 is that in IPv6 the *Don't Fragment* bit is always on, and because it's always on, it's not explicitly contained in the IPv6 packet header (Figure 3). There is only one fragmentation flag in the Fragmentation Header, the *More Fragments* bit, and the other two bits are reserved. The other change was that the packet-identifier size was doubled in IPv6, using a 32-bit packet identifier field.

Figure 3: IPv6 Packet Header and Fragmentation Header



An IPv6 router cannot fragment an IPv6 packet, so if the packet is too large for the next hop the router is required to generate an ICMPv6 type 2 packet, addressed to the source of the packet with a *Packet Too Big* (PTB) code, and also providing the MTU size of the next hop. While an IPv6 router cannot perform packet fragmentation, the IPv6 sender may fragment an IPv6 packet at the source.

Evaluating IPv6 Packet Fragmentation

The hope was that these IPv6 changes would fix the problems seen with IPv4 and fragmentation.

Our experience appears to point to a different conclusion.

The first problem is that there is widespread ICMP packet filtering in today's Internet. For IPv4 this approach was basically a reasonable defense tactic, and if you were willing to have a packet fragmented you cleared the *Don't Fragment* bit before sending the packet so that you didn't rely on receiving an ICMP message to indicate a path sender MTU problem. But in IPv6 the equivalent *Don't Fragment* bit function is jammed in the "on" position, and fragmentation can be performed only if the original sender receives the ICMPv6 PTB message and then resends the packet fragmented into a size that meets the specified MTU size. But when ICMPv6 PTB messages are filtered, the large packet is silently discarded within the network without any discernible trace. Attempts by the sender to time out and resend the large IPv6 packet will meet with the same fate, so this situation can lead to a wedged state.

This scenario has been seen in the context of the HTTP protocol, where the path MTU is smaller than the MTU of the host systems at either end. The TCP handshake completes because none of the opening packets is large. The opening HTTP GET packet also makes it through because this packet is normally not a large one.

However, the first response may be a large packet. If it is silently discarded because of the combination of fragmentation required and ICMPv6 filtering, then neither the client nor the server can repair the situation. The connection hangs.

The second problem is that the ICMPv6 PTB message is sent backwards to the source from the interior of a network path. Oddly enough, the IPv6 ICMP PTB message is perhaps the one critical instance in the entire IP architecture in which the IP source address is interpreted by anything other than the intended destination. The problems here include path asymmetry, in that the source address may be unreachable from the point of the generation of the ICMP packet. There is also the case of tunneling IP-in-IP. Because IPv6 fragmentation can be performed only at the source, should the ICMP message be sent to the tunnel ingress point or to the original source? Using the tunnel ingress assumes that the tunnel egress performs packet reassembly, potentially burdening the tunnel egress. This situation is further confounded in the cross protocol case of IPv6-in-IPv4 and IPv4-in-IPv6.^[6]

The third problem is the combination of IPv6 packet fragmentation and UDP. UDP is an unreliable datagram delivery service, so a sender of a UDP packet is not expected to cache the packet and be prepared to resend it. A UDP packet-delivery error can occur only at the level of the application, not at the IP or UDP protocol level. So what should a host do upon receipt of an ICMP PTB message if resending the IP packet is not an option? Given that the sender does not cache sent UDP packets, the packet header in the ICMPv6 message is unhelpful. Because the original packet was UDP, the sender does not necessarily have a connection state, so it is not clear how this information should be retained and how and when it should be used. How can a receiver even tell if an ICMPv6 PTB packet is genuine? If the sender adds an entry into its local IPv6 forwarding table, it is exposing itself to a potential resource starvation problem. A high volume flow of synthetic PTB messages has the potential to bloat the local IPv6 forwarding table. If the sender ignores the PTB message, the application is left to attempt to recover the transaction.

If it makes little sense in the context of an attempt to fragment a UDP packet, it makes less sense to fragment a TCP packet. In the context of a TCP session, a received ICMPv6 PTB message can be interpreted as a redefinition of the remote end MSS value, and the outgoing TCP segments can be reframed to conform to this MSS.

Wither Fragmentation?

The basic problem here is that the network was supposed to operate at the IP level and be completely unaware of transport, implying that IP-level fragmentation was meant to work in a manner that does not involve transport protocol interaction.

So much of today's network (firewalls, filters, etc.) is transport-aware and the trailing fragments have no transport context, meaning that transport-aware network middleware needs to reassemble the packet, and this process could represent a problem and a *Denial of Service* (DoS) vulnerability in its own right.

So is fragmentation worth it at all?

I'd still say that it's more useful to have it than not. But the IPv4 model of *forward fragmentation* in real time has proved to be more robust than the IPv6 model because the IPv4 model requires only that traffic flows in one direction and is an IP-level function. It has its problems, and no doubt the papers that warned that IP fragmentation was "harmful" were sincere in taking that view^[2]. But it is possible to make it worse, and the IPv6 model requiring a *backward* ICMPv6 message from the interior of the network was in retrospect a decision that did just that!

So what should we do now?

It is probably not a realistic option to try to alter the way that IPv6 manages fragmentation. There was an effort in 2013 in one of the IETF's IPv6 Working Groups to deprecate the IPv6 Fragment Header^[7]. That's possibly an overreaction to the problem of packet fragmentation and IPv6, but there is no doubt that the upper-level protocols simply should not assume that IPv6 fragmentation operates in the same manner as IPv4, or even operates in a reliable manner at all!

The implication is that transport protocol implementations, and even applications, should try to manage their behaviour on the assumption that ICMP message filtering is sufficiently prevalent that it is prudent to assume that all ICMP messages are dropped. The result is a default assumption that large IPv6 packets that require fragmentation are silently dropped.

How can we work around this problem and operate a network that uses variable-sized packets but cannot directly signal when a packet is too large? RFC 4281^[8] describes a *Path MTU Discovery* process that operates without relying on ICMP messages, and IPv6 TCP implementations should rely on this mechanism to establish and maintain a viable MTU size that can support packet delivery. In this way TCP can manage the path MTU and the application layer need not add explicit functions to manage persistent silent drop of large segments.

Path MTU Discovery

Path MTU discovery was specified in RFC 1191^[9]. The approach was to send packets with the *Don't Fragment* bit set. When a router on the path is unable to forward the packet because it is too large for the next hop, the *Don't Fragment* field directs the router to discard the packet and send a *Destination Unreachable* ICMP message with a code of "Fragmentation Required and DF set" (type 3, code 4).

RFC 1191 advocated the inclusion of the MTU of the next-hop network in the next field of the ICMP message.

A host receiving this form of ICMP message should store the new MTU in the local forwarding table, with an associated time to allow the entry to time out. Also the host should identify all active TCP sessions that are connected to the same destination address as given in the IP packet header fragment of the ICMP message, and notify the TCP session of the revised path MTU value.

RFC 1981^[10] defined much the same behaviour for IPv6, relying on the MTU information conveyed in the ICMPv6 PTB message in exactly the same manner as its IPv4 counterpart.

The problem of filtered ICMP messages is a difficult one, and attention has turned to path MTU Discovery ideas that do not rely on an ICMP message to operate correctly. RFC 4821 describes a mechanism that refines the RFC 1191 ICMP-based process by adding an alternate process that is based on detection and reporting of packet loss as an inference of path MTU problems when there is no ICMP feedback. This process uses a probe procedure that attempts to establish a working MTU size through probing the path with various sized packets to establish the upper-bound MTU. The trade-off here is the number of round-trip intervals taken to perform the probes and the accuracy of the path MTU estimate.

Because these probes take time, the entire exercise tends to be of value only in long-held TCP and TCP-like flows. For shorter sessions the pragmatic advice is to clamp the local MTU to a conservative value (1,280 is a good first choice for IPv6, and RFC 4821 also suggests 1,024 for IPv4) and try to avoid the entire issue of fragmentation in the first place.

UDP is a different story. The lightweight UDP protocol shim does not admit much in the way of additional functions, and one possible approach is to insist that UDP-based applications limit themselves to the local MTU size, or to be even more conservative, limit themselves to the 1,280-octet IPv6 minimum unfragmented packet size.

The major issue with such advice for UDP lies in the *Domain Name System* (DNS). Efforts to improve the security of the DNS with *Domain Name System Security Extensions* (DNSSEC) have added additional data into DNS responses. In addition, if you want to maintain the lightweight efficiency of the DNS, then it's not possible to keep DNSSEC responses under 1500 octets all the time, let alone under 1,280 octets. One option here is to insist that larger DNS responses use TCP, but this option imposes some considerable cost overhead on the operation of the DNS. What the DNS has chosen to do appears to represent a reasonable compromise.

The first part of the approach is that the management of the packet MTU is passed into the application layer. The application conventionally operates with a maximum UDP payload size that assumes that UDP fragmentation is working, and a DNS query normally offers an *Extension Mechanisms for DNS* (EDNS) buffer size of 4,096 octets. The responder uses this information to assemble its UDP response of up to 4,096 octets in length, a process that conventionally causes the source to perform UDP packet fragmentation for large responses. This fragmented response may not reach the querier for a variety of reasons, in which case the EDNS buffer size is dropped back to a more conservative value that is not expected to trigger UDP fragmentation, but may not be able to contain the complete response. The intended result is that if the network cannot complete a UDP transaction that entails a fragmented UDP response, the transaction is repeated using a smaller maximum UDP packet size, and the truncated response explicitly signals to the client to retry the query using TCP^[12]. This process is protocol-agnostic, in that it operates as intended in the case of IPv4 forward fragmentation, where trailing fragments are filtered out by middleware, and in the case of IPv6, where there is no forward fragmentation, and it operates whether or not the responder receives any ICMP PTB messages.

Conclusion

What we have learned through all this discussion is that packet fragmentation is extremely challenging, and is sensibly avoided if at all possible.

Rather than trying to bury packet fragmentation to an IP-level function performed invisibly at the lower levels of the protocol stack, a robust approach to packet fragmentation requires a more careful approach that lifts the management of Path MTU into the end-to-end transport protocol and even into the application.

IPv6 UDP-based applications that want a lightweight operation should look at keeping their UDP packets under the IPv6 1,280-octet unfragmented packet limit. And if that's not possible, then the application itself needs to explicitly manage Path MTU, and not rely on the lower levels of the protocol stack to manage it.

IPv6 TCP implementations should never assume that IPv6 PTB messages are reliably delivered. High-volume flows should use RFC 4821 Path MTU Discovery and management procedures to ensure that the TCP session can avoid Path MTU blackholing. For short flows, MSS clamping still represents the most viable approach.

I'm not sure that we should go as far as deprecating IP fragmentation in IPv6. The situation is not that dire. But we should treat Path MTU with a lot more respect, and include explicit consideration of the trade-offs between lightweight design and robust behaviour in today's network.

References

- [1] Jon Postel, “Internet Protocol,” RFC 791, September 1981.
- [2] Kent, C. and J. Mogul, “Fragmentation Considered Harmful,” Proc. SIGCOMM ’87 Workshop on Frontiers in Computer Communications Technology, August 1987.
- [3] Matt Mathis, Ben Chandler, and John W. Heffner, “IPv4 Reassembly Errors at High Data Rates,” RFC 4963, July 2007.
- [4] G. Ziemba, D. Reed, and P. Traina, “Security Considerations for IP Fragment Filtering,” RFC 1858, October 1995.
- [5] Geoff Huston, “Anatomy: A Look Inside Network Address Translators,” *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [6] Pekka Savola, “MTU and Fragmentation Issues with In-the-Network Tunneling,” RFC 4459, April 2006.
- [7] Bonica, R. W. Kumari, R. Bush, and H. Pfeifer, “IPv6 Fragment Header Deprecated,” Internet Draft, work in progress, **draft-bonica-6man-frag-deprecate**, July 2013.
- [8] Matt Mathis and John W. Heffner, “Packetization Layer Path MTU Discovery,” RFC 4821, March 2007.
- [9] Jeffrey C. Mogul and Stephen E. Deering, “Path MTU Discovery,” RFC 1191, November 1990.
- [10] Stephen E. Deering, Jack McCann, and Jeffrey Mogul, “Path MTU Discovery for IP Version 6,” RFC 1981, August 1996.
- [11] Mukesh Gupta, Stephen E. Deering, and Alex Conta, “Internet Control Message Protocol (ICMPv6) for the Internet Protocol Version 6 (IPv6) Specification,” RFC 4443, March 2006.
- [12] Paul Vixie, Joao Damas, and Michael Graff, “Extension Mechanisms for DNS (EDNS(0)),” RFC 6891, April 2013.
- [13] Godred Fairhurst and Lars Eggert, “Unicast UDP Usage Guidelines for Application Designers,” RFC 5405, November 2008.

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic

Resource Discovery in the Internet of Things

by Akbar Rahman and Chonggang Wang,
InterDigital Communications, Inc.

The *World Wide Web* (WWW or Web) is a global collection of connected documents and other resources that reside on the Internet. The introduction of the *Internet of Things* (IoT) is expected to dramatically increase the size of the Web in the near future and thus necessitates a fundamental change to the existing mechanisms of discovering resources. In IoT, the vision is that a significant number of new types of devices (or “things”) such as fridges, car sensors, traffic lights, and so on will be dynamically connected to the Web for communication and control. These IoT devices will have radically different characteristics from existing Web servers and users. This article looks at a key protocol development occurring in the *Internet Engineering Task Force* (IETF) for allowing IoT devices to discover resources via a new logical node called a *Resource Directory* (RD).

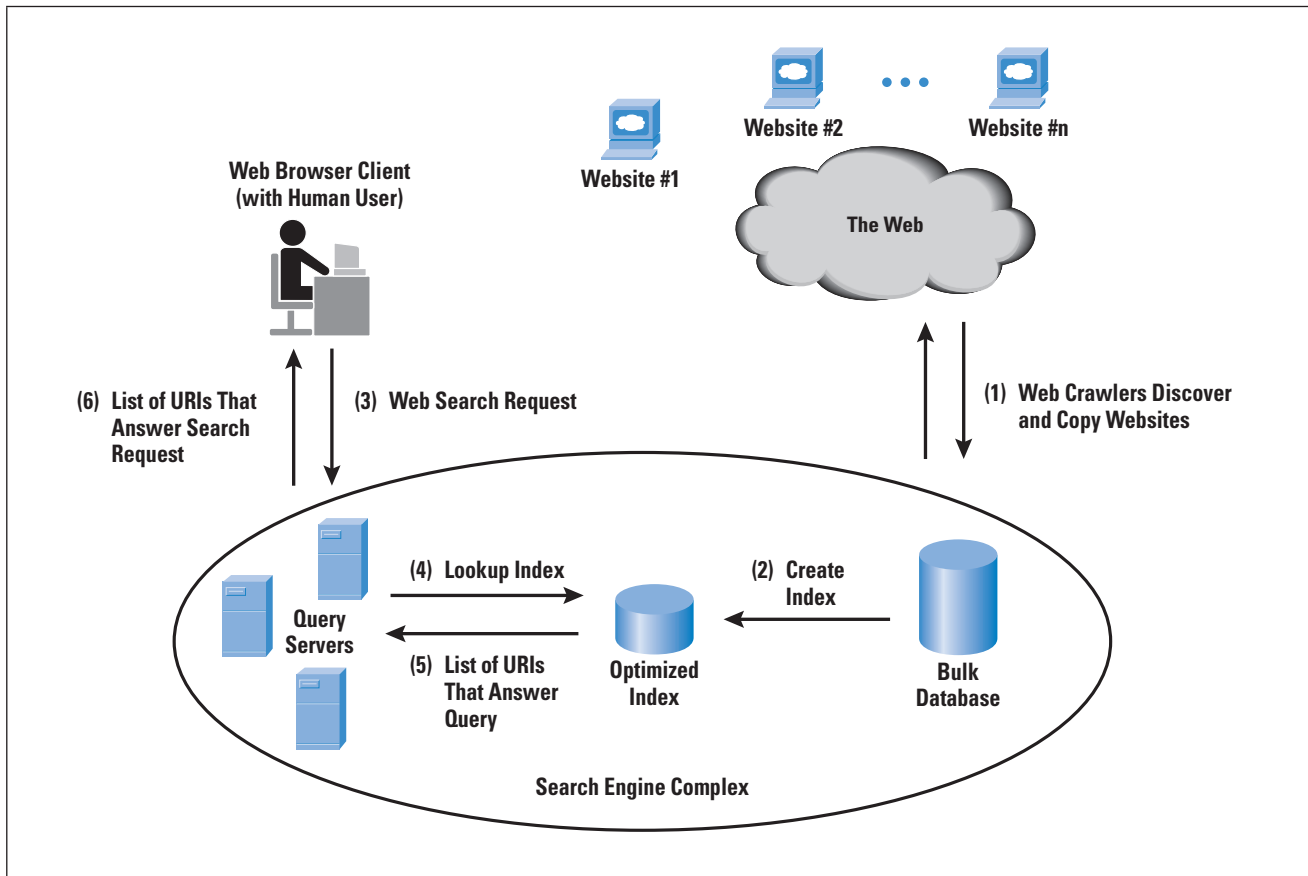
Resource Discovery in the Traditional Web

The basic unit of addressing on the Web is the *Uniform Resource Identifier* (URI), which identifies a resource^[1]. The resource may, for example, be a restaurant-review website page for a human user to read. Or in a more abstract form, the resource may be a software process to be triggered by a *Business-to-Business* (B2B) Web application as part of an automated stock market trading system. The key challenge in all cases is how users can quickly find the correct URI for the resource that they are interested in out of all possible URIs in the entire Web space. This process is referred to as *resource discovery*.

The most well-known and powerful resource discovery mechanism in the current Web is the one employed by Web search engines such as *Baidu*, *Bing*, *Google*, *Yahoo*, etc. Specifically, search engines use the mechanism of Web *crawlers* (also called spiders, ants, or robots) to periodically browse the Web to create a dynamic index of the resources of most publicly available websites. A website is defined as a server that hosts resources users can access with the *Hypertext Transfer Protocol* (HTTP)^[2]. Human users can then send a search request, via a Web browser client, to look up the specific resources that they are interested in.

Figure 1 shows the overall resource discovery process based on Web crawlers. Figure 1 is given in the context of a search engine, but academic researchers, market research companies, and others follow very similar processes. However, unlike a search engine, these other entities typically do not send crawlers to cover the entire Web to discover all possible resources. Instead, they send crawlers to cover parts of the Web to discover the specific type of resource that they are interested in.

Figure 1: Overview of Traditional Resource Discovery Process by Web Search Engines



For example, a market research company may send its Web crawlers to discover all the resources related to a specific type of product in a given geographic area as part of a pricing comparison study.

In Figure 1, Web crawlers start crawling out from the search engine server to an initially provisioned seed list of URIs. This seed list typically consists of very popular websites with a lot of URIs to other sites (that is, *hyperlinks*). From these initial websites, the Web crawlers then crawl outward to all connected hyperlinks. At each new website that it discovers, the Web crawler creates a copy of the website, which it sends back to the search engine^[3]. The search engine records all the received information in a bulk database and later processes it to create an optimized index for fast lookups. Then when a given search request comes from a Web browser client looking for some specific resource, the search engine can go quickly through its index using its own proprietary algorithm to find one or more matches.

Finally, the search engine will return to the client a list of URIs and selected application content pertaining to the resources that match the client's search parameters. This information is then displayed on the user's Web browser interface. Human users will then select ("click") the URI(s) that they want to visit.

Following are some key observations about resource discovery in the traditional Web:

- In terms of network configuration, the search engine functionality, or whichever entity dispatches the Web crawler, is typically located on a set of centralized servers and related databases with high-speed and large-bandwidth Internet connectivity. The resources that the Web crawlers discover may be widely distributed across the entire Internet. The Web browser-based clients that interface with the human users and send the search (lookup) requests are typically located at the edge of the network.
- Search engines primarily use a *pull model* to get resource information. In this approach, the receiving node (that is, search engine) goes out and explicitly requests information (via Web crawlers) from the sending node (content websites). However, a small number of URIs such as the initial seed list of URIs (for example, very popular websites) may be obtained without using the pull model, but these URIs are always a small fraction of the URIs in a search engine index.
- The list of resources the search engine returns for a given Web search (lookup) request may vary from a few URIs to potentially hundreds or even thousands of URIs. The order that these URIs are presented to the human user via the search engine Web interface is called the *ranking* of the resources. This ranking is critical because when a large number of URIs are returned to users for a given search request, users will typically select (“click”) only the top few ranked URIs.
- The ranking of resources is ultimately an algorithmic decision internal to the search engine. However, it can be affected by external input such as *Search Engine Optimization* (SEO) techniques that website developers use to try to get search engines to rank their specific URIs higher than other URIs with similar application content. For example, a simple SEO technique is to have website content clearly tagged (titles, section headings, etc.) and correlated to the website metadata. This metadata is an important input for the search index engine. A more sophisticated SEO technique is to have hyperlinks to a given website from as many other websites as possible because search engines consider this factor a measure of content popularity. There are many other SEO techniques^[4].

The Resource Discovery Problem in IoT

As mentioned previously, a key characteristic of current Web discovery technology is the use of Web crawlers to fan out and discover resources across the Internet. The implicit assumption in this approach is that Web servers are always active and available for Web crawlers that arrive in an unscheduled manner to discover easily. However, this assumption conflicts with the expected nature of many IoT devices that may have only intermittent connectivity to the Web.

The primary reason for this intermittent connectivity is that many IoT devices have a limited power supply (for instance, battery or solar power). To conserve their power they may “wake up” or become active only when required to perform a specific function. For example, a fire-detection sensor acting as a mini Web server may wake up and connect to the Web only to send a warning message to a remote controller when it senses a certain amount of smoke in its vicinity. At most other times, the fire-detection sensor is “asleep” (that is, in a low power state and not active) and unreachable via the Web. A secondary reason for intermittent connectivity is that many IoT devices are connected to the Web by low-power and lossy wireless networks. These wireless networks are more susceptible to interference and temporary loss of connectivity than traditional wired or cellular networks^[5].

Another key difference between IoT devices and other Web infrastructure is that most IoT devices may be deployed in semi-closed networks. For example, the IoT devices such as a lighting or heating control system in a home may have Internet connectivity only through a fire-walled home gateway. So the IoT devices and their associated resources may be accessible by the home owner through a smart phone control application with the proper security credentials from anywhere in the Internet. However, Web crawlers dispatched by a search engine will not discover the home IoT devices because they will not be able to traverse the fire-walled home gateway.

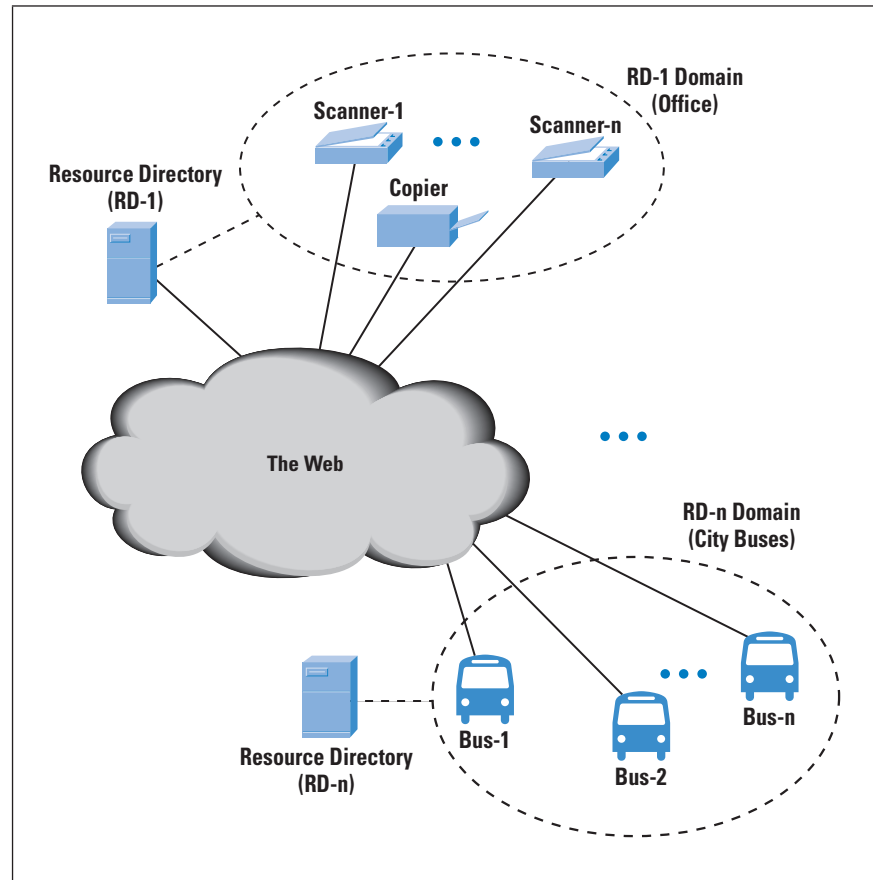
Therefore, the current pull model of Web discovery cannot be applied directly to the expected deployments of IoT networks. In other words, current Web crawler technology is unable to reliably discover a significant percentage of IoT devices that may be asleep or unconnected for significant periods of time, or may be located in semi-closed networks. The result is that traditional Web discovery techniques will not produce accurate discovery results for IoT scenarios.

Resource Directories to Solve the IoT Discovery Problem

The solution currently being standardized in the IETF to address the IoT resource discovery problem is based on a new logical network node called the *Resource Directory* (RD)^[6, 7]. The RD idea was originally conceived and validated in the *European Union* (EU)-funded *SENSEI* research program before coming to the IETF for standardization^[8]. The RD is defined in [6] to be applicable to a given *domain* and not the entire Web. The domain is a logical grouping of IoT devices that are related to an RD. An RD may support multiple domains. The details of defining the extent of a given domain boundary, however, are left to implementation and are not specified. Typically, the RD domains specified in IETF use cases are building-wide, campus-wide, or city-wide. The domain concept maps well into the expected deployment model of IoT devices in semi-closed networks.

In the simplest case, there would be a one-to-one mapping between each semi-closed network and a domain. The RD approach thus provides a distributed resource discovery mechanism for IoT scenarios. Figure 2 shows some typical RD domains.

Figure 2: Typical Resource Directory Domains



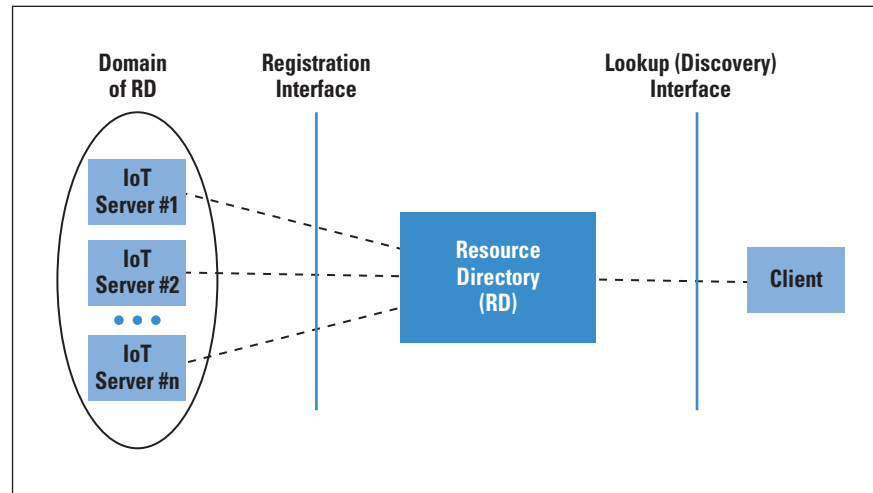
The resource registration step is done in a push fashion by IoT devices acting as mini Web servers pushing their resource information into the RD resource database. Figure 3 shows the architecture of a given RD. All the IoT devices acting as Web servers will first register their resources (URIs) via a registration interface.

Discovery can then be performed on the registered resources by an IoT client using the lookup interface. Mutual authentication, encryption, and access control are required for both the registration and lookup interfaces to ensure security and privacy of the entire resource discovery process.

A given device may use both the registration interface (as a Web server) and the lookup interface (as a client). The client may be located anywhere in the Web, but must have some knowledge regarding which specific RD to direct the resource discovery request to. For example, a newly installed home light controller may perform a lookup on its own home RD to find all the lights installed in the house.

Or, a national smart-grid controller may perform a lookup on a known RD in a remote city to find all the electric transformers located in that city.

Figure 3: IoT Resource Directory Architecture (adapted from [6])



IoT devices communicate with the RD using a *Representational State Transfer* (REST)-based protocol similar to HTTP but optimized for IoT. This protocol is referred to as the *Constrained Application Protocol* (CoAP)^[9]. The resource information pushed by the IoT servers into the RD uses CoAP messages with a specific payload format termed the *Link Format*^[7]. Only the URI, hyperlinks, and some meta-data are sent from the IoT device to the RD. Application content is not sent to the RD. Table 1 shows a comparison of the main resource discovery features of a traditional Web search engine and an IoT RD.

Resource Directory Protocol Considerations

As mentioned previously, CoAP is a Web transfer protocol, similar to HTTP, but optimized for IoT scenarios. CoAP provides a request/response interaction model between clients and servers. It supports key Web concepts such as URIs and Internet media types. CoAP messages are sent over *User Datagram Protocol* (UDP), and the CoAP header is encoded in a simple binary format. A CoAP request consists of a method (that is, GET, PUT, POST, and DELETE) that is applied to a resource identified by its URI, and a payload described by an Internet media type as well as other metadata.

CoAP messages may easily be interworked with HTTP in the forward or reverse directions via special cross-protocol proxies^[9]. In addition, CoAP uses *Datagram Transport Layer Security* (DTLS)^[10] to provide a secure session between the communicating parties.

In CoAP, every physical IoT device is assumed to have one or more resources, each identified by a URI. A resource may contain application information gathered by the IoT device (for example, temperature), or may be a method to control the device (for example, turn it ON/OFF). An example CoAP request and response pair is shown in Table 2.

Table 1: Comparison of Resource Discovery Features of Web Search Engine versus IoT Resource Directory

Characteristic	Traditional Web Search Engine (for example, Google)	IoT Resource Directory
(1) How is resource information initially received by node?	Mainly pulled from target website by Web crawlers after initial visit	Mainly pushed by target IoT devices directly to RD (usually after power-up)
(2) How is updated resource information transferred to node?	Pulled from target website by Web crawlers that revisit according to their search engine policy	Pushed by target IoT device directly to RD according to their own update policy
(3) What resource information is transferred to node?	The entire website (that is, URIs, hyperlinks, metadata, and most application content)	URIs, hyperlinks, and metadata (but no application content is transferred)
(4) What transfer protocols are supported?	HTTP	CoAP (Also some limited HTTP support exists. Further possible enhancements are discussed in [12].)
(5) What is the scope of a client discovery request for resources?	Global (that is, covers entire Internet)	Local within given RD domain (for example, city-wide)
(6) Typical end user that generates query for resources.	Human user (via a Web browser client) sends a search request	IoT device (that is, acting as both the client and end user) sends lookup request May also be used occasionally by human user (for example, via a CoAP-enabled Web browser client as part of management activities)
(7) Are resource discovery results ranked?	Yes	No (but being discussed as a future enhancement in [12])
(8) Are the resource discovery results machine readable?	No (but may support it in the future with further adoption of Semantic Web concept)	Yes (that is, results strictly follow Link Format ^[7])

Table 2: Example CoAP GET Request and Response

Request	GET coap://heater.net/temperature <u>Note:</u> Where Method = GET URI = coap://heater.net/temperature URI-Scheme component = coap:// URI-Host component = heater.net (or alternatively may be an IP address and Port Number) URI-Path component = /temperature
Response	2.05 Content "22.3 C" <u>Note:</u> Where Response code = 2.05 (indicating successful processing) Payload = 22.3 Celsius (C) temperature reading

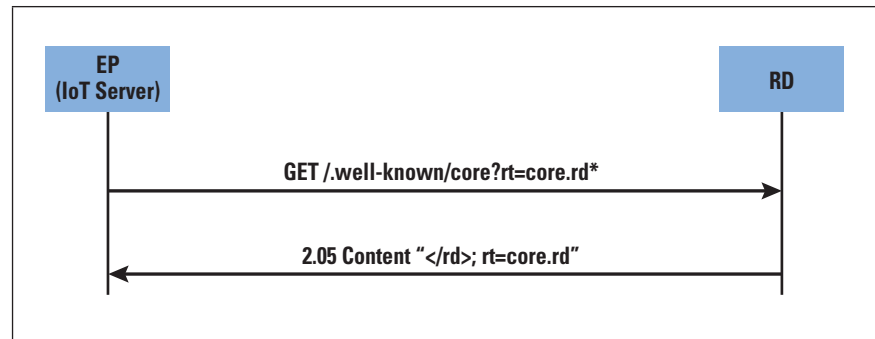
The following sections describe the key protocol steps and security characteristics related to RDs.

Finding the Resource Directory

The first step is for the IoT devices, or *End Points* (EPs) as they are called in [6], to find the appropriate RD. The most dynamic method for finding the RD is using IP multicast. Specifically, the device sends a CoAP multicast message to the CoAP IPv4 or IPv6 addresses reserved for this purpose^[11]. An alternative method would be, for example, factory preprovisioning of the RD information in the device.

Assuming the IP multicast method of finding the RD, each device (EP) sends a CoAP GET request to a specific URI-Path as shown in Figure 4. Specifically, the CoAP GET request is sent by multicast to the reserved “`/.well-known/core`” URI-Path. (Note that the URI-Scheme and URI-host components are not shown for simplicity in this and subsequent figures.) All the devices in the domain will then get this request because it is sent by IP multicast^[11]. However, only the RD will reply because the request URI has a query string for resource type (rt) added to the end (that is, `?rt=core.rd*`), indicating that the message is meant for the RD. The RD then responds indicating its URI-Path (that is, `/rd`) for subsequent registration or lookup requests^[6].

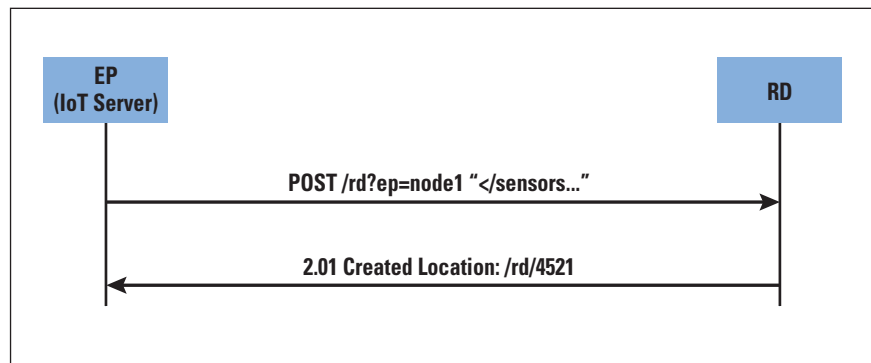
Figure 4: Finding a Resource Directory (adapted from [6])



Registering Resources

After finding the RD, each IoT device (EP) will register its own resources to the RD using the RD registration interface as shown in Figure 5. This registration is accomplished by each device sending a CoAP POST request directly to the RD with its list of URIs (that is, `/sensor...`) in the message payload, along with a query string identifying the registering device (that is, `?ep=node1`). The message payload containing the list of URIs being registered is formatted in the Link Format^[7]. The RD then responds with the resulting URI-Path (that is, `/rd/4521`) that it created to store the resources of the device^[6].

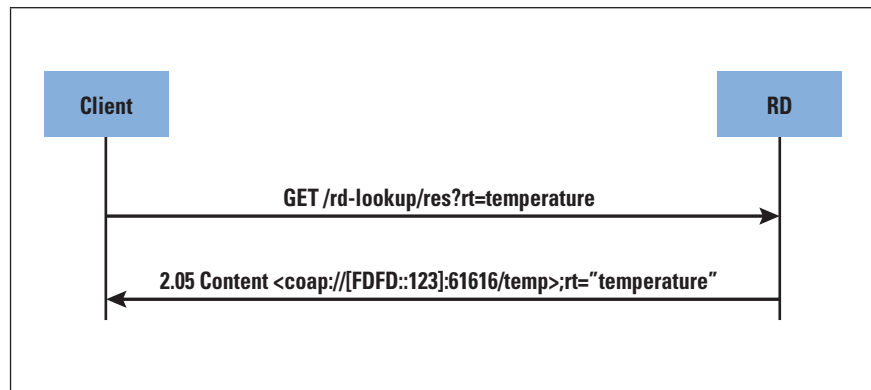
Figure 5: Registration of URIs to a Resource Directory (adapted from [6])



Resource Lookup (Discovery) by Client

The RD also supports a lookup interface for clients to make a discovery request on the RD database. The client may be located in the RD domain or may be outside of it. The client is aware of a given RD because it used the locating mechanism described previously, or it may have learned of the RD through other methods (for example, preprovisioning). Figure 6 shows a typical resource lookup request where a client is interested in finding all URIs related to “temperature.” Specifically, the client will send a GET request to the RD Lookup interface indicating that it is interested in the resource type of temperature in the query string (that is, `?rt=temperature`). The RD will then respond with a message containing the list of URIs of all the devices that it has in its registration database that match this criterion^[6]. The response message is formatted using the Link Format^[7].

Figure 6: Resource Lookup (Discovery) Request Sent to a Resource Directory (adapted from [6])



Other types of lookup requests may also be sent. For example, the RD may be queried to find out all the URIs supported by a given IoT device. Or, the RD may be queried to find out the identities of all the IoT devices in a given domain. The great majority of lookup requests to the RD will be sent by other IoT devices, without any human in the loop, for automated command and control. However, resource lookup requests may also be sent occasionally by humans via a CoAP-enabled Web browser interface for management activities.

Resource Directory Security Characteristics

Both the RD registration and lookup interfaces are protected by multiple layers of security to ensure that only authorized parties can access the RD. Specifically, mutual authentication is first required between the RD and any device attempting to access it. This authentication is accomplished using either preshared encryption keys, raw public keys, or X.509 security certificates as the security credentials^[9]. The appropriate credentials are used in the initial handshake of the DTLS session establishment to perform mutual authentication between the RD and the device or client accessing it. After the mutual authentication is completed, the cipher suite to be used for the DTLS session is negotiated. Then all subsequent messages exchanged between the RD and the device or client are securely encrypted via DTLS so that no unauthorized third party can decipher the communications^[6].

In addition to the DTLS security, the RD will also perform a fine-grained access control of any device attempting to communicate with it. Access control will be performed separately on the RD registration and lookup interfaces. Access control may be performed at the domain, device, or resource level^[6]. This control is especially important on the lookup interface for privacy and security reasons. For example, in a hospital setting many medical devices such as blood pressure monitoring devices may be registered to a RD, but only authorized medical staff should be able to discover a given device for privacy reasons. Or in a home setting, a visitor may be allowed to freely discover the television, but will be blocked from discovering the front door lock for security reasons.

Examples of Resource Directory Implementations

In parallel with the ongoing standardization efforts in the IETF for the RD protocol^[6, 12], there are several open source and commercial instances of RDs that have successfully interoperated with various IoT devices. Some examples are briefly described in the following paragraphs.

The *Californium* open source software project is a popular CoAP framework for IoT deployments. It is written in the *Java* programming language and specifically includes support for back-end infrastructure as part of its project scope. As such, it has released software loads that implement RD functionality that can be run on general-purpose servers^[13].

On the commercial front, ARM, the semiconductor and software company, has released several products for the IoT market. One of its products is a middleware offering called the “mbed Device Server.” This middleware includes support of RD functionality. This middleware software can run on various server hardware platforms^[14].

Another company that has done a lot of RD development work is Ericsson, the telecommunications equipment and service provider. Ericsson has done early prototyping and research^[15] in the RD concept starting from the initial EU SENSEI project days^[8]. The company has also participated in an open source software project for a cloud-based IoT gateway that includes RD functionality^[16].

Alternative Approaches to Discovery

The RD is not the only approach to the discovery problem for IoT networks. There are other methods such as *Domain Name Service – Service Discovery* (DNS-SD), which allows lookup of a given service via DNS^[17]. Another method is *Universal Plug and Play* (UPnP), which allows discovery of devices in home networks^[18].

The key difference between these other discovery methods and the RD approach is that the RD is geared towards resource discovery in the context of a REST-based Web model, meaning discovery of URIs and related metadata. The other existing discovery approaches are mainly oriented to discovering IP addresses, ports, and related parameters. So they are complementary to the URI discovery methods but cannot replace them. The only other widely used URI discovery scheme is the Web crawler approach described previously, which has the shortcomings in IoT deployments as described in Table 1.

Conclusion

The existing REST-based Web architecture and protocols have been extremely successful and a driving force behind the explosive growth of the Internet during the last 20 years. Search engines like *Google* and *Bing*, which use Web crawlers to discover resources (that is, URIs) efficiently, constitute a key part of the success of the Web. However, the existing model of resource discovery is expected to undergo radical changes with the addition in the future of an increasing number of IoT devices acting as both mini Web servers and clients. The IETF is currently standardizing protocol support for the Resource Directory, which will be optimized for distributed IoT resource discovery.

It is expected that an increasing number of discovery requests in the future will be handled by RDs for scenarios involving IoT devices. In parallel, human users will continue to heavily use traditional Web search engines like Google. There is also expected to be some cross-usage because traditional Web browsers may start to support CoAP software modules (plug-ins) and hence allow human users to make direct queries to RDs. However, a limiting feature of this interaction will be the security and privacy requirements of IoT deployments. Specifically, many IoT resources such as personal health-monitoring devices will have sensitive information that is not meant for public distribution, and they may also be located in semi-closed networks. Strong security and privacy is supported by the current RD model, which requires strict mutual authentication, encryption, and access control for both registration and discovery of IoT resources.

References

- [1] Tim Berners-Lee, Roy T. Fielding, and Larry Masinter, “Uniform Resource Identifier (URI): Generic Syntax,” RFC 3986, January 2005.
- [2] Roy Fielding and Julian Reschke, “Hypertext Transfer Protocol (HTTP/1.1): Message Syntax and Routing,” RFC 7230, June 2014.
- [3] The Web Robots Pages: <http://www.robotstxt.org/>
- [4] “What Is SEO, Search Engine Optimization?,”
<http://searchengineland.com/guide/what-is-seo>
- [5] Dominique Barthel, Mischa Dohler, Thomas Watteyne, and Tim Winter, “Urban WSNs Routing Requirements in Low Power and Lossy Networks,” RFC 5548, May 2009.
- [6] Z. Shelby, et al., “CoRE Resource Directory,” Internet Draft, work in progress, March 2016.
draft-ietf-core-resource-directory-07.txt
- [7] Zach Shelby, “Constrained RESTful Environments (CoRE) Link Format,” RFC 6690, August 2012.
- [8] S. Jokic, et al., “Evaluation of an XML Database Based Resource Directory Performance,”
[http://www.smartsantander.eu/downloads/
Presentations/XML_RD_Telfor_2011_v1.0Srdjan.pdf](http://www.smartsantander.eu/downloads/Presentations/XML_RD_Telfor_2011_v1.0Srdjan.pdf)
- [9] Zach Shelby, Carsten Bormann, and Klaus Hartke, “The Constrained Application Protocol (CoAP),” RFC 7252, June 2014.
- [10] Eric Rescorla and Nagendra Modadugu, “Datagram Transport Layer Security Version 1.2,” RFC 6347, January 2012.
- [11] Esko Dijk and Akbar Rahman, “Group Communication for the Constrained Application Protocol (CoAP),” RFC 7390, October 2014.
- [12] A. Rahman, “Advanced Resource Directory Features,” Internet Draft, work in progress, March 2016.
draft-rahman-core-advanced-rd-features-02.txt
- [13] Californiium (Cf) CoAP Framework:
[http://www.eclipse.org/proposals/technology.
californium/](http://www.eclipse.org/proposals/technology/californium/)
- [14] ARM mbed Device Server:
[https://www.mbed.com/en/development/cloud/
mbed-device-server/](https://www.mbed.com/en/development/cloud/mbed-device-server/)

- [15] Ericsson Research Blog: Having a headache using legacy IoT devices?
<https://www.ericsson.com/research-blog/internet-of-things/headache-using-legacy-iot-devices/>
- [16] Ericsson Research Blog: A Computational Engine for the Internet of Things.
<https://www.ericsson.com/research-blog/internet-of-things/computational-engine-internet-things/>
- [17] Stuart Cheshire and Marc Krochmal, “DNS-Based Service Discovery,” RFC 6763, February 2013.
- [18] Wikipedia, “Universal Plug and Play”:
https://en.wikipedia.org/wiki/Universal_Plug_and_Play
- [19] David Lake, Ammar Rayes, and Monique Morrow, “The Internet of Things,” *The Internet Protocol Journal*, Volume 15, No. 3, September 2012.
- [20] William Stallings, “The Internet of Things: Network and Security Architecture,” *The Internet Protocol Journal*, Volume 18, No. 4, December 2015.

AKBAR RAHMAN is a Principal Engineer at InterDigital Communications and is based in the company’s office in Montreal, QC, Canada. He has been closely involved in IoT protocol development at IETF for several years. He has multiple IETF RFCs published in the areas of IoT and Internet architecture. He has a BAsC degree from the University of Waterloo, Canada.

E-mail: Akbar.Rahman@InterDigital.com

CHONGGANG WANG is a Member of Technical Staff at InterDigital Communications and is based in the company’s office in King of Prussia, PA, USA. He is Editor-in-Chief of the *IEEE IoT Journal*, and a Distinguished Lecturer for the IEEE Communications Society. He has a PhD from the Beijing University of Posts and Telecommunications.

E-mail: Chonggang.Wang@InterDigital.com

The IANA Transition

by Vint Cerf, Google

In this article I will explore the notable proposal sent in March 2016^[0] by the *Internet Corporation for Assigned Names and Numbers* (ICANN) to the U.S. Department of Commerce, *National Telecommunication and Information Agency* (NTIA) to end the long-standing contractual relationship between ICANN and NTIA for the conduct of the *Internet Assigned Numbers Authority* functions (“IANA functions”)^[1, 2]. ICANN was formed in late 1998 in response to a White House “White Paper” issued by Ira Magaziner, then a senior advisor for policy to President Bill Clinton. ICANN would undertake to form a private sector entity to carry out the coordinated assignment of Internet domain names, Internet addresses, and the maintenance of parameter registries needed for the operation of the suite of protocols used in the Internet.

These functions had been managed by Jonathan Postel acting as the IANA at University of Southern California’s Information Sciences Institute (and other earlier institutions where Postel had worked) under various government contracts. By 1996, the Internet was experiencing its so-called “dot boom” and the potential scale and liabilities of carrying out the IANA functions led to a serious effort to institutionalize the operation. For lack of space, I will leave out two years of community debate and fast-forward to the creation of ICANN to fulfill these functions. ICANN was conceived as a multi-stakeholder organization drawing on input from the private sector, civil society, governments of the world, and the technical community for the development of policy for the IANA functions and for the coordination of the multiple parties having a role in managing these unique identifiers and parameters.

In 1998, many organizations were involved in the evolution and operation of the Internet, its *Domain Name System* (DNS), Internet address allocation, and standards development. The *Internet Society*, founded in 1991, housed the standards-oriented *Internet Architecture Board* (IAB) and the *Internet Engineering Task Force* (IETF). There were then three *Regional Internet Registries* (RIRs)^[3] for Internet address allocation—RIPE-NCC, APNIC, and ARIN—and two more to follow later (LACNIC and AFRINIC). There were nominally 13 DNS *Root Server* operators providing top-level domain name resolution. Verisign generated and distributed the official domain name root zone based on input from IANA and, under the terms of the NTIA/ICANN contract, authorization from NTIA. Many domain name registries and registrars were created to support DNS operation.

The original plan was for ICANN to operate under NTIA oversight for a few years and then operate as an independent organization. In fact, the contractual obligations extended from 1998 to the present.

In March 2014, however, NTIA proposed that this contractual relationship for the IANA functions should be ended and ICANN be allowed to perform the IANA functions independently. In March 2016, ICANN delivered to NTIA its consolidated proposal from all the constituent parties for the transition from the present contractual relationship to independent operation. The two-year effort leading to this comprehensive proposal was not without considerable debate among all the parties. Many ideas were surfaced, analyzed, argued over, adopted, adapted, or discarded, leading to a consolidated result. The Department of Commerce and the U.S. Congress will be evaluating the proposed new *modus operandi* in the weeks ahead.

Some fears have been voiced that the complex proposal poses risks that authoritarian governments within the ICANN *Governmental Advisory Committee* (GAC) or through some external means might wrest control of ICANN from its multi-stakeholder constituencies. While the proposal should be evaluated on all its merits, I am persuaded the terms and conditions of the proposed operating practices are well protected against such an outcome. A great many conditions must be satisfied before the more extraordinary powers of the *sole designator* can be exercised. The headquarters of ICANN will remain in the U.S. The many entities that cooperate with ICANN to manage core Internet identifier administration have expressed full support for the proposal.

If I have any trepidation about the proposal, it is associated with its general complexity. As the former chairman of ICANN, I am no stranger to the evolution of ICANN's structure and processes and their relative intricacy. The new proposal adds its own unique aspects to this tendency, and it remains to be seen how well the system will work. However, ICANN has shown a remarkable ability to reform and adapt when necessary, and I believe that capacity is preserved under the new proposal. There is still a good deal of work ahead to actually implement what is ultimately approved, but I am confident this community is capable of achieving a successful outcome.

[Ed.: An earlier version of this article appeared in *Communications of the ACM*, Volume 59, No. 5, May 2016.]

References

- [0] "Plan to Transition Stewardship of Key Internet Functions Sent to the U.S. Government,"
<https://www.icann.org/news/announcement-2016-03-10-en>
- [1] "NTIA Announces Intent to Transition Key Internet Domain Name Functions,"
<https://www.ntia.doc.gov/press-release/2014/ntia-announces-intent-transition-key-internet-domain-name-functions>

- [2] “NTIA Finds IANA Stewardship Transition Proposal Meets Criteria to Complete Privatization,”
<https://www.ntia.doc.gov/press-release/2016/iana-stewardship-transition-proposal-meets-criteria-complete-privatization>
- [3] Daniel Karrenberg, Gerard Ross, Paul Wilson, and Leslie Nobile, “Development of the Regional Internet Registry System,” *The Internet Protocol Journal*, Volume 4, No. 4, December 2001.
- [4] IANA Stewardship Transition Coordination Group:
<http://www.ianacg.org/>
- [5] NTIA IANA Functions’ Stewardship Transition:
<https://www.icann.org/stewardship>

VINTON G. CERF is Vice President and Chief Internet Evangelist for Google. He contributes to global policy development and continued spread of the Internet. Widely known as one of the “Fathers of the Internet,” Cerf is the co-designer of the TCP/IP protocols and the architecture of the Internet. He has served in executive positions at MCI, the Corporation for National Research Initiatives and the Defense Advanced Research Projects Agency (DARPA), and on the faculty of Stanford University.

Cerf served as Chairman of the Board of the Internet Corporation for Assigned Names and Numbers (ICANN) from 2000 to 2007 and has been a Visiting Scientist at the Jet Propulsion Laboratory since 1998. He served as founding President of the Internet Society (ISOC) from 1992 to 1995. Cerf is a Fellow of the IEEE, ACM, and American Association for the Advancement of Science, the American Academy of Arts and Sciences, the International Engineering Consortium, the Computer History Museum, the British Computer Society, the Worshipful Company of Information Technologists, and the Worshipful Company of Stationers, and he is a member of the National Academy of Engineering. He currently serves as Past President of the Association for Computing Machinery and Chairman of the American Registry for Internet Numbers (ARIN), and he has completed a term as Chairman of the Visiting Committee on Advanced Technology for the U.S. National Institute of Standards and Technology. President Obama appointed him to the National Science Board in 2012.

Cerf is a recipient of numerous awards and commendations in connection with his work on the Internet, including the U.S. Presidential Medal of Freedom, U.S. National Medal of Technology, the Queen Elizabeth Prize for Engineering, the Prince of Asturias Award, the Tunisian National Medal of Science, the Japan Prize, the Charles Stark Draper Award, the ACM Turing Award, Officer of the Legion d’Honneur, and 25 honorary degrees. In December 1994, *People* magazine identified Cerf as one of that year’s “25 Most Intriguing People.” His personal interests include fine wine, gourmet cooking, and science fiction. Cerf and his wife, Sigrid, were married in 1966 and have two sons, David and Bennett. vint@google.com

RACI

The *RIPE Academic Cooperation Initiative* (RACI) connects members of the academic community with the RIPE community by inviting students and researchers to present at meetings organized by the RIPE NCC. Successful applicants receive complimentary tickets, travel and accommodation to meetings and the opportunity to present their work to some of the leading technical figures in the Internet world. Examples of relevant topics include:

- Network Measurement and Analyses
- IPv6 Deployment
- BGP Routing
- Network Security
- Internet Governance
- Peering and Interconnectivity
- The Internet of Things

For more information about RACI, including the application process and deadlines, visit: <http://ripe.net/raci>

NTIA Issues IANA Transition Proposal Report

On 9 June 2016, *The National Telecommunications and Information Administration* (NTIA) issued its assessment report on the *IANA Stewardship Transition Proposal*. (Ed.: See article on page 26). In order to be accepted, the proposal needed to be shown to have broad community support and address the following four principles:

- Support and enhance the multistakeholder model
- Maintain the security, stability, and resiliency of the Internet DNS
- Meet the needs and expectations of the global customers and partners of the IANA services
- Maintain the openness of the Internet

The NTIA further stipulated that “it would not accept a proposal that replaces its role with a government-led or intergovernmental organization solution.” After thorough review the NTIA reports that it finds that “the IANA Stewardship Transition Proposal meets the criteria necessary to complete the long-promised privatization of the IANA functions.”

The full report is available at:

<https://www.ntia.doc.gov/report/2016/iana-stewardship-transition-proposal-assessment-report>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Publication of this journal is made possible by:

Supporters



Diamond Sponsors



Ruby Sponsor



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



Individual Sponsors

Lyman Chapin, Steve Corbató, Dave Crocker, Jay Etchings, Martin Hannigan, Hagen Hultzs, Dennis Jennings, Jim Johnston, Merike Kao, Bobby Krupczak, Richard Lamb, Tracy LaQuey Parker, Bill Manning, Andrea Montefusco, Tariq Mustafa, Mike O'Connor, Tim Pozar, George Sadowsky, Scott Seifel, Helge Skrivervik, Rob Thomas, Tom Vest, Rick Wesson.

For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

ADDRESS SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Fred Baker, Cisco Fellow
Cisco Systems, Inc.

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

Dr. Steve Crocker, Chairman
Internet Corporation for Assigned Names and Numbers

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol *Journal*

November 2016

Volume 19, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

In This Issue

From the Editor	1
Internet E-Mail Security.....	2
SDN Complexity and Reality	31
Thank You.....	42
Supporters and Sponsors	43

FROM THE EDITOR

Publication of this journal is made possible by numerous individuals and organizations. Every year in late August we initiate a sponsorship renewal campaign. This year our funding fell short of our sponsorship target, so we delayed publication of the September issue until now. If you are a subscriber, you should have received an e-mail asking for a donation. I am happy to say that many subscribers responded to that request (see page 42), and with the help of these individuals and our corporate sponsors we are now ready to deliver to you this *November* issue. This will be the third and final issue in 2016, but we hope to return to our regular quarterly publication schedule in 2017. We still need more individual and corporate sponsors, so please ask your company to sign up for a sponsorship or make a donation at <http://tinyurl.com/IPJ-donate>

We are pleased to welcome two new members of our Editorial Advisory Board. David Conrad is the Chief Technical Officer at the *Internet Corporation for Assigned Names and Numbers* (ICANN), and Cullen Jennings is a Cisco Fellow at Cisco Systems, Inc. We are grateful to Fred Baker, who has left Cisco Systems and our Editorial Advisory Board, and we wish him every success in the future.

A large percentage of Internet traffic is electronic mail. In our first article, William Stallings gives an overview of the many enhancements that are designed to make e-mail communication more secure and reliable in the face of an increasing amount of spam and other attack vectors.

Previous articles in IPJ have covered various aspects of *Cloud Computing* and *Software-Defined Networks* (SDNs). In our second article, Russ White and Shawn Zandi take a critical look at the complexity of these technologies.

As always, we welcome your feedback, suggestions, book reviews, articles, and sponsorship support. You can contact us by sending an e-mail to ipj@protocoljournal.org and visit our website for subscription information, back issues, author guidelines, sponsor information, and much more.

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

Comprehensive Internet E-Mail Security

by William Stallings, Independent Consultant

For both organizations and individuals, e-mail is pervasive and vulnerable to a wide range of security threats. In general terms, e-mail security threats can be classified as follows:

- *Authenticity-related threats*: Could result in unauthorized access to an enterprise's e-mail system. Another threat in this category is deception, in which the purported author isn't the actual author.
- *Integrity-related threats*: Could result in unauthorized modification of e-mail content.
- *Confidentiality-related threats*: Could result in unauthorized disclosure of sensitive information.
- *Availability-related threats*: Could prevent end users from being able to send or receive e-mail messages.

To assist in addressing these threat categories, the *National Institute of Standards and Technology* (NIST) has issued SP 800-177^[1], which recommends guidelines for enhancing trust in e-mail. The document is both a survey of available standardized protocols and a set of recommendations for using these protocols to counter security threats to e-mail usage.

For an understanding of the topics in this article, it is useful to have a basic grasp of the Internet mail architecture, which is currently defined in RFC 5598^[2]. The discussion now provides an overview of the basic concepts.

At its most fundamental level, the Internet mail architecture consists of a user world in the form of *Message User Agents* (MUA), and the transfer world, in the form of the *Message Handling Service* (MHS), which is composed of *Message Transfer Agents* (MTA). The MHS accepts a message from one user and delivers it to one or more other users, creating a virtual MUA-to-MUA exchange environment. This architecture involves three types of interoperability. One is directly between users: messages must be formatted by the MUA on behalf of the message author so that the message can be displayed to the message recipient by the destination MUA. There are also interoperability requirements between the MUA and the MHS—first when a message is posted from an MUA to the MHS and later when it is delivered from the MHS to the destination MUA. Interoperability is required among the MTA components along the transfer path through the MHS.

Figure 1: Function Modules and Standardized Protocols Used Between Them in the Internet Mail Architecture

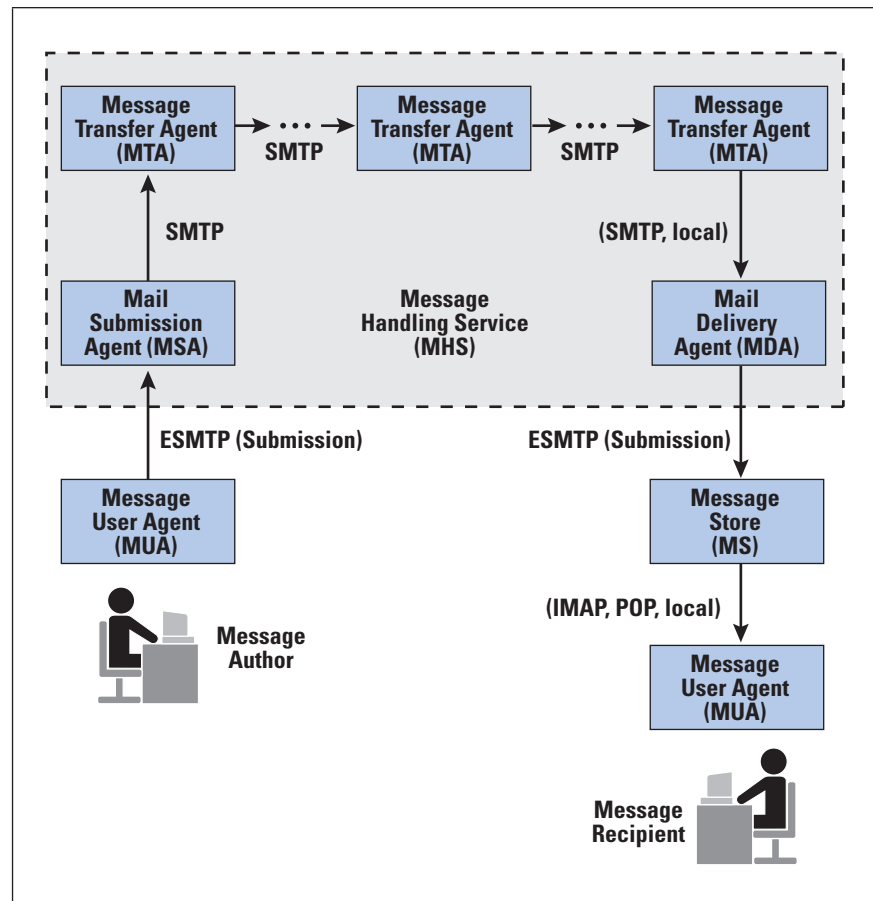


Figure 1 illustrates the key components of the Internet mail architecture, which include the following:

- *Message User Agent (MUA)*: Operates on behalf of user actors and user applications. It is their representative within the e-mail service. Typically, this function is housed in the user's computer and is referred to as a *client* e-mail program or a local network e-mail *server*. The author MUA formats a message and performs initial submission into the MHS via a *Mail Submission Agent (MSA)*. The recipient MUA processes received mail for storage and/or display to the recipient user.
- *Mail Submission Agent (MSA)*: Accepts the message submitted by an MUA and enforces the policies of the hosting domain and the requirements of Internet standards. This function may be located together with the MUA or as a separate functional model. In the latter case, the *Simple Mail Transfer Protocol (SMTP)* is used between the MUA and the MSA.
- *Message Transfer Agent (MTA)*: Relays mail for one application-level hop. It is like a packet switch or IP router in that its job is to make routing assessments and move the message closer to the recipients. Relaying is performed by a sequence of MTAs until the message reaches a destination MDA. An MTA also adds trace information to the message header. SMTP is used between MTAs and between an MTA and an MSA or MDA.

- *Mail Delivery Agent* (MDA): The MDA is responsible for transferring the message from the MHS to the *Message Store* (MS).
- *Message Store* (MS): An MUA can employ a long-term MS. An MS can be located on a remote server or on the same machine as the MUA. Typically, an MUA retrieves messages from a remote server using the *Post Office Protocol* (POP) or the *Internet Message Access Protocol* (IMAP).

As will be seen subsequently, an important element in securing e-mail is the use of *public-key cryptography*. In turn, the use of public-key cryptography depends on the use of *Public-key Certificates*. In essence, a public-key certificate consists of a public key plus a user ID of the key owner, with the whole block signed by a trusted third party. A common scheme for the creation and use of public key certificates is by means of a third party known as a *Certificate Authority* (CA). A CA is an entity that is trusted by the user community, such as a government agency or a financial institution. The essential elements in the CA scheme include:

1. Client software creates a pair of keys, one public and one private. The client prepares an unsigned certificate that includes a user ID and the user's public key. The client then sends the unsigned certificate to a CA in a secure manner.
2. A CA creates a signature by calculating the hash code of the unsigned certificate and encrypting the hash code with the CA's private key; the encrypted hash code is the signature. The CA attaches the signature to the unsigned certificate and returns the now-signed certificate to the client.
3. The client may send its signed certificate to any other user. That user may verify that the certificate is valid by calculating the hash code of the certificate (not including the signature), decrypting the signature using the CA's public key, and comparing the hash code to the decrypted signature.

If all users subscribe to the same CA, then there is a common trust of that CA. All user certificates can be placed in the directory for access by all users. In addition, users can transmit their certificate directly to other users. In either case, when B is in possession of A's certificate, B has confidence that messages it encrypts with A's public key will be secure from eavesdropping and that messages signed with A's private key are unforgeable.

If there is a large community of users, it may not be practical for all users to subscribe to the same CA. Because it is the CA that signs certificates, each participating user must have a copy of the CA's own public key to verify signatures. This public key must be provided to each user in an absolutely secure (with respect to integrity and authenticity) way so that the user has confidence in the associated certificates.

Thus, with many users it may be more practical for there to be numerous CAs, each of which securely provides its public key to some fraction of the users. In practice, there is not a single CA but rather a hierarchy of CAs. This setup complicates the problems of key distribution and trust, but the basic principles are the same.

Several issues with the use of CAs should be mentioned. As can be deduced from the preceding paragraph, a hierarchical CA system can become cumbersome and not scale well. Nevertheless, this scheme is still the preferred one, and it is recommended by SP 800-177. A separate issue is one of security. The global CA ecosystem has become subject to attack in recent years, and has been successfully compromised more than once. One way to protect against CA compromises is to use the *Domain Name System* (DNS) to allow domains to specify their intended certificates or vendor CAs. Such uses of DNS require that the DNS itself be secured with *Domain Name System Security Extensions* (DNSSEC) as described subsequently.

For the reader who needs an introduction or refresher on concepts of public-key cryptography, authentication, and digital signatures, a *Crypto Portal* white paper^[3] provides a quick and easy overview. A useful overview of CA and public-key certificate concepts is NIST SP 800-32^[4].

Trustworthy E-Mail

The following protocols and standards are described in and recommended by SP 800-177:

- *STARTTLS*: An SMTP security extension that enables an SMTP client and server to negotiate the use of *Transport Layer Security* (TLS) to provide private, authenticated communication across the Internet.
- *Secure Multipurpose Internet Mail Extensions* (S/MIME): Provides authentication, integrity, nonrepudiation (via digital signatures) and confidentiality (via encryption) of the message body carried in SMTP messages.
- *DNS-Based Authentication of Named Entities* (DANE): Designed to overcome problems in the *Certificate Authority* (CA) system by providing an alternative channel for authenticating public keys based on DNSSEC, with the result that the same trust relationships used to certify IP addresses are used to certify servers operating on those addresses.
- *Sender Policy Framework* (SPF): Enables a domain owner to specify the IP addresses of MTAs that are authorized to send mail on its behalf. SPF uses the DNS to allow domain owners to create records that associate the domain name with a specific IP address range of authorized MTAs. It is a simple matter for receivers to check the *SPF text record* (TXT) in the DNS to confirm that the purported sender of a message is permitted to use that source address and reject mail that does not come from an authorized IP address.

- *DomainKeys Identified Mail* (DKIM): Enables e-mail actors (authors or operators) to affix their domain name to the message reliably, using cryptographic techniques, so that filtering engines can develop an accurate reputation for the domain. The MTA can sign selected headers and the body of a message. This signature validates the source domain of the mail and provides message body integrity.
- *Domain-based Message Authentication, Reporting, and Conformance* (DMARC): Publishes a requirement for the author domain name to be authenticated by DKIM and/or SPF, for that domain's owner to request recipient handling of nonauthenticated mail using that domain, and for a reporting mechanism to send reports from recipients back to domain owners. DMARC lets senders know the proportionate effectiveness of their SPF and DKIM policies, and signals to receivers what action should be taken in various individual and bulk attack scenarios.

Figure 2 shows how these components interact to provide message authenticity and integrity. Not shown, for simplicity, is that S/MIME also provides message confidentiality by encrypting messages. Together, these protocols provide a comprehensive Internet e-mail security strategy. This article provides an overview of each.

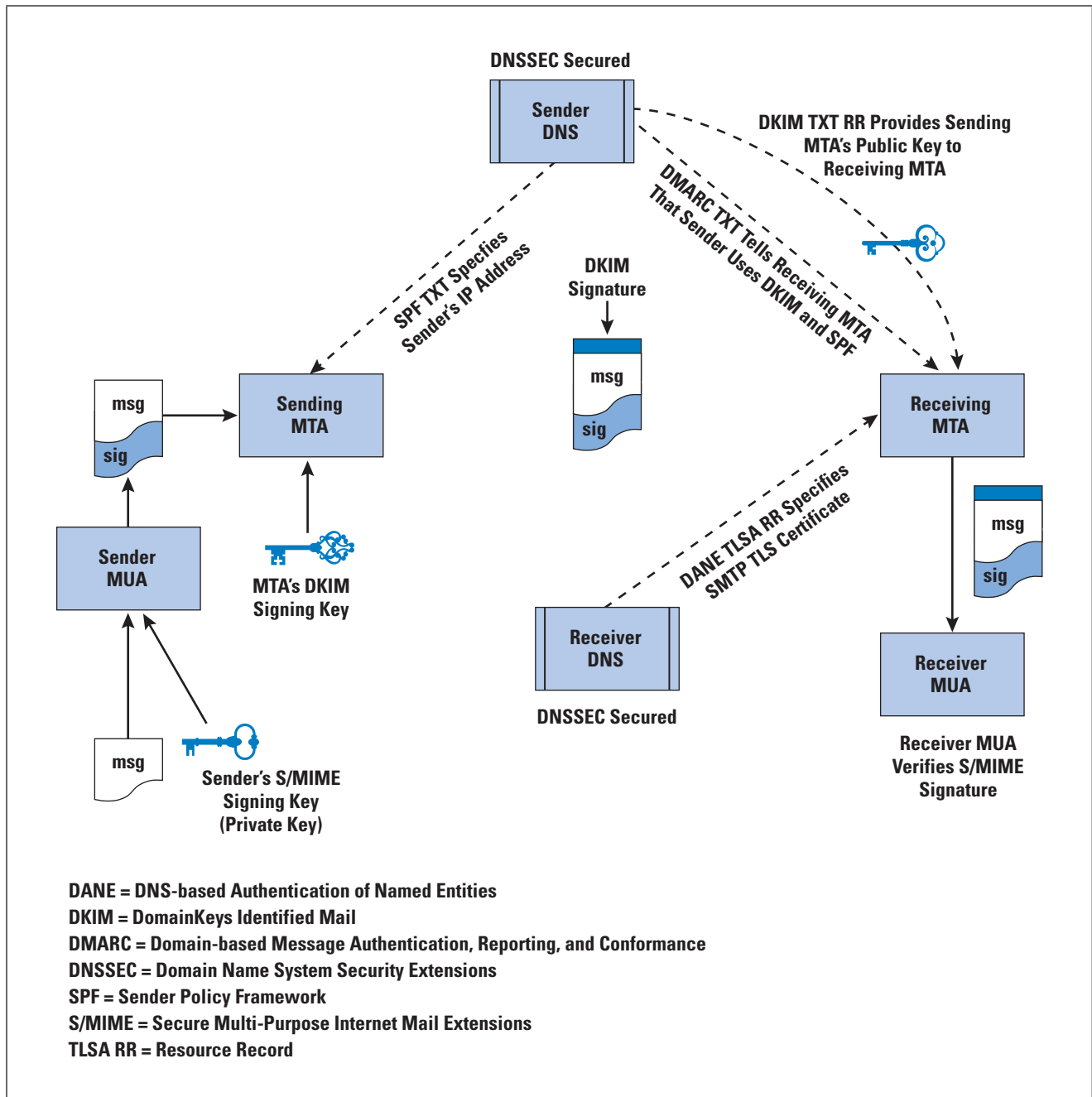
STARTTLS

A significant security-related extension for SMTP is STARTTLS, defined in RFC 3207^[5]. STARTTLS enables the addition of confidentiality and authentication in the exchange between SMTP agents. This addition enables SMTP agents to protect some or all of their communications from eavesdroppers and attackers by invoking a *Transport Layer Security* (TLS) session within the SMTP connection. STARTTLS has been widely deployed, and is supported by Amazon, Facebook, Google, Microsoft, Yahoo, and others^[6]. A 2014 study by Facebook, which sends several billions of e-mails daily, found that 76% of host names that receive Facebook e-mails support STARTTLS^[7].

TLS is a security layer implemented just above TCP. TLS is an Internet Standard that replaces *Secure Sockets Layer* (SSL) with essentially the same functionality^[8]. With TLS in place, an application has a TLS *socket address* and communicates to the TLS socket address at the remote application. These addresses are distinct from those used by the same application running directly over TCP. The security functions provided by TLS are transparent to the application and also to TCP. Thus, neither TCP nor the application needs to be modified to invoke the security features of SSL. TLS provides three categories of security, confidentiality, and authentication.

If the client does initiate the connection over a TLS-enabled port, the server may prompt with a message indicating that the STARTTLS option is available.

Figure 2: The Interrelationship of DNSSEC, SPF, DKIM, DMARC, DANE, and S/MIME for Assuring Message Authenticity and Integrity



The client can then issue the STARTTLS command in the SMTP command stream, and the two parties proceed to establish a secure TLS connection. Many e-mail providers and servers now have STARTTLS enabled^[9, 10], including Amazon, Comcast, Dropbox, Facebook, Google, Microsoft, and Yahoo.

As described in SP 800-177, STARTTLS may be vulnerable to a *Man-In-The-Middle* (MITM) attack when it is initiated as a request by the server. In this case, the MITM receives the STARTTLS request from the server reply to a connection request, and scrubs it out.

The initiating client sees no TLS upgrade request and proceeds with an unsecured connection. However, SP 800-177 takes the position that some security is better than no security and that until TLS is available everywhere and automatically invoked, TLS-capable servers must prompt clients to invoke the STARTTLS command. TLS clients should attempt to either use STARTTLS initially or issue the command when requested.

S/MIME

Secure/Multipurpose Internet Mail Extension (S/MIME) is a security enhancement to the MIME Internet e-mail format standard^[11]. S/MIME is a complex capability that is defined in many documents. The most important documents relevant to S/MIME include the following [12–19]:

- *RFC 5750, S/MIME Version 3.2 Certificate Handling*: Specifies conventions for X.509 certificate usage by S/MIME v3.2.
- *RFC 5751, S/MIME Version 3.2 Message Specification*: The principal defining document for S/MIME message creation and processing.
- *RFC 4134, Examples of S/MIME Messages*: Gives examples of message bodies formatted using S/MIME.
- *RFC 2634, Enhanced Security Services for S/MIME*: Describes four optional security service extensions for S/MIME.
- *RFC 5652, Cryptographic Message Syntax (CMS)*: Describes CMS. This syntax is used to digitally sign, digest, authenticate, or encrypt arbitrary message content.
- *RFC 3370, CMS Algorithms*: Describes the conventions for using several cryptographic algorithms with the CMS.
- *RFC 5752, Multiple Signatures in CMS*: Describes the use of multiple, parallel signatures for a message.
- *RFC 1847, Security Multiparts for MIME—Multipart/Signed and Multipart/Encrypted*: Defines a framework within which security services may be applied to MIME body parts. The use of a digital signature is relevant to S/MIME, as explained subsequently.

S/MIME functionality is built into most modern e-mail software and interoperates between them. S/MIME provides four message-related services: authentication, confidentiality, compression, and e-mail compatibility.

Authentication is provided by means of a digital signature. Most commonly RSA with SHA-256 is used. The sequence is as follows:

1. The sender creates a message.
2. SHA-256 is used to generate a 256-bit message digest of the message.
3. The message digest is encrypted with RSA using the sender's private key, and the result is appended to the message. Also appended is identifying information for the signer, which will enable the receiver to retrieve the signer's public key.
4. The receiver uses RSA with the sender's public key to decrypt and recover the message digest.
5. The receiver generates a new message digest for the message and compares it with the decrypted hash code. If the two match, the message is accepted as authentic.

The combination of SHA-256 and RSA provides an effective digital signature scheme. Because of the strength of RSA, the recipient is assured that only the possessor of the matching private key could have generated the signature. Because of the strength of SHA-256, the recipient is assured that no one else could generate a new message that matches the hash code and, hence, the signature of the original message.

Although signatures normally are found attached to the message or file that they sign, it is not always the case: Detached signatures are supported. A detached signature may be stored and transmitted separately from the message it signs. This option is useful in several contexts. A user may wish to maintain a separate signature log of all messages sent or received. A detached signature of an executable program can detect subsequent virus infection. Finally, detached signatures can be used when more than one party must sign a document, such as a legal contract. Each person's signature is independent and is therefore applied only to the document. Otherwise, signatures would have to be nested, with the second signer signing both the document and the first signature, and so on.

S/MIME provides *confidentiality* by encrypting messages using conventional encryption with a secret key, also known as a *symmetric key*. Most commonly, *Advanced Encryption Standard* (AES) with a 128-bit key is used, with the *Cipher Block Chaining* (CBC) mode. The key itself is also encrypted, typically with RSA, as explained subsequently.

As always, one must address the problem of key distribution. In S/MIME, each symmetric key, referred to as a *content-encryption key*, is used only once. That is, a new key is generated as a random number for each new message. Because it is to be used only once, the content-encryption key is bound to the message and transmitted with it. To protect the key, it is encrypted with the receiver's public key.

The sequence can be described as follows:

1. The sender generates a message and a random 128-bit number to be used as a content-encryption key for this message only.
2. The message is encrypted using the content-encryption key.
3. The content-encryption key is encrypted with RSA using the recipient's public key and is attached to the message.
4. The receiver uses RSA with its private key to decrypt and recover the content-encryption key.
5. The content-encryption key is used to decrypt the message.

As Figure 3 illustrates, both confidentiality and encryption may be used for the same message. The figure shows a sequence in which a signature is generated for the plaintext message and appended to the message. Then the plaintext message and signature are encrypted as a single block using symmetric encryption and the symmetric encryption key is encrypted using public-key encryption.

S/MIME allows the signing and message encryption operations to be performed in either order. If signing is done first, the identity of the signer is hidden by the encryption. Plus, it is generally more convenient to store a signature with a plaintext version of a message. Furthermore, for purposes of third-party verification, if the signature is performed first, a third party need not be concerned with the symmetric key when verifying the signature.

If encryption is done first, it is possible to verify a signature without exposing the message content. This option can be useful in a context in which automatic signature verification is desired, as no private-key material is required to verify a signature. However, in this case the recipient cannot determine any relationship between the signer and the unencrypted content of the message.

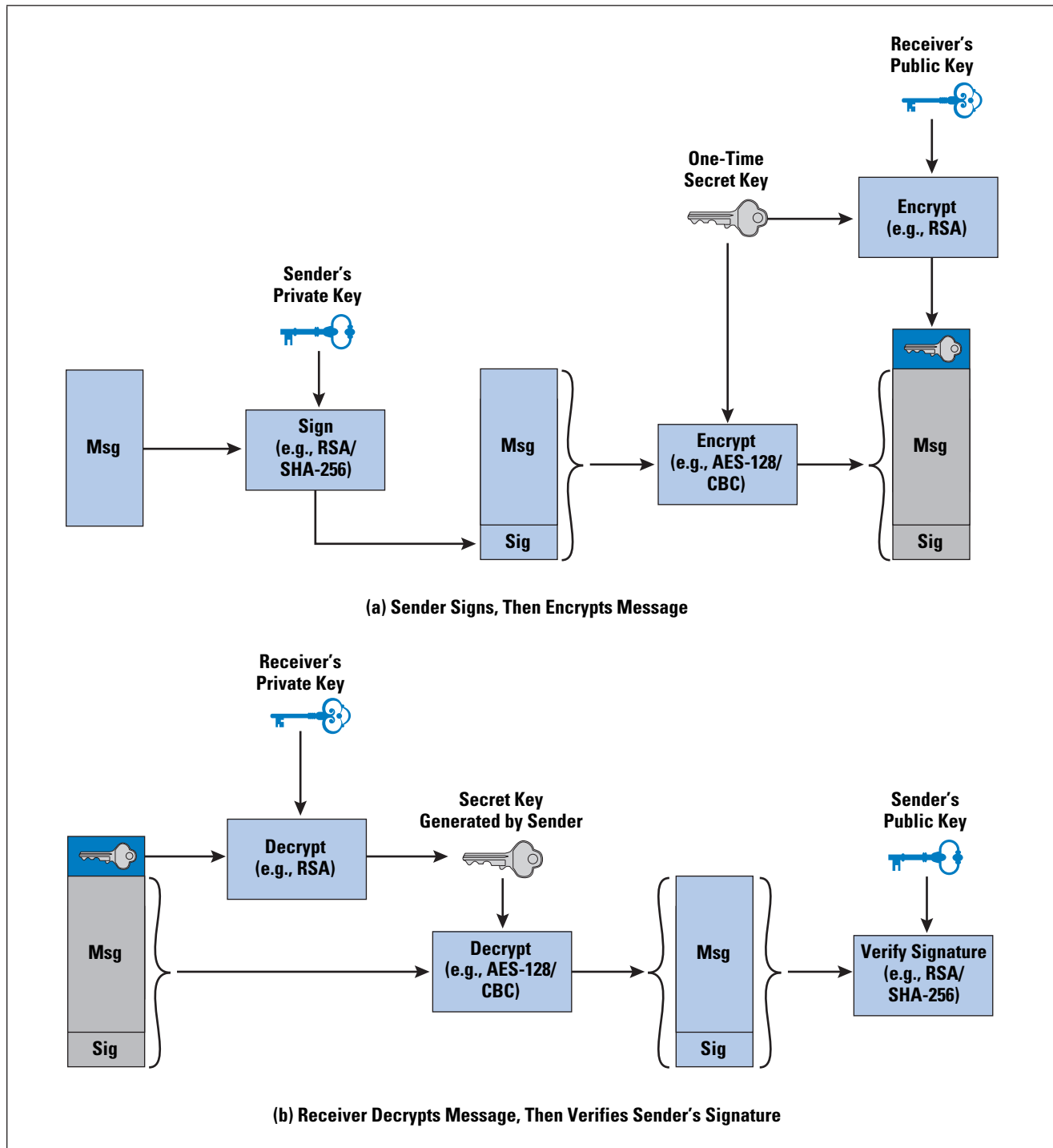
When S/MIME is used, at least part of the block to be transmitted is encrypted. If only the signature service is used, then the message digest is encrypted (with the sender's private key). If the confidentiality service is used, the message plus signature (if present) are encrypted (with a one-time symmetric key). Thus, part of or the entire resulting block consists of a stream of arbitrary 8-bit octets. However, many electronic mail systems only permit the use of blocks consisting of ASCII text. To accommodate this restriction and provide *compatibility*, S/MIME provides the service of converting the raw 8-bit binary stream to a stream of printable ASCII characters, a process referred to as 7-bit encoding. The scheme typically used for this purpose is Base64 conversion. Each group of three octets of binary data is mapped into four ASCII characters.

S/MIME also offers the ability to *compress* a message. Message compression has the benefit of saving space for both e-mail transmission and file storage. Compression can be applied in any order with respect to the signing and message encryption operations.

RFC 5751 provides the following guidelines:

- Compression of binary encoded encrypted data is discouraged, since it will not yield significant compression. Base64 encrypted data could very well benefit, however.
- If a lossy compression algorithm is used with signing, you will need to compress first, then sign.

Figure 3: Simplified S/MIME Functional Flow



SP 800-177 recommends the use of certificate chain authentication against a known certificate authority. Further, SP 800-177 indicates that users who want more assurance that the public key supplied is bound to the sender's domain may use a work-in-progress DANE-S/MIME mechanism^[20], in which the certificate and key can be independently retrieved from the DNS and authenticated per the DANE mechanism described subsequently.

In addition, SP 800-177 notes that MUAs typically use S/MIME private keys to decrypt the e-mail message each time it is displayed, but leave the message encrypted in the e-mail store. This mode of operation is not recommended, as it forces recipients of the encrypted e-mail to maintain their private key indefinitely. Instead, the e-mail should be decrypted prior to being stored in the mail store. The mail store, in turn, should be secured using an appropriate cryptographic technique (for example, disk encryption), extending protection to both encrypted and unencrypted e-mail.

OpenPGP

Pretty Good Privacy (PGP) was developed by Phil Zimmermann as a publicly-available freeware package to enable individuals to exchange secure e-mails without the need to rely on any institution. Efforts began early on to develop Internet standards for PGP^[21], culminating in *OpenPGP*. OpenPGP^[22, 23] is a proposed Internet Standard for providing authentication and confidentiality for e-mail messages. Although it is similar in purpose and functionality to S/MIME, OpenPGP uses different message and key formats and a different approach to establishing and using certificates. SP 800-177 cites many difficulties with OpenPGP, including lack of usability, scalability issues related to key distribution, and lack of authentication of key owners. Further discussion can be found in [24] and [25]. Accordingly, SP 800-177 recommends the use of only S/MIME and deprecates the use of OpenPGP.

DNS and DNSSEC

As background for the following sections, this section briefly reviews DNS and DNSSEC. The *Domain Name System* (DNS) is a directory lookup service that provides a mapping between the name of a host on the Internet and its numerical Internet address. Four elements comprise the DNS. The domain name space is a tree-structured name space to identify resources on the Internet. The DNS database is a collection of resource records organized into a distributed database; conceptually, each node and leaf in the name-space tree structure names a collection of information (for example, IP address, name server for this domain name) that is contained in *Resource Records* [RRs]). *Name Servers* are server programs that hold information about a portion of the domain-name tree structure and the associated RRs. *Resolvers* are programs that extract information from name servers in response to client requests. A typical client request is for an IP address corresponding to a given domain name.

The DNS database is divided into thousands of separately managed *zones*, which are managed by separate administrators. Using this database, DNS servers provide a name-to-address directory service for network applications that need to locate specific application servers.

Domain Name System Security Extensions (DNSSEC)^[26] is used by several protocols that provide e-mail security. DNSSEC provides end-to-end protection through the use of digital signatures that are created by responding zone administrators and verified by a recipient's resolver software. In particular, DNSSEC avoids the need to trust intermediate name servers and resolvers that cache or route the DNS records originating from the responding zone administrator before they reach the source of the query. DNSSEC consists of a set of new resource record types and modifications to the existing DNS protocol.

In essence, DNSSEC is designed to protect DNS clients from accepting forged or altered DNS resource records. It protects these clients by using digital signatures to provide: (1) data origin authentication to ensure that a RR has originated from the correct source; and (2) data integrity verification to ensure that the content of a RR has not been modified. The DNS zone administrator digitally signs every *Resource Record set* (RRset) in the zone, and publishes this collection of digital signatures, along with the zone administrator's public key, in the DNS itself.

In DNSSEC, trust in the public key (for signature verification) of the source is established not by going to a third party or a chain of third parties (as in *Public-Key Infrastructure* [PKI] chaining), but by starting from a trusted zone (such as the root zone) and establishing the chain of trust down to the current source of response through successive verifications of the signature of the public key of a child by its parent. The public key of the trusted zone is called the *trust anchor*.

DANE

DNS-Based Authentication of Named Entities (DANE)^[27, 28] is a protocol that provides mechanisms for domains to specify which X.509 certificates, which are commonly used for *Transport Layer Security* (TLS), should be trusted for the domain. DANE enables certificates to be bound to DNS names using DNSSEC. It is proposed in RFC 6698^[29] as a way to authenticate TLS client and server entities without a *Certificate Authority* (CA).

Briefly, DANE is an alternative mechanism for securely distributing information about domain names by using DNS. DANE defines a new type of DNS record that enables a domain to sign statements specifying which entities are authorized to represent it. Applications can use these records either to augment the existing system of CAs or to create a new chain of trust, rooted in the DNS.

The rationale for DANE is the vulnerability of the use of CAs in a global *Public-Key Infrastructure* (PKI) system. Every browser developer and operating system supplier maintains a list of CA root certificates as trust anchors. These certificates are called the *root certificates* of the software and are stored in its root certificate store. The PKI scheme allows a certificate recipient to trace a certificate back to the root. So long as the root certificate remains trustworthy and the authentication concludes successfully, the client can proceed with the connection. However, if any of the hundreds of CAs operating on the Internet is compromised, the effects can be widespread. The attacker can obtain the private key of the CA, be issued certificates under a false name, or introduce new bogus root certificates into a root certificate store. There is no limitation of scope for the global PKI, and a compromise of a single CA damages the integrity of the entire PKI system. In addition, some CAs have engaged in poor security practices. For example, some CAs have issued wildcard certificates that allow the holder to issue sub-certificates for any domain or entity, anywhere in the world.

The purpose of DANE is to replace reliance on the security of the CA system with reliance on the security provided by DNSSEC. This protocol is well expressed in RFC 6698:

“DNS-Based Authentication of Named Entities (DANE) offers the option to use the DNSSEC infrastructure to store and sign keys and certificates that are used by TLS. DANE is envisioned as a preferable basis for binding public keys to DNS names, because the entities that vouch for the binding of public key data to DNS names are the same entities responsible for managing the DNS names in question. While the resulting system still has residual security vulnerabilities, it restricts the scope of assertions that can be made by any entity, consistent with the naming scope imposed by the DNS hierarchy. As a result, DANE embodies the security “principle of least privilege” that is lacking in the current public CA model.”

DANE defines a new DNS record type, TLSA, which can be used for a secure method of authenticating *Secure Sockets Layer/Transport Layer Security* (SSL/TLS) certificates. The TLSA provides for:

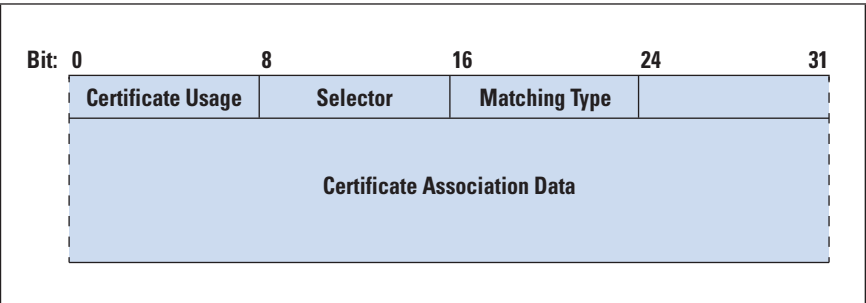
- Specifying constraints on which CA can vouch for a certificate, or which specific PKI end-entity certificate is valid.
- Specifying that a service certificate or a CA can be directly authenticated in the DNS itself.

The TLSA RR enables certificate issue and delivery to be tied to a given domain. A server domain owner creates a TLSA resource record that identifies the certificate and its public key. When a client receives an X.509 certificate in the TLS negotiation, it looks up the TLSA RR for that domain and matches the TLSA data against the certificate as part of the client’s certificate validation procedure.

Figure 4 shows the format of a TLSA RR as it is transmitted to a requesting entity. It contains four fields. The *Certificate Usage* field defines four different usage models, to accommodate users who require different forms of authentication. The usage models follow:

- *PKIX-TA (CA constraint)*: Specifies which CA should be trusted to authenticate the certificate for the service. This usage model limits which CA can be used to issue certificates for a given service on a host. The server certificate chain must pass PKIX validation that terminates with a trusted root certificate stored in the client.
- *PKIX-EE (service certificate constraint)*: Defines which specific end-entity service certificate should be trusted for the service. This usage model limits which end-entity certificate can be used by a given service on a host. The server certificate chain must pass PKIX validation that terminates with a trusted root certificate stored in the client.
- *DANE-TA (trust anchor assertion)*: Specifies a domain-operated CA to be used as a trust anchor. This usage model allows a domain-name administrator to specify a new trust anchor—for example, if the domain issues its own certificates under its own CA that is not expected to be in the end users’ collection of trust anchors. The server certificate chain is self-issued and does not need to verify against a trusted root stored in the client.
- *DANE-EE (domain-issued certificate)*: Specifies a domain-operated CA to be used as a trust anchor. This certificate usage allows for a domain-name administrator to issue certificates for a domain without involving a third-party CA. The server certificate chain is self-issued and does not need to verify against a trusted root stored in the client.

Figure 4: TLSA RR
Transmission Format



The first two usage models are designed to coexist with and strengthen the public CA system. The final two usage models operate without the use of public CAs.

The *Selector* field indicates whether the full certificate or just the value of the public key will be matched. The match is made between the certificate presented in TLS negotiation and the certificate in the TLSA RR. The *Matching Type* field indicates how the match of the certificate is made. The options are exact match, SHA-256 hash match, or SHA-512 hash match. The *Certificate Association Data* is the raw certificate data in hex format.

DANE can be used in conjunction with SMTP over TLS, as provided by STARTTLS, to more fully secure e-mail delivery. DANE can authenticate the certificate of the SMTP submission server that the user's mail client (MUA) communicates with. It can also authenticate the TLS connections between SMTP servers (MTAs). The use of DANE with SMTP is documented in RFC 7672^[30].

As discussed previously, SMTP can use the STARTTLS extension to run SMTP over TLS, so that the entire e-mail message plus SMTP envelope are encrypted. This option is used if both sides support STARTTLS. Even when TLS is used to provide confidentiality, it is vulnerable to attack in the following ways:

- Attackers can strip away the TLS capability advertisement and downgrade the connection to not use TLS.
- TLS connections are often unauthenticated (for example, the use of self-signed certificates as well as mismatched certificates is common).

DANE can address both these vulnerabilities. A domain can use the presence of the TLSA RR as an indicator that encryption must be performed, thus preventing malicious downgrade. A domain can authenticate the certificate used in the TLS connection setup using a DNSSEC-signed TLSA RR.

DNSSEC can be used in conjunction with S/MIME to more fully secure e-mail delivery, in a manner similar to the DANE functionality. This use is documented in an Internet Draft^[21], which proposes a new SMIMEA DNS RR. The purpose of the SMIMEA RR is to associate certificates with DNS domain names.

S/MIME messages often contain certificates that can assist in authenticating the message sender and can be used in encrypting messages sent in reply. This feature requires that the receiving MUA validate the certificate associated with the purported sender. SMIMEA RRs can provide a secure means of doing this validation.

In essence, the SMIMEA RR will have the same format and content as the TLSA RR, with the same functionality. The difference is that it is geared to the needs of MUAs in dealing with domain names as specified in e-mail addresses in the message body, rather than domain names specified in the outer SMTP envelope.

Sender Policy Framework

Sender Policy Framework (SPF) is the standardized way for a sending domain to specify a list of MTAs that are authorized to send on behalf of the domain. The problem that SPF addresses is the following: with the current e-mail infrastructure, any host can use any domain name for each of the various identifiers in the mail header, not just the domain name where the host is located.

Two major drawbacks of this freedom follow:

- It is a major obstacle to reducing *Unsolicited Bulk E-mail* (UBE), also known as *spam*. It makes it difficult for mail handlers to filter out e-mails on the basis of known UBE sources.
- *Administrative Management Domains* (ADMDs) are understandably concerned about the ease with which other entities can use their domain names, often with malicious intent.

However, a basic limitation of SPF is that it forces mail to follow a specific path and breaks when legitimate mail deviates from this path, such as a message that goes through a mailing list.

RFC 7208 defines the SPF^[31]. It provides a protocol by which ADMDs can authorize hosts to use their domain names in the **MAIL FROM** or **HELO** identities. (It is worth noting that this domain name is the return address for error messages, rather than being required to be the same as the author's address.) Compliant ADMDs publish SPF records in the DNS specifying which hosts are permitted to use their names, and compliant mail receivers use the published SPF records to test the authorization of sending MTAs using a given **HELO** or **MAIL FROM** identity during a mail transaction.

SPF works by checking a neighboring, upstream client MTA IP address against the policy encoded in any SPF record found at the sending domain. The sending domain is the domain used in the SMTP connection, not the domain indicated in the Author From field in the message header as displayed in the MUA. Thus SPF checks can be applied before the message content is received from the sender.

Figure 5 on the following page is an example in which SPF would come into play. Assume that the sender's IP address is **192.168.0.1**. The message arrives from the MTA with domain **mta.example.net**. The sender uses the envelope **MAIL FROM** tag of **alice@example.org**, indicating that the message originates in the **example.org** domain. But the message header specifies **alice.sender@example.net**. The receiver uses SPF to query for the SPF RR that corresponds to **example.org** to check if the IP address **192.168.0.1** is listed as a valid sender, and then takes appropriate action based on the results of checking the RR.

A sending domain needs to identify all the senders for a given domain and add that information into the DNS as a separate resource record. Next, the sending domain encodes the appropriate policy for each sender using the SPF syntax. The encoding is done in a TXT DNS resource record as a list of mechanisms and modifiers. Mechanisms are used to define an IP address or range of addresses to be matched, and modifiers indicate the policy for a given match. The SPF syntax is fairly complex and can express complex relationships between senders. For more details, see RFC 7208.

Figure 5: Example in Which SMTP
Envelope Header Does Not
Match Message Header

```

S: 220 foo.com Simple Mail Transfer Service Ready
C: HELO mta.example.net
S: 250 OK
C: MAIL FROM:<alice@example.org>
S: 250 OK
C: RCPT TO:<Jones@foo.com>
S: 250 OK
C: DATA
S: 354 Start mail input; end with <CRLF>.<CRLF>
C: To: bob@foo.com
C: From: alice.sender@example.net
C: Date: Today
C: Subject: Meeting Today
. . .

```

If SPF is implemented at a receiver, the SPF entity uses the SMTP envelope **MAIL FROM:** address domain and the IP address of the sender to query an SPF TXT RR. The SPF checks can be started before the body of the e-mail message is received, possibly resulting in blocking the transmission of the e-mail content. Alternatively, the entire message can be absorbed and buffered until all the checks are finished. In either case, checks must be completed before the mail message is sent to the end user's inbox.

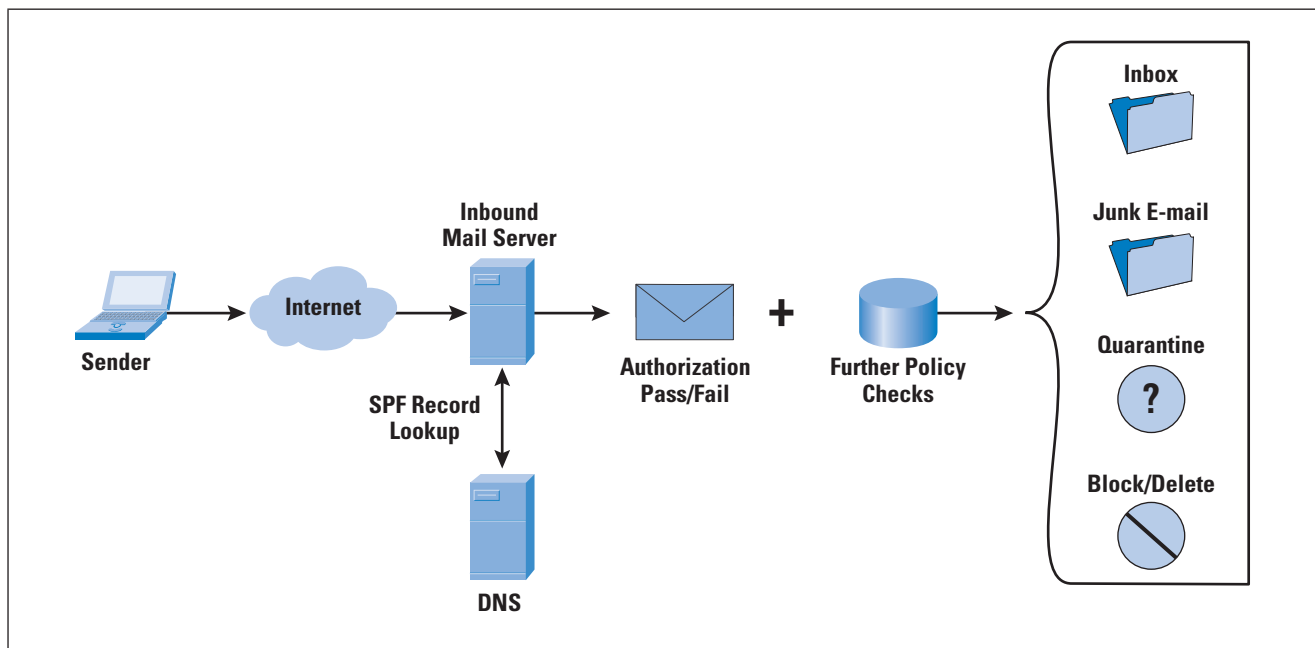
The checking involves the following rules:

1. If no SPF TXT RR is returned, the default behavior is to accept the message.
2. If the SPF TXT RR has formatting errors, the default behavior is to accept the message.
3. Otherwise the mechanisms and modifiers in the RR are used to determine disposition of the e-mail message.

With respect to SPF alone, to say in step 1, preceding, that the default behavior is to accept the message is correct. However, it should be noted that SPF is usually working within a mixture of anti-abuse tools and the aggregate filtering engine typically does not accept a message based on the results of only one of its tools, such as SPF.

Figure 6 illustrates SPF operation. As of 2016, more than 27% of all Internet domains implement SPF^[32].

Figure 6: Sender Policy Framework Operation



DKIM

DomainKeys Identified Mail (DKIM) permits a person, role, or organization that owns the signing domain to claim some responsibility for a message by associating the domain with the message^[33]. The domain can be an author's organization, an operational relay, or one of their agents. DKIM separates the question of the identity of the signer of the message from the purported author of the message. Assertion of responsibility is validated through a cryptographic signature and by querying the signer's domain directly to retrieve the appropriate public key.

The qualifier *some* in the first sentence of the preceding paragraph is important. In particular, the text directly "covered" by the signature is not vetted for authenticity.

Message recipients (or agents acting in their behalf) can verify the signature by querying the signer's domain directly to retrieve the appropriate public key and thereby can confirm that the message was attested to by a party in possession of the private key for the signing domain. DKIM is an Internet Standard defined in RFC 6376^[34]. DKIM has been widely adopted by a range of e-mail providers, including corporations, government agencies, Gmail, Yahoo, and many *Internet Service Providers* (ISPs). As of 2016, an estimated 40% of Internet sites deploy DKIM^[35].

An *Administrative Unit* (AU) is that portion of the path of an e-mail message that is under a single administration. DKIM focuses primarily on attackers located outside of the AUs of the claimed originator and the recipients, indirectly, by creating a verifiable signature of valid mail from the administrative unit.

It is with these external AUs that the trust relationships required for authenticated message submission may not exist and do not scale adequately to be practical. Conversely, within these AUs, there are other mechanisms (such as authenticated message submission) that are easier to deploy and more likely to be used than DKIM. External bad actors are usually attempting to exploit the “any-to-any” nature of e-mail that motivates most recipient MTAs to accept messages from anywhere for delivery to their local domain. They may generate messages without signatures, with incorrect signatures, or with correct signatures from domains with little traceability. They may also pose as mailing lists, greeting cards, or other agents that legitimately send or resend messages on behalf of others.

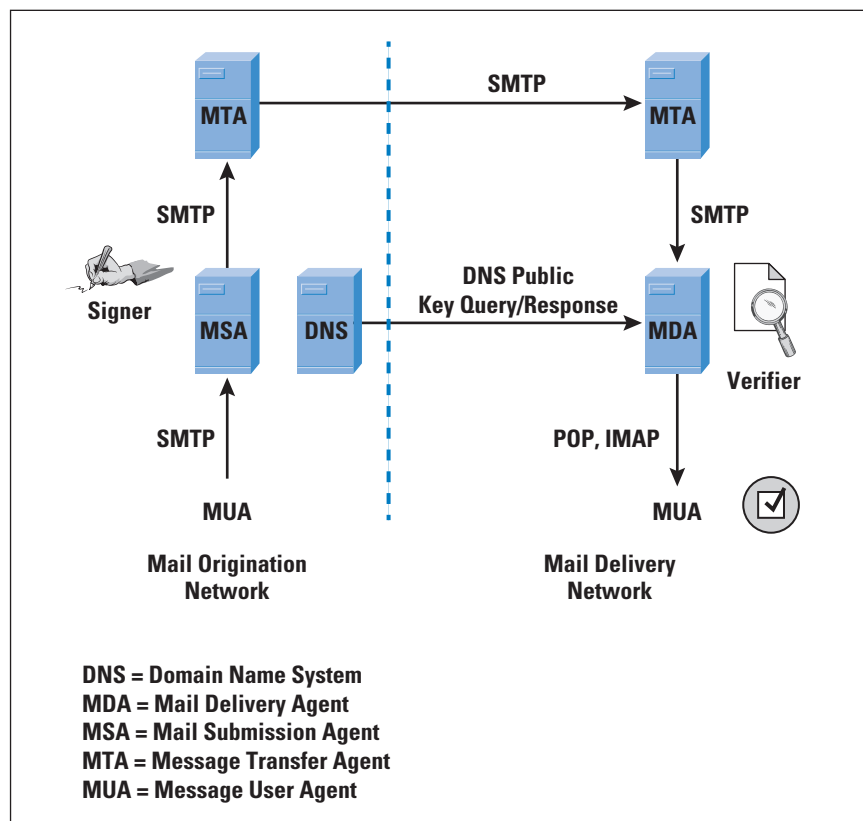
DKIM is designed to provide an e-mail authentication technique that is transparent to the end user. In essence, a user’s e-mail message is signed by a private key of the administrative domain from which the e-mail originates. The signature covers none, some, or all of the content of the message and some of the e-mail message headers.

Note that the signature is not validating any of what is signed, as digital signatures usually do. Rather, the choice of what to cover is meant as a means of gluing the **d=domain name** to the overall message in a way that is difficult to spoof. At the receiving end, the Message Delivery Agent can access the corresponding public key via a DNS and verify the signature, thus authenticating that the message comes from the claimed administrative domain. Thus, DKIM allows an enterprise to vouch for an e-mail message sent from a domain it does not control. This approach differs from that of S/MIME, which uses the originator’s private key to sign the content of the message. The motivation for DKIM is based on the following reasoning:

- S/MIME depends on both the sending and receiving users employing S/MIME. For almost all users, the bulk of incoming mail does not use S/MIME, and the bulk of the mail the user wants to send is to recipients not using S/MIME.
- S/MIME signs only the message content. Thus, RFC 5322^[36] header information concerning origin can be compromised.
- DKIM is not implemented in client programs (MUAs) and is therefore transparent to the user; the user doesn’t need to take any action.
- DKIM can be configured to apply to all mail from cooperating domains.
- DKIM allows good senders to prove that they did send a particular message and to prevent forgers from forging the DKIM signature.

Figure 7 is a simple example of the operation of DKIM. We begin with a message generated by a user and transmitted into the *Message Handling Service* (MHS) to an MSA that is within the user’s administrative domain. An e-mail message is generated by an e-mail client program.

Figure 7: Simple Example of DKIM Deployment



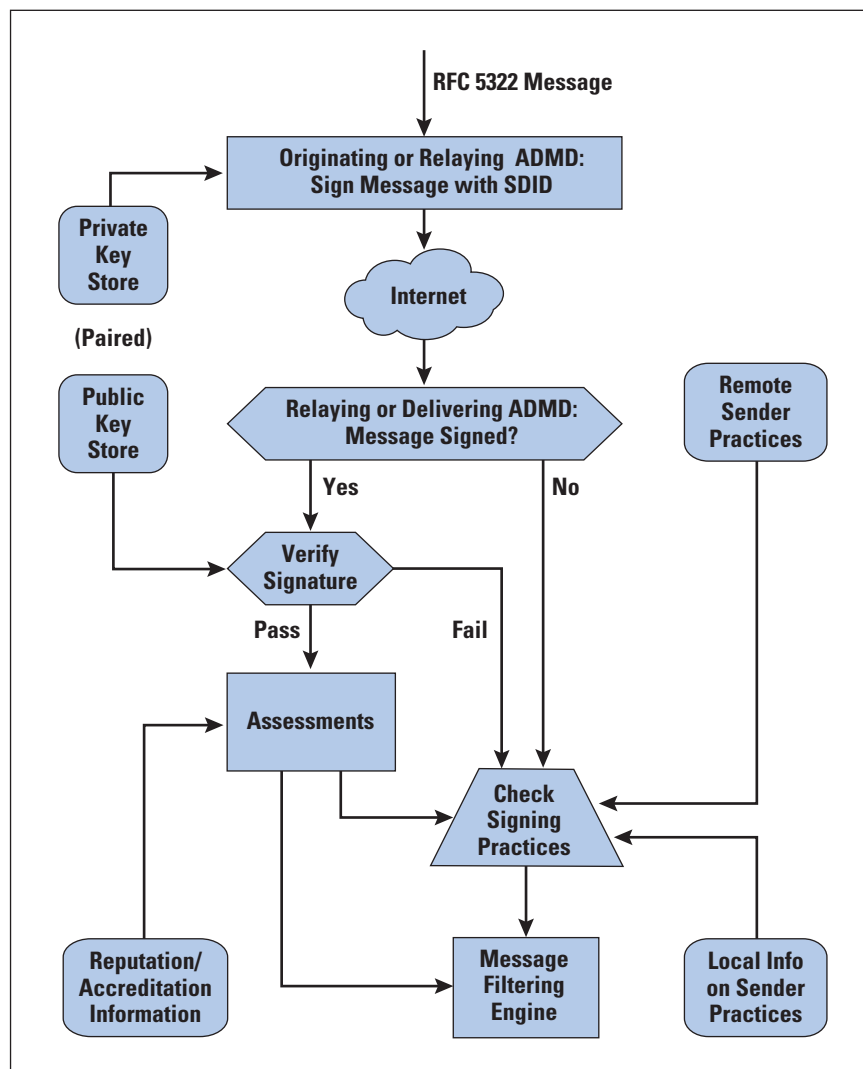
The content of the message, plus selected RFC 5322 headers, is signed by the e-mail provider using the provider's private key. The signer is associated with a domain, which could be a corporate local network, an ISP, or a public e-mail facility such as Gmail. The signed message then passes through the Internet via a sequence of MTAs. At the destination, the MDA retrieves the public key for the incoming signature and verifies the signature before passing the message on to the destination e-mail client. The default signing algorithm is RSA with SHA-256. RSA with SHA-1 also may be used.

Figure 8 on the following page, from RFC 5585^[37], provides a more detailed look at the elements of DKIM operation. Basic message processing is divided between a signing *Administrative Management Domain* (ADMD) and a verifying ADMD. At its simplest, this processing is between the originating ADMD and the delivering ADMD, but it can involve other ADMDs in the handling path.

Signing is performed by an authorized module within the signing ADMD and uses private information from a Key Store. Within the originating ADMD, this signing might be performed by the MUA, MSA, or an MTA. Verifying is performed by an authorized module within the verifying ADMD. Within a delivering ADMD, verifying might be performed by an MTA, MDA, or MUA. The module verifies the signature or determines whether a particular signature was required.

Verifying the signature uses public information from the Key Store. If the signature passes, reputation information is used to assess the signer and that information is passed to the message filtering system.

Figure 8: DKIM Functional Flow



If the signature fails or there is no signature using the author's domain, information about signing practices related to the author can be retrieved remotely and/or locally, and that information is passed to the message filtering system. For example, if the sender (for example, Gmail) uses DKIM but no DKIM signature is present, then the message may be considered fraudulent.

The signature is inserted into the RFC 5322 message as an additional header field, starting with the keyword Dkim-Signature. You can view examples from your own incoming mail by using the "View Long Headers (or similar wording) option for an incoming message.

Before a message is signed, a process known as *canonicalization* is performed on both the header and body of the RFC 5322 message. Canonicalization is necessary to deal with the possibility of minor changes in the message made en route, including character encoding, treatment of trailing white space in message lines, and the “folding” and “unfolding” of header lines. The intent of canonicalization is to make a minimal transformation of the message (for the purpose of signing; the message itself is not changed, so the canonicalization must be performed again by the verifier) that will give it its best chance of producing the same canonical value at the receiving end. DKIM defines two header canonicalization algorithms (“simple” and “relaxed”) and two for the body (with the same names). The simple algorithm tolerates almost no modification, while the relaxed tolerates common modifications.

DMARC

Domain-Based Message Authentication, Reporting, and Conformance (DMARC), defined in RFC 7489^[38], allows e-mail senders to specify policy on how their mail should be handled, the types of reports that receivers can send back, and the frequency of those reports.

DMARC works with SPF and DKIM. SPF enables senders to advise receivers, via DNS, whether mail purporting to come from the sender is valid, and whether it should be delivered, flagged, or discarded. However, neither SPF nor DKIM includes a mechanism to tell receivers if SPF or DKIM is in use, nor do they have a feedback mechanism to inform senders of the effectiveness of the anti-spam techniques. For example, if a message arrives at a receiver without a DKIM signature, DKIM provides no mechanism to allow the receiver to learn if the message is authentic but was sent from a sender that did not implement DKIM, or if the message is a spoof. In essence, DMARC addresses these issues by indicating whether SPF and/or DKIM will be used, what a receiver should do when they aren’t, and how receivers should report aggregate results for the domain.

DKIM, SPF, and DMARC authenticate various aspects of an individual message. DKIM authenticates the domain that affixed a signature to the message. SPF focuses on the SMTP envelope, defined in RFC 5321^[39]. It can authenticate either the domain that appears in the **MAIL FROM** portion of the SMTP envelope or the **HELO** domain, or both. These domains may be different, and they are typically not visible to the end user.

DMARC authentication deals with the From domain in the message header, as defined in RFC 5322. This field is used as the central identity of the DMARC mechanism because it is a required message header field and therefore guaranteed to be present in compliant messages, and most MUAs represent the RFC 5322 From field as the originator of the message and render some or all of this content of the header field to end users. The e-mail address in this field is the one used by end users to identify the source of the message and therefore is a prime target for abuse.

DMARC requires that the From address match (be aligned with) an Authenticated Identifier from DKIM or SPF. In the case of DKIM, the match is made between the DKIM signing domain and the From domain. In the case of SPF, the match is between the SPF-authenticated domain and the From domain.

A mail sender that uses DMARC must also use SPF or DKIM, or both. The sender posts a DMARC policy in the DNS that advises receivers on how to treat messages that purport to originate from the sender's domain. The policy is in the form of a DNS TXT resource record associated with the sender's domain name. The sender also needs to establish e-mail addresses to receive aggregate and forensic reports. Because these e-mail addresses are published unencrypted in the DNS TXT RR, they are easily discovered, leaving the poster subject to unsolicited bulk e-mail. Thus, the poster of the DNS TXT RR needs to employ some kind of abuse countermeasures.

Similar to SPF and DKIM, the DMARC policy in the TXT RR is encoded in a series of tag=value pairs separated by semicolons. Once the DMARC RR is posted, messages from the sender are typically processed as follows:

1. The domain owner constructs an SPF policy and publishes it in its DNS database. The domain owner also configures its system for DKIM signing. Finally, the domain owner publishes via the DNS a DMARC message-handling policy.
2. The author generates a message and hands the message to the domain owner's designated mail submission service.
3. The submission service passes relevant details to the DKIM signing module in order to generate a DKIM signature to be applied to the message.
4. The submission service relays the now-signed message to its designated transport service for routing to its intended recipient(s).

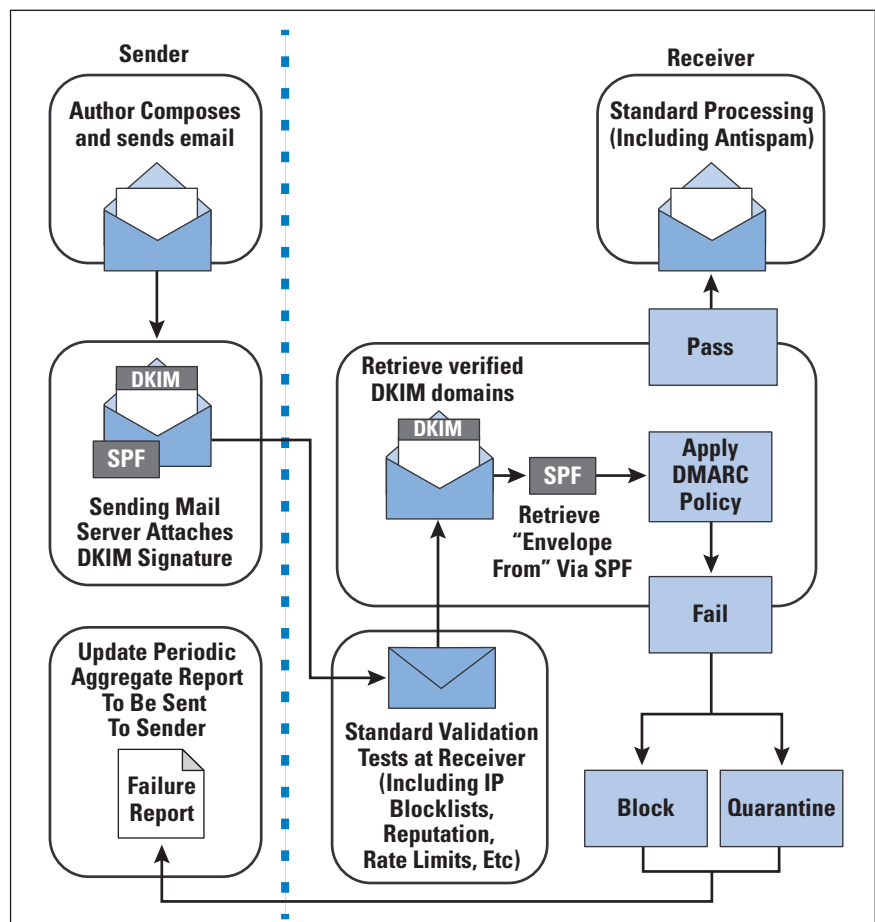
A message generated on the sender side may pass through other relays but eventually arrives at a receiver's transport service. The typical processing order for DMARC on the receiving side follows:

1. The receiver performs standard validation tests, such as checking against IP blocklists and domain reputation lists, as well as enforcing rate limits from a particular source.
2. The receiver extracts the RFC 5322 From address from the message. This address must be a single, valid address or else the mail is refused as an error.
3. The receiver queries for the DMARC DNS record based on the sending domain. If none exists, DMARC processing is terminated.
4. The receiver performs DKIM signature checks. If more than one DKIM signature exists in the message, one must verify.

5. The receiver queries for the SPF record of the sending domain and performs SPF validation checks.
6. The receiver conducts Identifier Alignment checks between the RFC 5321 From and the results of the SPF and DKIM records (if present).
7. The results of these steps are passed to the DMARC module along with the Author's domain. The DMARC module attempts to retrieve a policy from the DNS for that domain. If none is found, the DMARC module determines the organizational domain and repeats the attempt to retrieve a policy from the DNS.
8. If a policy is found, it is combined with the Author's domain and the SPF and DKIM results to produce a DMARC policy result (a "pass" or "fail") and can optionally cause one of two kinds of reports to be generated.
9. Recipient transport service either delivers the message to the recipient inbox or takes other local policy action based on the DMARC result.
10. When requested, Recipient transport service collects data from the message delivery session to be used in providing feedback.

Figure 9, based on one at [DMARC.org](https://dmarc.org), summarizes the sending and receiving functional flow.

Figure 9: DMARC Functional Flow



DMARC reporting provides the senders feedback on their SPF, DKIM, Identifier Alignment, and message disposition policies, which enables the sender to make these policies more effective. Two types of reports are sent: *Aggregate Reports* and *Forensic Reports*.

Aggregate Reports are sent by receivers periodically and include aggregate figures for successful and unsuccessful message authentications, including:

- The sender's DMARC policy for that interval
- The message disposition by the receiver (that is, delivered, quarantined, rejected)
- SPF result for a given SPF identifier
- DKIM result for a given DKIM identifier
- Whether identifiers are in alignment or not
- Results classified by sender subdomain
- The sending and receiving domain pair
- The policy applied, and whether it is different from the policy requested
- The number of successful authentications
- Totals for all messages received

This information enables the sender to identify gaps in e-mail infrastructure and policy. SP 800-177 recommends that a sending domain begin by setting a DMARC policy of **p=none**, so that the ultimate disposition of a message that fails some check is determined by the receiver's local policy. As DMARC aggregate reports are collected, the sender will have a quantitatively better assessment of the extent to which the sender's e-mail is authenticated by outside receivers, and will be able to set a policy of **p=reject**, indicating that any message that fails the SPF, DKIM, and alignment checks really should be rejected. From their own traffic analysis, receivers can determine whether a sender's **p=reject** policy is sufficiently trustworthy to act on.

A *Forensic Report* helps senders refine the component SPF and DKIM mechanisms as well as alerting them that their domain is being used as part of a phishing/spam campaign. Forensic reports are similar in format to aggregation reports, with these changes:

- Receivers include as much of the message and message header as is reasonable to allow the domain to investigate the failure. Add an *Identity-Alignment* field, with DKIM and SPF DMARC-method fields as appropriate.
- Optionally add a *Delivery-Result* field. Add DKIM Domain, DKIM Identity, and DKIM selector fields, if the message was DKIM signed. Optionally also add DKIM Canonical header and body fields.
- Add an additional DMARC authentication failure type, for use when some authentication mechanisms fail to produce aligned identifiers.

Since its introduction, DMARC has seen rapid acceptance. Thousands of companies use it to prevent billions of messages fraudulently using their Internet domains from reaching inboxes, thereby protecting their customers and employees from phishing and other abuse. Recently, two of the largest mailbox providers in the world—Google and Yahoo—have announced that they are extending that protection to cover more of their Internet domains^[40].

Summary

The IETF has developed a suite of protocols that provide comprehensive Internet e-mail security. Many of these protocols have been widely deployed, and the entire suite is recommended by NIST.

Acknowledgment

The author would like to express his gratitude to the reviewer for the many detailed and helpful comments.

References

- [1] National Institute of Standards and Technology, “Trustworthy Email,” NIST Special Publication 800-177, September 2016.
- [2] Dave Crocker, “Internet Mail Architecture,” RFC 5598, July 2009.
- [3] Crypto Portal, “Cryptography with Cryptool: Practical Introduction to Cryptography and Cryptanalysis,” August 2010.
<https://www.cryptool.org/images/ctl/presentations/CrypToolPresentation-en.pdf>
- [4] National Institute of Standards and Technology, “Introduction to Public Key Technology and the Federal PKI Infrastructure,” NIST Special Publication 800-32, February 2001.
- [5] Paul Hoffman, “SMTP Service Extension for Secure SMTP over Transport Layer Security,” RFC 3207, February 2002.

- [6] ZDNet, “Google, Microsoft, Yahoo: We want to stop e-mail snooping by fixing these encryption flaws,” March 21, 2016.
<http://www.zdnet.com/article/google-microsoft-yahoo-we-want-to-stop-e-mail-snooping-by-fixing-these-encryption-flaws/#!>
- [7] Facebook, “The Current State of SMTP STARTTLS Deployment,” May 13, 2014.
<https://www.facebook.com/notes/protect-the-graph/the-current-state-of-smtp-starttls-deployment/1453015901605223/>
- [8] William Stallings, “SSL: Foundation for Web Security,” *The Internet Protocol Journal*, Volume 1, No. 1, June 1998.
- [9] Andrea Peterson, “Facebook’s security chief on the Snowden effect, the Messenger app backlash and staying optimistic,” *The Washington Post*, August 12, 2014.
- [10] David Cohen, “Facebook: 95% of Notification Emails Encrypted Thanks to Providers’ STARTTLS Deployment,” *AdWeek*, August 19, 2014.
- [11] Marshall Rose and David Strom, “Secure E-Mail: Problems, Standards, and Prospects,” *The Internet Protocol Journal*, Volume 2, No. 1, March 1999.
- [12] Sean Turner and Blake Ramsdell, “Secure/Multipurpose Internet Mail Extensions (S/MIME) Version 3.2 Certificate Handling,” RFC 5750, January 2010.
- [13] Sean Turner and Blake Ramsdell, “Secure/Multipurpose Internet Mail Extensions (S/MIME) Version 3.2 Message Specification,” RFC 5751, January 2010.
- [14] Paul Hoffman, “Examples of S/MIME Messages,” RFC 4134, July 2005.
- [15] Paul Hoffman, “Enhanced Security Services for S/MIME,” RFC 2634, June 1999.
- [16] Russ Housley, “Cryptographic Message Syntax (CMS),” RFC 5652, September 2009.
- [17] Russ Housley, “Cryptographic Message Syntax (CMS) Algorithms,” RFC 3370, August 2002.
- [18] Jim Schaad and Sean Turner, “Multiple Signatures in S/MIME,” RFC 5752, January 2010.
- [19] Sandy Murphy, Jim Galvin, Steve Crocker, and Ned Freed, “Security Multiparts for MIME: Multipart/Signed and Multipart/Encrypted,” RFC 1847, October 1995.

- [20] Paul Hoffman and Jakob Schlyter, “Using Secure DNS to Associate Certificates with Domain Names for S/MIME,” Internet Draft, work in progress, **draft-ietf-dane-smime-10**, February 24, 2016.
- [21] Philip Zimmermann, Derek Atkins, and William Stallings, “PGP Message Exchange Formats,” RFC 1991, August 1996.
- [22] Hal Finney, Lutz Donnerhacke, Jon Callas, Rodney Thayer, and David Shaw, “OpenPGP Message Format,” RFC 4880, November 2007.
- [23] Dave Del Torto, Michael Elkins, Raph Levien, and Thomas Roessler, “MIME Security with OpenPGP,” RFC 3156, August 2001.
- [24] A. Whitten and J. Tygar, “Why Johnny can’t encrypt: a usability evaluation of PGP 5.0,” *Proceedings of the 8th conference on USENIX Security Symposium - Volume 8 (SSYM’99)*, 1999.
- [25] Matthew Green, “What’s the Matter with PGP?” Cryptography Engineering Blog, August 13, 2014.
<http://blog.cryptographyengineering.com/2014/08/whats-matter-with-pgp.html>
- [26] Miek Gieben, “DNSSEC: The Protocol, Deployment, and a Bit of Development,” *The Internet Protocol Journal*, Volume 7, No. 2, June 2004.
- [27] Richard Barnes, “Use Cases and Requirements for DNS-Based Authentication of Named Entities (DANE),” RFC 6394, October 2011.
- [28] Richard Barnes, “Let the Names Speak for Themselves: Improving Domain Name Authentication with DNSSEC and DANE,” *The Internet Protocol Journal*, Volume 15, No. 1, March 2012.
- [29] Jakob Schlyter and Paul Hoffman, “The DNS-Based Authentication of Named Entities (DANE) Transport Layer Security (TLS) Protocol: TLSA,” RFC 6698, August 2012.
- [30] Viktor Dukhovni and Wesley Hardaker, “SMTP Security via Opportunistic DNS-Based Authentication of Named Entities (DANE) Transport Layer Security (TLS),” RFC 7672, October 2015.
- [31] Scott Kitterman, “Sender Policy Framework (SPF) for Authorizing Use of Domains in Email, Version 1,” RFC 7208, April 2014.
- [32] “SPF-all Domain Survey,” <http://spf-all.com/stats.html>

- [33] Barry Leiba and Jim Fenton, “DomainKeys Identified Mail (DKIM): Using Digital Signatures for Domain Verification,” CEAS 2007—*Fourth Conference on E-mail and Anti-Spam*, August 2–3, 2007.
- [34] Murray Kucherawy, Dave Crocker, and Tony Hansen, “DomainKeys Identified Mail (DKIM) Signatures,” RFC 6376, September 2011.
- [35] “Global DKIM Deployment,”
<https://eggert.org/meter/dkim>
- [36] Peter W. Resnick, “Internet Message Format,” RFC 5322, October 2008.
- [37] Tony Hansen, Dave Crocker, and Phillip Hallam-Baker, “DomainKeys Identified Mail (DKIM) Service Overview,” RFC 5585, July 2009.
- [38] Murray Kucherawy and Elizabeth Zwicky, “Domain-based Message Authentication, Reporting, and Conformance (DMARC),” RFC 7489, March 2015.
- [39] John C. Klensin, “Simple Mail Transfer Protocol,” RFC 5321, October 2008.
- [40] Dmarc.org, “Global Mailbox Providers Deploying DMARC to Protect Users,” Dmarc Press Release, October 19, 2015.

WILLIAM STALLINGS is an independent consultant and author of numerous books on security, computer networking, and computer architecture. His latest book is *Cryptography and Network Security* (Pearson, 2016). He maintains a computer science resource site for computer science students and professionals at **ComputerScienceStudent.com** and is on the editorial board of *Cryptologia*. He has a Ph.D. in computer science from M.I.T. He can be reached at **ws@shore.net**

Cloudy-Eyed: Complexity and Reality with Software-Defined Networks

by Russ White and Shawn Zandi, LinkedIn

Software-Defined Networks (SDN) are promoted as a way to eliminate the complexity of distributed control planes, increase network responsiveness to specific applications and business requirements, and reduce operational and equipment cost. If this description sounds like the classic “too good to be true” situation, that’s because it might just be. Just like you can’t build a database that has ideal consistency, accessibility, and partitionability, you can’t build a cheap network with optimal routing and minimal control-plane state. It’s just a reality of the complexity built into the physical shape of the universe that everything has a tradeoff—*cheap, fast, and high quality, choose two*.

When we reach the top of the SDN hype cycle, what will our options be? Perhaps the best place to start in answering this question is by considering why the “big promise” of SDN hasn’t been really successful in the real world.

Defining SDN: Then and Now

To really understand the hype and promise of SDNs, it’s important to go back to the beginning and consider what the original promise really was. There were originally three crucial elements to the SDN story.

First, SDNs were supposed to remove the intelligence from distributed control planes, replacing them with the centralized calculation of network paths in a controller. While an individual autonomous router has only a localized view of network conditions, a centralized controller can gain a more global view. A global view would allow the controller to more efficiently manage and direct traffic through the network in a way that improves both the efficiency of the network and the performance of applications running across the network.

Second, SDNs were supposed to provide a much more granular level of control—down to the flow level. This added level of control would enable much better policy control in various ways, including the discovery and direction of elephant flows, quality of service on a per-application/per-user basis, and other options.

Third, SDN would enable the network to be programmable, thereby reducing operational costs, enabling a more lean/agile view of the network, and allowing applications to interact directly with the network.

The definition of SDN has changed over the years, broadening so that it now includes just about any network technology that allows programmatic access to information about and control over the network.

An SDN, in more recent terms, seems to include everything from the ability of an application to schedule bandwidth (which is a rather more complicated problem than it seems) to gaining better telemetry data. The centralized controller, flow-based forwarding, and commoditization of hardware are still in scope, but they appear to be mixed in with a much more limited view of the “core components” of the SDN message. Why has the concept of the SDN changed across time?

It’s possible to argue that this definitional change is just a matter of the marketing departments at a wide array of vendors grabbing hold of the term, but there seems to be something deeper here. Perhaps the “something deeper” is the original ideals have proven more difficult to achieve than were first thought. A short overview of the challenges of deploying the original SDN ideal might be useful in understanding the historical flow of these changes. Three larger areas are considered in the following sections: centralizing the calculation of network paths, flow-based forwarding, and network programmability.

Centralizing the Calculation of Network Paths

Distributed control planes, such as *Intermediate System-to-Intermediate System* (IS-IS) and *Border Gateway Protocol* (BGP), are often (rightly) seen as one of the most complex components of a network. In fact, entire networks are designed around the operation of these routing protocols, including the consideration of topics like:

- Splitting up failure domains through information hiding
- Managing complex policies through communities, tags, and metrics
- Choosing topologies based on fast convergence characteristics
- The interaction of multiple distributed control planes running on a single network

Further, in order to support the complex processing and data handling of distributed control planes, network devices are typically large, expensive devices, with fast processors and large memory pools. In particular, as the need for policy-driven path selection (which generally means choosing a path through the network that is less optimal than the shortest path from a metrics perspective, but more optimal from a network usage or quality-of-service perspective) increased, the processing power and memory requirements of individual routers ramped up.

If distributed control planes could be eliminated and replaced by a controller (or set of controllers), the complexity of each forwarding device could be reduced dramatically, because the jobs of discovering the local topology and calculating the best path per destination would be offloaded from the individual boxes, and pushed onto the controller. By removing this processing from the routers, small, cheap, lightweight forwarding-only devices could be used instead of the traditional router.

Hence the world could move to *white-box* devices that would be available off the shelf and require little configuration.

Complexity, however, is not so easy to slay. The centralized controller approach presents numerous problems that will, most likely, forever limit it in scale and scope to something smaller than what was originally envisioned, such as two or three controllers providing forwarding information for tens of thousands of switches running at scale. Some of these problems include the relationship between centralized computation and reactive control planes, remote reactions to local topology and reachability changes, and what can fairly be described as the halo effect around software engineering.

Centralized Control and Reactive Forwarding

Distributed control planes, such as IS-IS, are proactive in their discovery of topology and reachability. Before the first packet is transmitted across the network, the routing protocol must discover a set of loop-free paths that can reach every destination in the network. Since this discovery and calculation process typically involves flooding, processing, and managing a lot of information, distributed control planes often rely on information hiding through aggregation to manage the amount and speed of state being carried in the protocol. For instance, in IS-IS intermediate systems in the level 2 flooding domain don't have any information about the topology of the outlying level 1 flooding domains. In a similar way, intermediate systems in a level 1 flooding domain know only about the topology within the flooding domain and which intermediate systems are connected to the level 2 flooding domain.

When the calculation of routes is centralized, there must still be some form of information hiding to scale the control plane. Instead of aggregation at specific topological points in the network, SDN control planes most often opt for moving to a reactive control plane, meaning the forwarding devices discover reachability information only when they receive the first packet in a flow. While this does reduce the amount of forwarding state in any particular device, it also has many drawbacks.

Specifically, reactive control planes disconnect the apparent state of the network from the perspective of any attached device from the actual state of the network. From the host's perspective, the network is up, and therefore there is a path to most destinations that begins with the first packet in a flow. In a reactive control plane, however, there is some amount of lag between the first packet in a flow being transmitted and the path actually being available. One objection to this observation is that the *Domain Name System* (DNS) is also reactive in much the same way. However, end devices generally participate in the DNS system, and hence know the state of their ability to forward in terms of name resolution.

Further, while it's always possible for the network to change state in the middle of a flow being transmitted, reactive control planes suffer from a wider set of causes for these changes. This situation is always true, of course, but while proactive control planes treat a disconnect between apparent and actual states as an error condition to be resolved, reactive control planes treat such a disconnect as a normal state of affairs. In a larger sense, disconnects between the actual and perceived states of the network are seen by attached devices as network instability; the stronger the disconnect, the more unstable the network appears to be. This condition can have an adverse effect on applications and host behavior. Local cache timeouts, cache failures, and other problems need to be included in the more general topology changes and problems common to distributed control planes for path failures.

Centralized Control and Fast Reaction to Changes in the Network

Centralized control planes disconnect local state from recalculation of the best path. If a local node or link fails, information about the state change must be transmitted to a remote device (the controller), which must recalculate a new set of paths, and then distribute those paths throughout the network. These operations can be made very quickly using techniques such as calculating and installing a backup route, but there is no simple way for a centralized controller to react more quickly, and with less chance of an unanticipated failure mode than with a distributed control plane.

The centralized/decentralized decision isn't necessarily a *better-versus-worse* decision, it's just a *different* decision with a *distinct set of tradeoffs*. Each path has its own complexities and problems to address; no set of problems seems to be much less complex to solve than any other set in this case.

The Halo of Software Development

Distributed control planes, as mentioned previously, are very complex, and they require a lot of configuration to deploy, design, troubleshoot, and manage. It seems simpler, in many ways, to just replace all the people who do this configuring, troubleshooting, and managing, with a small team of coders who can build and maintain a controller. The code would be simpler because it's all "in one place," and can be more customized to fit a particular business environment. The reality, however, is far different.

But a single controller simply won't do when it comes to scaling out a network. Even if you could run a network of thousands of routers with a single controller, it goes against every foundational concept of solid system design to do so. There must be at least two controllers, in topologically diverse locations within the network, to provide redundancy in support of overall system resilience. Moving from one controller to two inevitably means providing some way to distribute reachability and policy information between the controllers.

Ultimately, then, a distributed control plane must be built to allow communication between the controllers.^[1] Couldn't this distribution just be some standard distributed database? It could, but there's a difference between distributing a database and distributing the meaning contained in the database. To distribute the meaning, you must have an agreed-on format, encoding, and other things. If you examine existing distributed routing protocol specifications, you'll find they spend a lot of time describing not only how to carry information, but also how to specify what sort of information is being carried, and consistent ways to interpret and use that information. To make multiple controller configurations successful (especially across multiple controller vendors), either it all will need to be rebuilt in an inter-controller protocol, or—perhaps simpler—the controllers could just use an existing routing protocol. Regardless of the solution chosen, the problems involved in a distributed control plane haven't been removed from the network, they've just been moved to someplace else in the network.

Further, the distributed protocols the SDN controller is designed to replace are really just other software. The complexity in these protocols comes from the propensity of engineers to push functionality into them to address an ever-expanding array of use cases. As time passes and the larger (or perhaps more obvious) use cases are handled, protocol developers chase smaller problems, finally reaching into large amounts of code for what is really a set of corner cases. But moving the development of the control plane from one place in the network to another place in the network isn't going to solve this problem—the process of accretion of new features and an ever-larger code base and inter-controller protocol specifications to support an ever-increasing set of use cases is going to remain the same.

Flow-Based Forwarding

Standard IP headers contain at least five fields of interest to network devices:

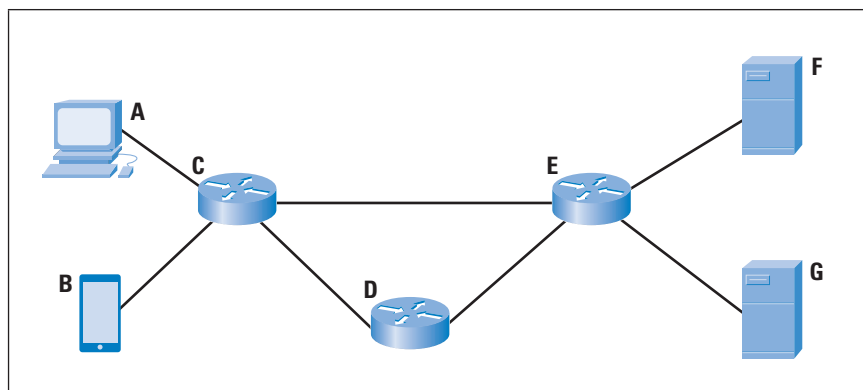
- Source IP address
- Source port number
- Destination IP address
- Destination port number
- Protocol number (or identifier)

Some transport protocols also include more information that might be of interest, such as the *Transmission Control Protocol* (TCP) socket number, which can indicate a particular application, and information about the expected quality of service for this packet (and ultimately flow). Much of this information (and more) could be used to forward traffic into multiple paths through the network based on policy and available bandwidth. However, existing routing protocols are designed to provide reachability, and hence forwarding, information based on the destination address only.

There has long been a desire to forward traffic based on much more than the destination address, so that individual applications can be independently routed through the network, and information other than the destination address can be used to deny access to specific network resources. These requirements have led to a string of work in the area of accounting for the IP source, at the least, when making forwarding decisions.^[2]

Figure 1 provides an example.

Figure 1: Flow-Based Forwarding Control Example



In this example, Host A is sending a large file to Server F, while the user at Host B, a small handheld mobile device, is participating in a video conference through Server G. Assuming the two server addresses are shared among numerous different services, destination-based addresses cannot be used to differentiate between the large file transfer and the video conference. In this case, if the network administrator knows about the file transfer, the source addresses of A and B, along with the source and destination protocol information, can be used to differentiate the two traffic streams. This solution would allow the file-transfer traffic to be directed along the [C,D,E] link, while the video conferencing traffic would be directed along the [C,E] link. This traffic separation can be used to allow the video conference traffic to pass along links that aren't being heavily used by the file transfer.

Flow-based forwarding, however, presents many problems.

First, the amount of control-plane state required to forward every flow in a large network individually would be far beyond reasonable. The problem is not just the number of flows, but also the flow setup rate. To put this idea in real terms, if there are 10,000 hosts such as A and B in the illustration, and each attempts to open a different website, and each website requires 20 TCP connections, the network is required to calculate and install 2,000,000 flows in a matter of seconds. Few controllers could handle flow setup rates at this level.^[3]

Second, the hardware costs of implementing such a scheme would be very high. The amount of flow state required in each device would be incredibly large—larger than most commodity hardware can support.

In addition, the cost of examining the full header on each packet at each hop in the network to achieve correct routing would be very high as well. The cost includes not only the capital expenditures (CapEx) of acquiring hardware that can support full header examination at every hop, has the table space to hold per-flow tables across millions of flows, and can support the flow setup rate required in a large fabric, but also the operating expenses (OpEx) in terms of power use for such devices. Power drives much more of large-scale design than is often considered; however small the energy cost per packet to examine the entire header at each device, it can still add up, over billions of packets switched, to significant numbers.

Third, the use cases for such flow-based forwarding, in the real world, tend to be rather narrow. Replacing the control plane that manages millions of flows through a large-scale data center fabric to support custom routing for a few thousand flows at any given time doesn't appear to be a good tradeoff in terms of complexity and network manageability.

Of course, SDNs can operate in a mode where most traffic is forwarded based on the packet destination, and the small number of flows that need special routing are handled by examining the full packet header (the five tuples noted previously or deeper), but this solution is a compromise with reality, rather than the original ideal of SDNs. The concluding section of this article considers the more realistic option of compromising with reality, so it is not covered here.

Making the Network Programmable

Finally, SDNs have promised a great deal in terms of network programmability. The breakdown involves three different areas: dynamic provisioning, and dynamic interactions between applications and the network. These topics are considered in the sections that follow.

Dynamic Provisioning

If there's one point virtually every network engineer agrees on, it's that large-scale networks are difficult to provision, monitor, and troubleshoot. It would certainly be a boon to network operations, particularly in large networks, if there were a single, unified interface into every vendor's platform, and every control-plane implementation deployed across the network, to facilitate provisioning and management. While the idea of a single interface is noble, the reality of the market is probably going to intercede—as it has many times in the past—because vendors must be able to differentiate themselves somehow in order to actually sell hardware, software, and services. This reality isn't an indictment of vendor business models, it's just reality as it exists. There are two ways to express this problem.

First, vendors try to differentiate themselves with new features, architectures, and ideas their competitors don't have. New ideas, however, require new models that can be used to configure and manage newly designed and/or modified hardware and software.

If the vendor publishes standardized models for managing these things before they are completed, they lose competitive advantage.

Second, vendors tend to be able to command higher returns on vertically integrated solutions that are easy to deploy and manage *as a unit*. Building vertically integrated solutions, however, tends to thrive on well-integrated, single-vendor interfaces between the parts.

Both of these factors place vendors in the position of trying to balance openness with profit margins. The market demands openness, but it also demands simplicity and innovation, and these goals are sometimes (or even often) contradictory from the vendor's perspective.

The most likely result of these two factors is that SDN interfaces tend to be restricted in their scope and scale to the “lowest common denominator” of available features. Some level of configuration and trace information might be available through vendor-specific extensions, but not on the “common model.” Models such as *OpenFlow* tend to start with clean implementations, and then tend to fragment over time as vendors rush to build product. There is little incentive to consider additions to the base work, along with the rework such additions would require on a per-vendor basis, over time.

There is tension around automated provisioning from the network operator's side, as well. On the positive side, dynamic provisioning does take humans out of the repetitive action loop of quickly provisioning network devices and virtual topologies. Thus the speed and accuracy of configuration, provisioning, and fault isolation can be improved dramatically; in other words, automation can reduce the *Mean Time Between Mistakes* (MTBM). However, automating processes also introduce a level of brittleness into the operational cycle that can be undesirable.

Brittleness, in this context, can be seen as a set of systems that react to a wide array of situations with a small set of behaviors. Just as there can be monocultures in bacteria colonies, there can be monocultures in networks. To give a specific example, if every implementation of IS-IS in the network reacts the same way to a given situation, then it's possible for a single defect to cause every router in the network to fail under a single (though perhaps unusual) set of conditions.

The same sorts of situations can arise in provisioning or managing a network; an event that “slips through the cracks” of the automation system, or an attacker who can feel out the perimeters of defense, can take an entire system down very quickly. Another term for this situations is “robust yet fragile”:

At some point, any complex system becomes brittle—robust yet fragile is one phrase you can use to describe this condition. A system is robust yet fragile when it is able to react resiliently to an expected set of circumstances, but an unexpected set of circumstances will cause it to fail.^[4]

The best ways to counter are to intentionally avoid monocultures where possible, and intentionally inject human decision points in the process. Reducing repetitive human work is good, but removing humans from the entire decision process is bad. This brittleness can end up replacing a large number of smaller failures due to human error and replace them with large systemic failures.

Application Interaction

Combining dynamic provisioning and dynamic policy results in what can be called an *Application Programming Interface* (API) for the network itself. Treating the network as a programmable entity allows applications to directly interact with the network as a system. The general idea can look something like this:

- An application needs a certain amount of bandwidth with specific quality-of-service parameters at a particular time.
- The application uses an interface into a controller to reserve this bandwidth, providing the controller with the impacted endpoints, etc.
- The controller uses some means to build the right network conditions to accommodate the needs of the application.

Another example might be offloading the processing of packets for security reasons into the network. Applications and operating system security are becoming more widely deployed as encryption of data in motion becomes more common. For instance, LinkedIn currently deploys *Transport Layer Security* (TLS) on all external-facing connections, and is in the midst of deploying TLS across the data center fabric among internal applications. This type of encryption reduces the usefulness of firewalls as network appliances (or a “bump in the wire”) for blocking various types of attacks. The movement towards application and operating system security, however, means the host must perform all filtering, and must also forward traffic that needs to be forwarded to a honeypot or collection point for further processing. If the network has a policy interface, however, the host could instruct the controller to install policy at any point in the network that makes sense to either block or redirect attacker flows. This model would take security-related packet processing off the host and place it into the network, where specialized hardware can be deployed, and traffic can be optimally redirected or dropped more optimally.

The same objections that can be raised for dynamic provisioning can be raised for direct interaction between applications and the network, such as brittleness. To such interactions can be added the potential for feedback loops between various applications and network conditions (the main reason live measurement of network conditions was removed from the *Enhanced Interior Gateway Protocol* [EIGRP], soon after its first deployments, and replaced with relatively static metrics).

Summary of Network Programmability

Once a dynamic interface to the network as a network is in place, this abstraction can breed complexity beyond what the engineers responsible for maintaining and troubleshooting the network can readily understand. This complexity leads to several different problems, such as the “magic-button effect,” where no one really knows why “doing x” solves a particular problem, but since no one can figure it out (and no one has time to figure it out), someone writes a script that “pushes the button” every time “x” happens.

Overall, then, the promise of SDNs in the provisioning space is great—but parallel complexities must be managed. At this point, there is little sense that our understanding of SDN complexity has matured to the point of being able to use the full potential of the technology in the provisioning space.

Conclusion: Looking to the Future

While SDNs aren’t poised to “consume the world” in their original form because of issues surrounding centralized controllers, scale, and speed, the concepts involved are beginning to be applied to many different problem spaces. A hybrid-mode approach that allows a more standard distributed control plane to provide forwarding information for the bulk of the traffic based on the destination address, but allows overriding forwarding decisions based on other factors for a small percentage of the traffic, is gaining traction in data center fabrics of all sizes. Programmability is being used in long-haul networks, particularly in conjunction with optical transport, to handle customized forwarding as well. Essentially, the model that’s being adopted in the real world is splitting policy from base reachability, leaving the base reachability under the control of distributed control planes, while moving policy-based forwarding into a controller.

Leaving the the proven scalable distributed control plane in place and using SDN to take advantage of the perks such as traffic engineering, bandwidth optimizations, intelligent routing, special policies, and other uses seems to be the most practical path forward. Network operators may find themselves deploying different mixes of SDN-type controls and distributed control planes based on application support and business strategy, but there’s little doubt that both distributed control planes and what will be called SDN—programmability layered on top of the distributed control plane—will both continue to be used into the foreseeable future.

SDN is not a product. Rather, it’s a methodology or tool; not a destination, goal, or product to sell, or sometimes to market, and should not be considered a target to reach but a strategy to perform certain tasks depending on real needs and if certain requirements apply.

References

- [1] See, for instance, Liron Schiff, Stefan Schmid, and Petr Kuznetsov, “In-Band Synchronization for Distributed SDN Control Planes,” *ACM SIGCOMM Computer Communication Review*, Volume 46, Number 1, January 2016, <http://www.sigcomm.org/sites/default/files/ccr/papers/2016/January/0000000-0000004.pdf>
- [2] Such as the Internet Drafts **draft-baker-ipv6-isis-dst-src-routing**, **draft-baker-ipv6-ospf-dst-src-routing**, and **draft-ietf-spring-segment-routing**.
- [3] OpenFlow 1.3 moves towards the proactive installation of forwarding-table information in recognition of the timing issues involved in reactive control planes. This ability does resolve some components of this problem, but not others.
- [4] Russ White and Jeff Tantsura, *Navigating Network Complexity: Next-Generation Routing with SDN, Service Virtualization, and Service Chaining*, Addison-Wesley Professional, 2015, ISBN-13: 978-0133989359.

RUSS WHITE has more than 20 years of experience in designing, deploying, breaking, and troubleshooting large-scale networks. Across that time, he has co-authored more than 40 software patents, has spoken at venues throughout the world, has participated in the development of several Internet standards, has helped develop the *Cisco Certified Design Expert* (CCDE) and *Cisco Certified Architect Certification* (CCAR) certifications, and has worked in Internet governance with the *Internet Society* (ISOC). Russ is currently a member of the Architecture Team at LinkedIn, where he works on next-generation data center designs, complexity, and security. His most recent books are *The Art of Network Architecture* and *Navigating Network Complexity*. Russ holds an MSIT from Capella University; an MACM from Shepherds Theological Seminary; CCIE #2635, CCDE 2007:001, and CCAR certifications, and is currently working on a PhD at Southeastern Theological Seminary. You can find Russ at <http://ntwrk.guru/> and [linkedin.com/in/rw777](https://www.linkedin.com/in/rw777)

SHAWN ZANDI is a lead infrastructure architect with LinkedIn, where he builds large-scale data center and core networks. Shawn currently lives in San Francisco, California. For the past 15 years, he has worked as network and security architect for consulting firms from Dubai to Silicon Valley. In addition to a bachelor's degree in computer science from ATS University of Technology, Shawn holds more than 40 industry certifications including triple *Cisco Certified Internetwork Expert* (CCIE) and CCDE certifications. He can be reached via [linkedin.com/in/szandi](https://www.linkedin.com/in/szandi)

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Fabrizio Accatino	Octavio Alfageme Gorostiaga	David Martin	Scott Sandefur
Scott Aitken	Barry Greene	Timothy Martin	Arturas Satkovskis
Matteo D'Ambrosio	Geert Jan de Groot	Gabriel Marroquin	Phil Scarr
Danish Ansari	Gulf Coast Shots	Carles Mateu	Jeroen Van Ingen
John Bigrow	Martin Hannigan	Juan Jose Marin Martinez	Schenau
Axel Boeger	John Hardin	Brian McCullough	Roger Schwartz
Kevin Breit	Headcrafts SRLS	Carsten Melberg	SeenThere
Ilia Bromberg	Edward Hotard	Kevin Menezes	Scott Seifel
Christophe Brun	Bill Huber	Bart Jan Menkveld	Yaron Sheffer
Gareth Bryan	Hagen Hultzsc	William Mills	Tj Shumway
Scott Burleigh	Karsten Iwen	Charles Monson	Thorsten Sideboard
Jon Harald Bøvre	David Jaffe	Andrea Montefusco	Helge Skrivervik
Olivier Cahagne	Dennis Jennings	Fernando Montenegro	Darren Sleeth
Roberto Canonico	Jim Johnston	Tariq Mustafa	Mark Smith
Lj Cemerar	Jonatan Jonasson	Stuart Nadin	Job Snijders
Dave Chapman	Daniel Jones	Mazdak Rajabi Nasab	Peter Spekreijse
Stefanos Charchalakakis	Amar Joshi	Krishna Natarajan	Thayumanavan Sridhar
Greg Chisholm	Merike Kaeo	Darryl Newman	Matthew Stenberg
Narelle Clark	David Kekar	Ovidiu Obersterescu	Adrian Stevens
Steve Corbató	Shan Ali Khan	Mike O'Connor	Clinton Stevens
Brian Courtney	Nabeel Khatri	Carlos Astor Araujo	Viktor Sudakov
Dave Crocker	Henry Kluge	Palmeira	Edward-W. Suor
John Curran	Alexader Koga	Alexis Panagopoulos	Roman Tarasov
Morgan Davis	John Kristoff	Manuel Uruena Pascual	Phil Tweedie
Freek Dijkstra	Terje Krogdahl	Ricardo Patara	Unitek Engineering AG
Geert Van Dijk	Bobby Krupczak	Alex Parkinson	John Urbanek
Karlheinz Dölger	Warren Kumari	Dipesh Patel	Martin Urwaleck
Andrew Dul	Darrell Lack	Dan Paynter	Betsy Vanderpool
Peter Robert Egli	Yan Landriault	Chris Perkins	Surendran Vangadasalam
George Ehlers	Markus Langenmair	Rob Pirnie	Alejandro Vennera
Torbjörn Eklöv	Fred Langham	Blahoslav Popela	Luca Ventura
Peter Eisses	Richard Lamb	Tim Pozar	Tom Vest
ERNW GmbH	Tracy LaQuey Parker	David Raistrick	Dario Vitali
ESdatCo	Robert Lewis	Priyan R Rajeevan	Andrew Webster
Mikhail Evstiounin	Sergio Loreti	Paul Rathbone	Tim Weil
Paul Ferguson	Guillermo a Loyola	Justin Richards	Jd Wegner
Christopher Forsyth	Hannes Lubich	Mark Risinger	Rick Wesson
Tomislav Futivic	Dan Lynch	Ron Rockrohr	Peter Whimp
Edward Gallagher	Alexis Madriz	Carlos Rodrigues	Jurrien Wijlhuizen
Chris Gamboni	Michael Malik	Boudhayan Roychowdhury	Pindar Wong
Xosé Bravo Garcia	Yogesh Mangar	RustedMusic	Bernd Zeimetz
Serge Van Ginderachter	Bill Manning	Babak Saberi	
Greg Goddard	Harold March	George Sadowsky	

Supporters and Sponsors

Supporters



Diamond Sponsors



Ruby Sponsor



Sapphire Sponsors

Your logo here!

Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

ADDRESS SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

David Conrad, Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, Chairman
Internet Corporation for Assigned Names and Numbers

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

March 2017

Volume 20, Number 1

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
BGP Large Communities	2
Internet of Insecure Things ...	12
DNS Privacy	20
Fragments	31
Thank You.....	32
Call for Papers.....	34
Supporters and Sponsors	35

The *Regional Internet Registries* (RIRs) form an important part of the administrative ecosystem of the Internet. APNIC and RIPE are already sponsors of this journal, and we are pleased to announce that another RIR, the *Latin America and Caribbean Network Information Centre* (LACNIC), is now a sponsor of IPJ. LACNIC has additionally agreed to translate selective articles and provide article summaries from IPJ in Spanish.

Selective articles from IPJ are also available in Russian through the publication Интернет изнутри (*Internet Inside*), available at:
<http://www.ccni.ru/publications/>

Our individual donors and organizational sponsors make publication of this journal possible. We are especially thankful for the support of Rabbi Rob Thomas and Lauren Thomas who agreed to match individual donations from October 2016 until April 2017.

The *Border Gateway Protocol* (BGP) is a core component of the Internet routing system. Like most Internet protocols, BGP was developed in the *Internet Engineering Task Force* (IETF). The IETF has been criticized for its slow process, and in many cases for developing protocols without proper input from those who actually run networks, collectively referred to as the “operators.” Job Snijders describes how a group of dedicated operators were able to develop specifications for *BGP Large Communities* within the IETF process in record time.

Various aspects of *The Internet of Things* (IoT) have previously been covered in this journal. This time Bob Hinden looks at the problem of securing IoT devices in light of some large-scale attacks that exploited security weaknesses in common devices such as IP cameras.

In our final article, Geoff Huston and Joao Luis Silva Dama discuss *privacy* in the context of the *Domain Name System* (DNS). The *DNS PRIVate Exchange* (DPRIVE) Working Group of the IETF has been working on this topic, considering ways in which the interaction between a DNS client and a DNS resolver can be protected.

Visit our website for subscriptions, back issues, author guidelines, sponsor information, and much more, and send us your feedback via e-mail to ipj@protocoljournal.org.

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

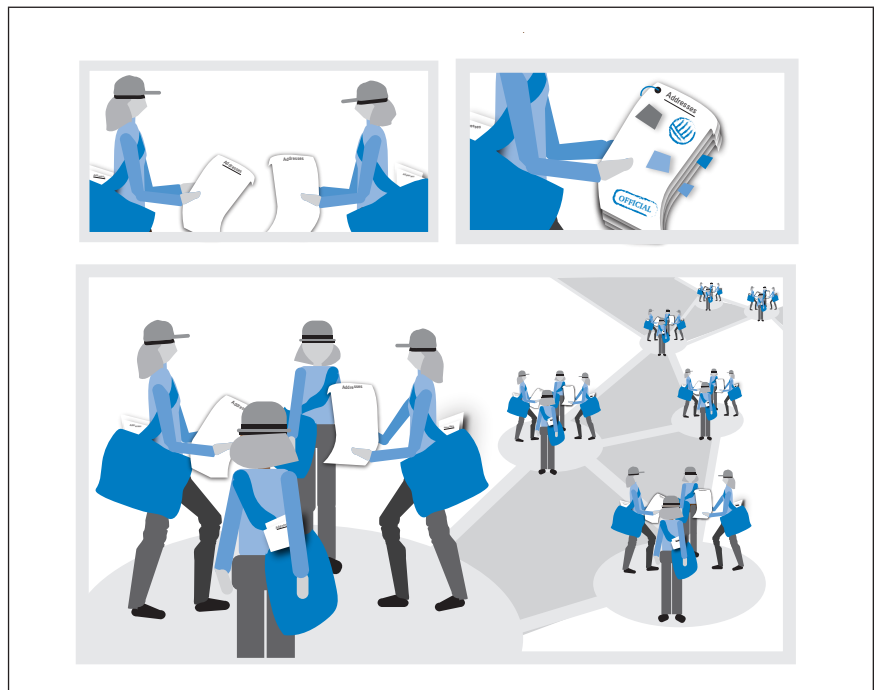
BGP Large Communities

by Job Snijders, NTT Communications

If the *Domain Name System* (DNS) is the phonebook of the Internet, used to translate names to numeric addresses, then the *Border Gateway Protocol* (BGP)^[1] is the map of the Internet: how to reach those addresses. This article focuses on a minute, but oh-so-critical aspect of BGP operations called *BGP Communities*. We will cover what BGP Communities are, why “done” is better than “perfect,” and how disaster was avoided.

A BGP crash course: The Internet is an assemblage of independent networks. The technical term for such a network is an *Autonomous System* (AS), and each is assigned a globally unique identifier called the *Autonomous System Number* (ASN). All these networks exchange routing information with other networks using BGP, a path vector protocol. This exchange of routing information is composed of announcements such as “this specific set of addresses can be reached through my network.” A set of addresses is called a *route*. Routes are often decorated with meta-information known as BGP Communities. These BGP Communities can be considered marker colors of sorts. Such marker colors are used as an additional input in the route evaluation process, which is a function of a network routing policy (Figure 1).

Figure 1: The Stamps on the List of Addresses Symbolize the Function of BGP Communities



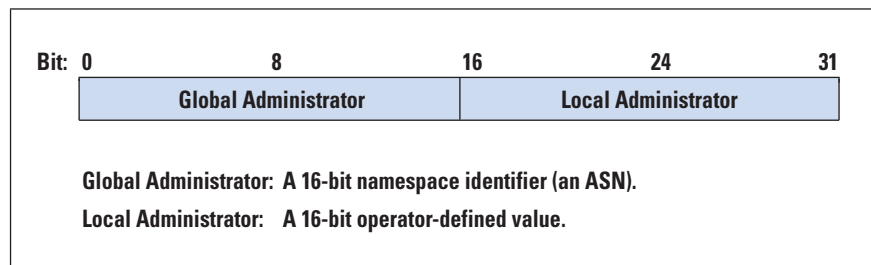
Over the last 20 years, BGP Communities have become the tool of choice to facilitate operations in all BGP networks. BGP Communities are used, for instance, to group routes together and ensure that only a specific group of routes is announced at a given interconnection point with another network.

However, along with the advent of a new type of network identifier in 2007, 32-bit ASNs (previously ASNs were 16 bits) came an operational problem: by embracing 32-bit ASNs, we had suddenly outgrown classic BGP Communities.

How Can BGP Communities Be Too Small?

Let's first explore what BGP Communities actually mean for operations. BGP Communities are defined in RFC 1997^[2]. Each BGP Community is a fixed-width 32-bit entity; you can attach multiple BGP Communities to a route. The convention is that the first 16 bits are the ASN in which the last 16 bits have a meaning. For human consumption, a BGP Community is usually represented as two 16-bit values separated by a colon (Figure 2).

Figure 2: A Classic BGP Community



An example BGP Community and its application would be **2914:664**; the first part, **2914**, is NTT Communications' ASN. The second part (called the *Local Administrator*), which carries the value **664**, is defined by NTT and has meaning only within NTT's network. According to NTT's documentation^[3], the value **664** within the **2914** namespace means "only blackhole outside the country the announcement originated." BGP Communities are the lingua franca of inter-domain routing; however, (oddly enough) its vocabulary is exchanged through out-of-band means such as published documentation. Thousands of networks have defined their own routing functions and associated BGP Communities.

In the early 2000s work began to extend the BGP protocol so that it could accommodate the ever-growing Internet. The exhaustion of the 16-bit ASN pool we are facing right now was already anticipated then. Through RFC 4893^[4] the range of possible ASNs was extended from 65,535 to 4,294,967,295 ($2^{32} - 1$). The observant reader will have noted the friction between the previously described application of a BGP Community and the existence of 32-bit ASNs: you simply cannot fit a 32-bit value in a 16-bit field!

Thus, operators of networks with a 32-bit ASN have been forced to work around this problem. Operators have used kludges ranging from using private 16-bit ASNs in the "ASN" field (those first 16 bits), to the ultimate rejection of the concept of 32-bit ASNs: returning the assigned 32-bit ASN to its respective *Regional Internet Registry* (RIR) and requesting a fresh 16-bit ASN. However, the *Internet Assigned Numbers Authority* (IANA) has been depleted of its supply of 16-bit ASNs. The RIRs are making impossible searches for 16-bit ASNs.

Rumor has it that some RIRs were considering reclamation strategies to increase their pool of 16-bit ASNs. A dire situation!

The Road to Perfection Is Always Under Construction

Surely the *Internet Engineering Task Force* (IETF) thought of this situation when updating the BGP-4 specification for 32-bit ASNs? Absolutely! But not in a way that would match operational practices. RFC 4893^[4] deferred the issue of BGP Communities as follows: “... the high-order two-octets of the community attribute [...] encode the Autonomous System number.” Quite clearly this would not work for BGP speakers that use 4-octets Autonomous System Numbers. Such BGP speakers should use the four-octet AS-Specific *Extended Communities* instead (see Figure 3).

Yes, Extended Communities, as defined in RFC 4360^[5], are bigger than regular BGP Communities: they contain a *Type* and *Subtype* field followed by 48 bits of data. However, the 48-bit length of the Extended Community value precludes the common operational practice of having the ability to encode ASNs in both the Global Administrator and the Local Administrator subfields. You can either encode a 16-bit value in the first part and a 32-bit value in the second part (Figure 3), or the other way around, but not a 32-bit value in both fields (Figure 4).

Figure 3: A BGP Extended Community Flavor #1

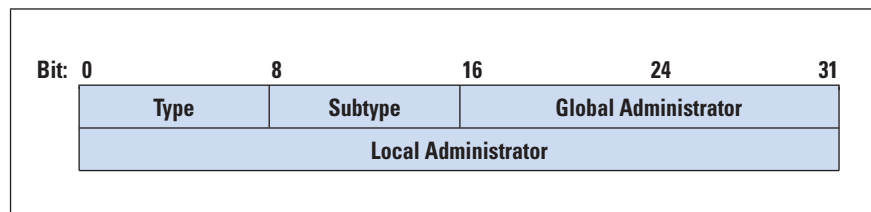
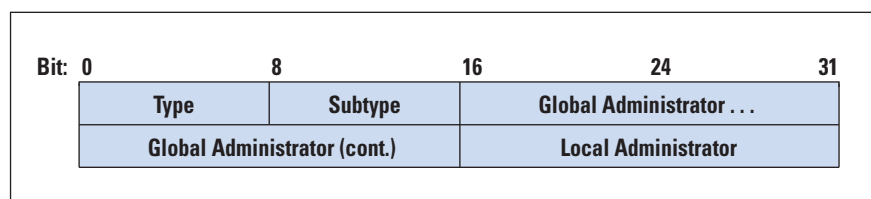


Figure 4: A BGP Extended Community Flavor #2



The *Type* and *Subtype* fields define the kind of Extended Community. The Value field can be split into either a two-octet *Global Administrator* subfield and a four-octet *Local Administrator* subfield, or a four-octet *Global Administrator* and a two-octet *Local Administrator*.

Even if your network does not have a 32-bit ASN, you might have to interact with 32-bit ASN networks. If we go back to the example of NTT’s routing policy^[3], one of the traffic engineering features is exposed through BGP Communities **65501:nnn**, where nnn is to be replaced with the Peer ASN and **65501** means “prepend AS 2914 once on outbound.” Ideally this works for 32-bit ASNs too!

The previous example highlights two issues: apparently AS 2914 ran out of BGP Community space and resorted to using a Private ASN in the first 16 bits: 65501; this situation is considered a form of namespace pollution. And secondly, had 2914 instead been a 32-bit ASN (such as AS 199036), you would not be able to encode a 32-bit ASN in the second field, even with Extended Communities. In other words, you cannot fit 64 bits worth of information into 48 bits of room. You can find many examples^[6] of networks offering traffic engineering features through the verbatim reference of a target Peer ASN in the BGP Community itself. Extended Communities have excellent use cases, but simply put, they aren't big enough for Internet routing operations.

So, if both Communities and Extended Communities weren't suitable, a third reimagination of the Communities *concept* was needed. Since these technologies tend to last for decades, and the previous two iterations had proven not to be usable if the technology changes, the bar was set pretty high.

Notable efforts in this problem space include *Flexible Communities* (started in 2003) and *Wide Communities* (started in 2010). These efforts have highlighted a disconnect between the IETF and operator communities. Not only were these two efforts moving forward on vastly different timescales than the operational community required (keep the impending doom of 16-bit ASN exhaustion in mind), but both efforts presented a tendency towards *feature creep*. The limitless extensibility was both a virtue and a curse: every possible use case would be consumed effortlessly in the specification, so the schedule overran, the specification complexity increased, and as a result no actual implementations had been produced. "A bird in the hand is worth two in the bush."

Another anti-pattern was at work in the IETF *Inter-Domain Routing* (IDR) Working Group. The anti-pattern is commonly called *Cookie Licking*. Cookie Licking is a reference to the metaphorical situation where someone takes a cookie, licks it, puts it back on the tray, and does not eat it! In volunteer communities this phenomenon can be noticed when a certain topic is discussed and someone says "I've already got this covered in my Internet-Draft." This statement might defer others from working on the topic, since you could easily assume the problem will be taken care of. But with the 32-bit ASN specification approaching the respectable age of 10 years, the pool of 16-bit ASNs running out, and no commonly acceptable successor available for BGP Communities, it became increasingly clear that a catastrophe was imminent. The issue had become a matter of extreme urgency: if we didn't start turning the tanker right now, it would crash into a wall two years down the road.

Origin Story: BGP Large Communities

In March 2016 in Sweden, Ignas Bagdonas (Equinix) presented his perspective on “some form of larger BGP communities.” In this presentation he iterated over challenges that 32-bit ASNs network operators face and provided an up-to-date overview of current and past attempts to mitigate those issues. After this presentation, I approached Bagdonas and proposed to team with him and jointly drive this effort forward.

In April 2016, I flew to the United States, where my friend Jared Mauch introduced me to Jakob Heitz (Cisco) over lunch. Heitz was intrigued by the effort and committed to providing running code for some form of “a simple, bigger BGP Community.” This agreement meant the first router vendor got on board, even if we didn’t know what the results would look like!

In May 2016 at the RIPE 72 meeting, Bagdonas presented again in the *Routing Working Group*. In the Q&A session, Ruediger Volk (Deutsche Telekom) made one of the most formative comments on the “larger communities” effort. Volk said: “The discussions in IETF about extending this [communities] have been around for many years, and no progress has been made. Any proposal that has an extended functionality comes with the problem that discussion of the additional functionality does not have a proof of termination.”

In other words, the IETF suffers from *bikeshedding*^[7]. Only something that was purposefully specified to be as narrow and as simple as possible would meet the network operator community’s immediate needs. Volk went on to say that something similar to the opaque approach of RFC 1997^[2] should be done, where the first bits indicate “Who owns the namespace” followed by some extra bits for the actual routing policy work.

It was this specific argument about namespace that led to defining Large BGP Communities as a 96-bit entity: All operators whether they have a 16-bit or a 32-bit ASN, would have 64 bits (8 bytes) of room to signal information or trigger actions in their network. With 64 bits available to the network operator, there even is enough to target a 32-bit ASN and still have space for an action such as “prepend” or “do not export.”

These conversations in the operational community led directly to Large Communities (by design) not being extendable. Extensibility comes at a cost. Also, knowing that the amount of noise generated by an idea is inversely proportional to the complexity of the idea, the IETF was asked to consider the simplicity of the Large Community a virtue, not a disadvantage. Of course, that request was ignored: over the span of just four months, almost 5% of all emails ever sent through the IDR mailing list in the last 20 years addressed the topic of Large Communities.

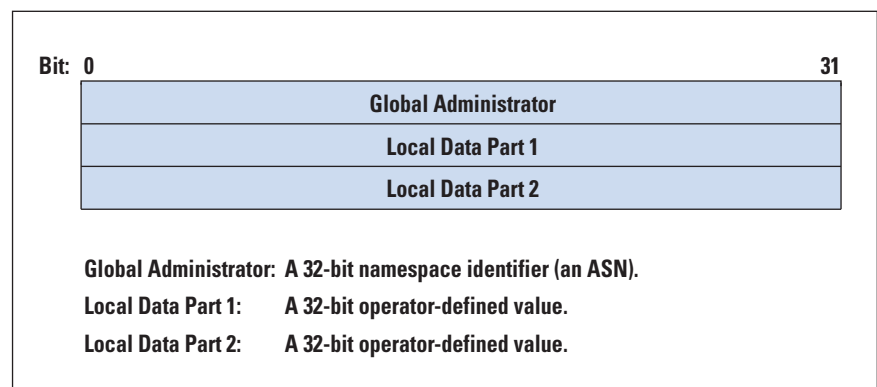
Aside from the fierce debate on the actual technology itself, the IANA *Early Allocation Procedure* uncovered very serious deployment issues. Initially Large Communities were assigned *BGP Path Attribute* value 30 by the IANA. Immediately following this allocation, BGP Beacons with the Large Community attached were brought online^[8] to test for problems of any kind. We expected the beacon prefix at least to propagate to all corners of the Internet, and optionally with its *Transitive Large BGP Communities* attribute still attached and intact. The beacon was launched from my home, but unfortunately the assigned codepoint did not pass the “Family Acceptance Factor.” My significant other noticed how certain websites that were reachable the day before could be reached no more. Collective sleuthing of the NLNOG^[9] operators brought to light that certain Huawei software was using BGP Path Attribute value 30 for something entirely different than Large Communities, and as such legitimate Large Communities were considered invalid and treated such prefixes with a withdraw. After this knowledge became public, Cisco and Juniper also emerged and a grand total of six squatted BGP Path Attribute values were uncovered^[10]. Squatting in this context means that the code point is used for a different purpose than the one designated by the IANA. In the end BGP Large Communities was assigned a new value: 32 — a very befitting number for the problem it solves.

To end this origin story: Ruediger Volk provided the effort with a challenge: “If we go really fast-track, we’ll be there in 5 years’ time, not before.” We all know that nothing motivates as much as a solid race against the clock: In exactly 6 months the Large Communities team produced 18 versions of the technical specification^[11], patched 2 packet analyzers, developed 7 implementations^[12], and obtained approval from the *Internet Engineering Steering Group* (IESG) to publish the document as an IETF Standards Track RFC.

So What Is This “BGP Large Community”?

An example of a Large Community is **2914:65400:38016**. Each BGP Large Community value is encoded as a 96-bit quantity: three unsigned 32-bit integers, separated by a colon. The **2914** part is called the *Global Administrator*, and the second and third fields (**65400** and **38016**, respectively) are called *Local Data Part 1* and *Local Data Part 2* (Figure 5).

Figure 5: A BGP Large Community



The Global Administrator is a 32-bit namespace identifier that allows different Autonomous Systems to define Large Communities without collision. The recommendation is to put your RIR-assigned ASN in the Global Administrator field. Using a 32-bit field allows full parity and fairness between 16-bit and 32-bit ASNs. Everyone can use their own ASN in Large Communities. The “Local Data Parts” are to be interpreted as defined by the owner of the ASN.

Special care was taken to define a canonical representation for Large Communities. Especially in our international community where we face communication challenges because of language barriers, it is important that Large Communities are easy to remember and easy to communicate by phone or email.

Use of BGP Large Communities

A design pattern promoted by the Internet Draft **draft-ietf-grow-large-communities-usage**^[13] specifies a **ASN:Function:Parameter** pattern to fill the three Large Community fields.

In existing deployments of Communities RFC 1997^[2] and preliminary deployments of Large Communities, two categories of Communities exist: *Informational Communities* and *Action Communities*.

Informational Communities serve as markers regarding, for instance, the origin of the route announcement, the relation with the *External Border Gateway Protocol* (EBGP) neighbor, or the intended propagation audience. Informational Communities also assist in network operations such as debugging. The Global Administrator field is set to the ASN that marks the routes with the Informational Communities. As an example: on a route that AS 64497 announces to AS 64498, AS 64497 might add Large BGP Community **64497:100:31** to signal to AS 64498 that the route was learned in the Netherlands. In this instance, the **100** value in *Local Data Part 1* is an indicator for the function “in which country a route originated” and the value **31**, as parameter, symbolises the Netherlands. In general, the intended audience of Informational Communities is downstream networks, but any Autonomous System could benefit from receiving these communities.

Action Communities are attached to routes to request nondefault behaviour in an Autonomous System. For instance, Action Communities are used to change the propagation characteristics of the route, or BGP Path attributes such as *LOCAL_PREFERENCE*, *AS_PATH*, and so forth. The Global Administrator field is set to the ASN value of the AS that has defined the meaning of the remaining fields and is expected to perform the action upon receiving the route. For instance, if AS 64499 wants to request AS 64497 to lower the *LOCAL_PREFERENCE* to **50** (below the default of **100**), AS 64499 could tag the route with **64497:20:50**. In general the intended audience of Action Communities is an upstream provider.

A real-life example of the application of Large BGP Communities can be found at the Route Servers operated by the *Internet Neutral Exchange Association* (INEX), Ireland. INEX^[14] extended its Route Server routing policy to support control over which Peer ASN receives what routing information through a trivial suppress/unsuppress mechanism. For instance, Large Community **43760:1:peer-as** means “Announce prefix to a peer-as,” whereas **43760:0:peer-as** means “Prevent announcement of a prefix to the ASN filled in as peer-as.” With BGP Communities such signaling was exclusively applicable to 16-bit ASNs.

Implementations

An IETF attitude from the past is “rough consensus and running code,” and especially for a technology like Large Communities (designed to be used in an inter-domain context), the more implementations, the more stable and commonly accepted the standard is. At the moment of writing there are nine confirmed implementations in various stages of general availability.

It is clear that the open source projects have taken a lead in implementing Large Communities. I attribute this situation in part to the fact that any contributor can dedicate time to create a patch. Secondly, the size of an ecosystem and the expected self-reliance (if any) of the user base affect the size of the effort. After all, most open source projects won’t need, or won’t have to train pre- and post-sales staff and technical assistance centers, update volumes of internal and external documentation, and, last but not least, figure out how to port the feature into the many concurrently supported releases.

These open source projects benefited from an open regression testing suite: *The Large Communities Playground*^[15]. This open source cross-vendor effort brought easy specification compliance, and functional testing to any Open Source Software project runnable inside Docker. The availability of this tool prevented duplicate efforts in establishing the appropriate environments, and ensured consistent interoperability of the implementations.

Considering the information available to me right now: *ExaBGP*, *GoBGP*, *OpenBGPD*, *rtbrick*, and *pmacct* are shipping Large Communities in their stable releases. *Quagga* and *frr* ship today or will do so shortly. *Arista EOS*, *Cisco IOS XR*, Juniper’s *Junos OS*, and Nokia *SR OS* are expected to publish software in the second half of 2017. Ancillary parts of the ecosystem such as *tcpdump*, *Wireshark*, *pbgpp*, *zebra-dump-parser*, *bgpdump*, and *mrtparse* have also been updated. Support for Large Communities in just the BGP Speakers simply wouldn’t be enough: to be a viable alternative to BGP Communities, every element in the ecosystem has to be updated, from packet analysers to research tools to statistics backends.

The future looks bright...2018 will be the year of Large Communities.

Conclusion

This challenging crusade through the IETF process has changed me from being an outside operator to a participant in the standardisation community. As much as Large Communities are a part of me, I am now a part of the IETF. Any criticism against the IETF institution displayed in this article applies equally to me as it does to others. We need to be vigilant to prevent another occurrence where operational reality and standardisation efforts grow apart so far that the only recourse left is to round up all the operators and show up with pitchforks and torches. In the end, it's just the results that matter; anything else is unimportant. Large Communities are here now, and will be so for the next decades.

Large Communities closed the final feature-parity gap between 16-bit and 32-bit ASNs. Now 32-bit ASNs are first-class citizens too.

Acknowledgements

I would like to thank Arjen Zonneveld, Saku Ytti, Russ White, Ruediger Volk, Teun Vink, Gunter van de Velde, Jeff Tantsura, Sander Steffann, Richard Steenbergen, Wesley Steehouwer, Adam Simpson, Rob Shakir, Shyam Sethuram, Julian Seifert, Mark Schouten, Tom Scholl, Martijn Schmidt, Alvaro Retana, Kay Rechthien, Gaurab Raj Upadhaya, Joe Provo, Stefan Plug, Tom Petch, Jussi Peltola, Keyur Patel, Barry O'Donovan, Arnold Nipper, Christopher Morrow, Remco van Mook, Martin Millnert, Jared Mauch, Marco Marzetti, Ben Maddison, Joel M. Halpern, Duncan Lockwood, Acee Lindem, Kristian Larsson, Warren Kumari, Thomas King, Grzegorz Janoszka, Geoff Huston, Paul Hoogsteder, Nick Hilliard, Wim Henderickx, John Heasley, Markus Hauschild, Richard Hartmann, Greg Hankins, Jeffrey Haas, David Freedman, Bill Fenner, David Farmer, Matt Griswold, Bertrand Duvivier, Linda Dunbar, Brad Dreisbach, Gert Doering, Peter van Dijk, Brian Dickson, Ian Dickinson, Marco Davids, Adam Davenport, Tom Daly, Nabeel Cocker, Mach Chen, Adam Chappell, Pier Carlo Chiodi, Randy Bush, Jay Borkenhagen, Theodore Baschak, Niels Bakker, Jan Baggen, Ignas Bagdonas, Alexander Azimov, Eduardo Ascenco Reis, Mikael Abrahamsson, John Scudder and Susan Hares for their support in developing BGP Large Communities. Special thanks to Corinne Pritchard for the imagery.

References

- [1] Rekhter, Yakov, Hares, Susan, and Li, Tony, "A Border Gateway Protocol 4 (BGP-4)," RFC 4271, January 2006.
- [2] Chandra, R., Traina, P., and Li, T., "BGP Communities Attribute," RFC 1997, August 1996.
- [3] "NTT Routing Policies,"
<https://www.us.ntt.net/support/policy/routing.cfm>
- [4] Chen, Enke and Vohra, Quaizar, "BGP Support for Four-octet AS Number Space," RFC 4893, May 2007.

- [5] Sangli, Srihari R. and Rekhter Yakov, “BGP Extended Communities Attribute,” RFC 4360, February 2006.
- [6] “BGP Communities Guide,”
<https://onestep.net/communities/>
- [7] “Why Should I Care What Color the Bikeshed Is?”
<http://bikeshed.com/>
- [8] Snijders, J., “Large BGP Communities Beacon,” October 2016.
<http://largebgpcommunities.net/2016/beacon/>
- [9] Netherlands Network Operator Group, <https://nlnog.net/>
- [10] Snijders, J., “Deprecation of BGP Path Attribute Values 30, 31, 129, 241, 242, and 243,” RFC 8093, February 2017.
- [11] Heitz J., Snijders, J., Patel, K., Bagdonas I., and Hilliard H., “BGP Large Communities,” RFC 8092, February 2017.
- [12] “BGP Large Communities Implementations,”
<http://largebgpcommunities.net/implementations/>
- [13] Snijders, J. and Schmidt, M., “Usage of BGP Large Communities,” Internet Draft, Work In Progress, December 2016.
<https://tools.ietf.org/html/draft-ietf-grow-large-communities-usage>
- [14] “INEX Deploys Large BGP Communities In Production,” November 2016,
<http://largebgpcommunities.net/2016/inex-first-in-production/>
- [15] Chiodi, P. C. and Snijders, J., “The BGP Large Communities Playground,” January 2017,
https://labs.ripe.net/Members/pier_carlo_chiodi/the-large-bgp-communities-playground

JOB SNIJDERS is IP Development Engineer at NTT Communications, where he analyzes and architects NTT’s global IP network for future growth. He has been actively involved in the Internet community in an operational capacity, as a frequent presenter at network operator events, and in numerous community projects for over 10 years. Job is founder of the NLNOG Foundation and vice president of PeeringDB. Job’s special interests are routing policy, routing security, and large-scale BGP deployments. He maintains several tools such as *irrtree* and *irrexplorer*, and is active in the IETF, where he has coauthored or contributed to RFCs and Internet-Drafts. E-mail: job@ntt.net, Twitter: [@JobSnijders](https://twitter.com/JobSnijders)

The Internet of Insecure Things

by Bob Hinden, Check Point Software

We have a problem. As we have learned recently, most *Internet of Things* (IoT) devices are not secure. They have numerous security weaknesses, including default login/passwords and fixed firmware login/passwords (not easily changeable). In addition, the devices do not get regular software updates to fix security problems, and there is no one to call if problems arise. To make matters worse, these devices exist in large numbers. Gartner says that 6.4 billion IoT devices are deployed now, and forecasts 20.8 billion in 2020. The IoT contains billions of insecure things.

The security problem with these devices is not theoretical; the devices are used today to launch very large-scale *Distributed Denial of Service* (DDoS) attacks. Malware that takes advantage of the security weaknesses of these devices to create botnets exists.

IoT-Based DDoS Attacks

The first large IoT-based attack I am aware of was the DDoS attack on the popular online security publication *KrebsOnSecurity* blog^[1]. This attack happened on September 13, 2016. The blog had recently published a series of articles about an organization called *vDOS*, a DDoS-for-hire service. The blog stated that *vDOS* made approximately \$600k in 2 years by knocking sites offline. Two weeks after the articles were published, *KrebsOnSecurity* was attacked with 620 Gbps of traffic^[2]. Akamai was providing pro-bono DDoS protection to *KrebsOnSecurity*, but it couldn't continue to handle the traffic load^[3]. Fortunately, *Google Project Shield*^[4] is now protecting the website. Apparently, the source of the attack was IoT devices that included a myriad of technologies composed of IP cameras, *Digital Video Recorders* (DVRs), home routers, and other embedded computers.

The next large IoT-based DDoS attack occurred in early October of 2016. OVH^[5] is a large international web-hosting provider. The company was attacked by a botnet comprising more than 145,000 compromised IP cameras and digital video recorders. The attack peaked at 1 Tbps, and fortunately, OVH was able to withstand the attack. This attack was the largest DDoS attack at the time.

A very visible attack occurred on Friday, October 21, 2016. This attack was directed against DYN, a large provider of *Domain Name System* (DNS) services^[6]. The attack limited access to many of DYN's customers, including some of the biggest sites on the Internet, like Twitter, Amazon, Tumblr, Reddit, Spotify, and Netflix. For many Internet users, the Internet was down. Table 1 lists the major sites that were affected.

Table 1: Major Sites Affected in the Dyn Attack

Airbnb	Amazon.com	Ancestry.com	The A.V. Club	BBC
Boston Globe	Box	Business Insider	CNN	Comcast
CrunchBase	DirecTV	Elder Scrolls	Electronic Arts	Etsy
EQAO	FiveThirtyEight	Fox News	Guardian	GitHub
Grubhub	HBO	Heroku	HostGator	iHeartRadio
Imgur	Indiegogo	Mashable	NHL	Netflix
NYT	Overstock.com	PayPal	Pinterest	Pixlr
PlayStation	Qualtrics	Quora	Reddit	Roblox
Ruby Lane	RuneScape	SaneBox	Seamless	Second Life
Shopify	Slack	SoundCloud	Squarespace	Spotify
Starbucks	Storify	Swedish Civil Contingencies Agency	Swedish Government	Tumblr
Twilio	Twitter	Verizon	Visa	Vox Media
Walgreens	WSF	Wikia	Wired	Wix.com
WWE	Xbox Live	Yammer	Yelp	Zillow

We are still learning more about this attack, but it appears to be similar to the two earlier attacks. It's been reported that it peaked at 1.2 Tbps, another record. However, unlike the first two attacks, the effects of this one were visible to many Internet users, and it made the news cycle. We don't know for certain the motivations for this attack, but there is speculation that it was due to DYN's support of Brian Krebs, the man behind KrebsOnSecurity. It certainly got everyone's attention.

IoT Devices Involved in DDoS Attacks

The devices used in these attacks include a range of IP cameras, DVRs, and home routers. What these devices share is that they all ship with default login/passwords and can be enabled without changing these defaults. Some even have fixed firmware login/passwords that can't be changed. They are all widely available and fairly low-cost. For example, a search on Amazon for "IP camera" will result in several hundreds of pages of IoT devices for sale, most less than \$100 USD. While there clearly has been a lot of innovation to build so many different types and styles of IP cameras, innovation clearly hasn't addressed making them secure.

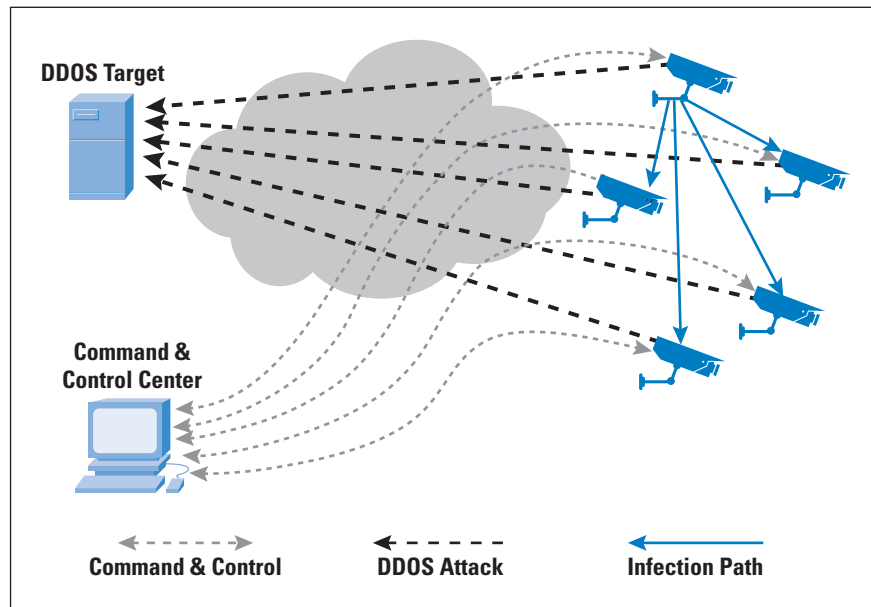
IoT Malware

The malware used in the KrebsOnSecurity attack was released recently. It is called *Mirai*^[7]. It works by scanning the Internet for vulnerable devices, looking for systems with factory default usernames and passwords. It installs itself in these IoT devices, turning devices into "bots" used to launch these DDoS attacks. *Bashlight* is another form of IoT botnet malware. It's similar to Mirai, as it infects IoT devices with default usernames and passwords.

According to research from Level3 Communications, Bashlight is responsible for infecting almost a million IoT devices and now competes with other botnets based on Mirai.

Mirai has two phases; in the first phase it infects IoT devices like IP cameras. In the second phase the attack is directed from a *Command and Control System*. These phases are illustrated in Figure 1.

Figure 1: Mirai Attack Mechanism



In the first phase, it scans broad ranges of IP addresses looking for devices that allow it to log in with a set of default login name and password combinations. These default login and password combinations are derived from the documentation of IoT devices. These attacks are very simple, not very sophisticated.

Examples of these login and password combinations follow:

```
admin/admin
root/admin
root/88888
root/root
ubnt/ubnt
```

The problem is that while these IoT devices allow the user to change the default, many users don't, and the setup of these IoT devices does not require it. Thus they are very vulnerable to being taken over by malware like Mirai.

One interesting part of the address scanning is that Mirai has a list of IP address ranges it does not scan. Some of these ranges belong to large companies like General Electric and Hewlett Packard, the U.S. Department of Defense, and the U.S. Postal Service. It's hardly the complete set of addresses for these organizations; I assume the Mirai authors had some reason to avoid these addresses. While the reason is unknown, they may have been trying to avoid detection.

A group named *Malware Must Die*^[8] has discovered new variants of IoT malware. This discovery isn't too surprising since the source code for Mirai was released and a lot of malware is a variant of earlier malware. One of these variants is called *Linux/IRCTelnet*^[9]. Its attacks are similar to those of Mirai, but it uses *Internet Relay Chat* (IRC) to communicate with compromised Linux-based IoT devices. It also has the capability to attack IoT devices running IPv6.

We Should Be Worried

We should be worried for three fundamental reasons. First, the sheer scale of these attacks is enormous. As shown in Table 2, they are growing:

Table 2: Growth of Attacks

Site	Number of Attackers	Traffic
KrebsOnSecurity	1.2 million	620 Gbps
OVH	145 thousands	1 Tbps
DYN	Millions	1.2 Tbps

At this level, every organization is vulnerable. The current number and the growth rate of new IoT devices show no evidence of slowing down. And finally, the nature of most IoT devices makes attacks on some of these new devices much more difficult to stop than attacks on conventional IT devices because they don't come with any kind of effective support.

Can We Fix These Problem?

Fixing these problems will be a challenge. Technically, some of the problems are straightforward to fix. For example:

- Do not allow default login/passwords, and require users to change them before the device is enabled.
- Do not have fixed unchangeable firmware login/passwords.
- Provide automatic updates of software.

Beyond that, matters get more difficult. How does the industry provide support for low-cost IoT devices? How long will they be supported? What happens when support ends; do the devices turn themselves off? How are attacks detected and contained? Harder still, how do we fix the currently very large deployed bases of IoT devices? Is anyone in the supply chain of the IoT devices concerned?

The economics of providing real security for IoT devices raises even more questions. Who is responsible when an IoT device is taken over by malware?

The user, the retail outlet where the device was purchased, the manufacturer, or the component vendors that the manufacturer used? How do you provide long-term support for a device that sells for less than \$100 USD? What happens when support ends?

Unfortunately, there are many more questions than answers.

What Can We Do?

While there isn't a simple single solution to these problems, there are some things we can do. Possibilities range from what IoT device owners, companies that sell the devices, companies that design and build the devices, and the IoT Industry need to do.

The first and simplest is that people who own IoT devices should not run them with default passwords. That precept applies to all Internet devices. You always should use unique and hard-to-guess login usernames and passwords. You should use different passwords for each device. This common sense advice will greatly reduce the likelihood that IoT malware will infect your devices.

Note that using a password manager is recommended for IoT devices and, of course, all places where you have accounts on the Internet. These password managers allow you to use impossible-to-guess passwords. This approach is clearly much better than using "admin/admin."

The next step after setting new login/password information is to reboot the devices. The IoT malware lives in memory and will be erased if the device is rebooted. Then check to determine if the new login/passwords are still there. This step won't help with default unchangeable firmware accounts, but is still worth doing.

Companies that sell IoT devices need to ensure that the devices they sell are reasonably secure. This can include things like requiring login/passwords to be set upon installation, not having fixed unchangeable firmware passwords, and allowing users to get updates that fix known security problems. Retail channels like Amazon, Best Buy, NewEgg, etc., and companies that specialize in these devices like FLIR Systems may have the most leverage over companies that design and build these devices because the device manufacturers don't sell directly themselves. They all sell through well-known retail channels. These channels may also be liable if they are selling insecure IoT devices. Will someone file a class action suit against one of these companies? Time will tell.

Companies that design, build, and manufacture IoT devices should make their products more secure. They need to require that login/passwords account information needs to be configured as part of the installation procedure. They should also provide automatic security updates.

I am not sure that without any other motivation these companies will do this procedure, given how far in the channel they are from the end user. I think with serious encouragement from the retail channel, for example, telling them “we will stop selling your products until their security is improved,” it will happen. There may even be a market for devices that are more secure than the others. On a positive note, it has been reported that the Chinese electronics firm Hangzhou Xiongmai Technology Co Ltd has initiated a recall of its IP cameras (webcams)^[11]. It is hard to tell how effective this effort will be, but it clearly will not be as good as pushing software updates.

The IoT industry needs to develop an approach to make the IoT secure. If it continues on the current path, IoT is going to be a disaster for this industry. I don’t think we are there yet, but the signs are not encouraging. We are not going to get the grand visions of IoT everywhere if people think the devices are insecure. The current IoT malware doesn’t attack the owners of these devices, and they may not even be aware their devices are compromised, but this is only the beginning.

The potential for the collection of data on the users of these devices is staggering. If the device is compromised, its camera and sensors can be used against the user. Who wouldn’t be susceptible to a ransom demand that threatened to share what you did in your house? While I am sure that most of us don’t have anything serious to hide, would you want the world to see everything you have done or said in your house? We already have DDoS for hire firms, so what happens when we get surveillance companies hiring IoT surveillance for hire firms? It would certainly make divorce court proceedings more interesting.

One bit of good news based on these recent attacks, a lot of work is going on to create security frameworks for IoT devices. Organizations doing this experimentation include parts of the U.S. Government like the *National Institute of Standards and Technology* (NIST)^[12] and the *Department of Homeland Security* (DHS)^[13], and a variety of organizations including the *Internet Architecture Board* (IAB)^[10], the *Broadband Internet Technical Advisory Group* (BITAG), the *GSM Alliance*, the *Industrial Internet Consortium*, and the *Open Web Application Security Project*. Bruce Schneier has compiled a good list of resources^[14]. While this is all good news, it’s not clear how these efforts will cause current and future IoT devices to change to be more secure. Some mix of consumer guidance, retail channel control, and/or government regulation may be needed.

Conclusion

The problem we have with the “Internet of Insecure Things” is only going to get worse for the immediate future. Many things need to happen to improve the situation, and it isn’t going to happen quickly. I believe there are a lot of benefits if we can start building the Internet of Secure Things, but we have a long way to go. It’s important that everyone involved take security more seriously. “Everyone” includes users, the retail channel, manufacturers of the devices, and the IoT industry. It will get better only if everyone gets involved.

References and Further Reading

- [1] KrebsOnSecurity Blog: <https://krebsonsecurity.com>
- [2] “KrebsOnSecurity Hit with Record DDoS,”
<https://krebsonsecurity.com/2016/09/krebsonsecurity-hit-with-record-ddos/>
- [3] “Akamai on the Record KrebsOnSecurity Attack,”
<https://krebsonsecurity.com/2016/11/akamai-on-the-record-krebsonsecurity-attack/>
- [4] “Protecting news from digital attacks,”
<https://projectshield.withgoogle.com/public/>
- [5] The DDoS that didn’t break the camel’s VAC,”
<https://www.ovh.com/us/news/articles/a2367.the-ddos-that-didnt-break-the-camels-vac>
- [6] Dyn Statement on 10/21/2016 DDoS Attack:
<http://dyn.com/blog/dyn-statement-on-10212016-ddos-attack/>
- [7] “Mapping Mirai: A Botnet Case Study,”
<https://www.malwaretech.com/2016/10/mapping-mirai-a-botnet-case-study.html>
- [8] Malware Must Die: <http://blog.malwaremustdie.org>
- [9] “A DDoS botnet aims IoT w/ IPv6 ready,”
<http://blog.malwaremustdie.org/2016/10/mmd-0059-2016-linuxirctelnet-new-ddos.html>
- [10] Internet Architecture Board (IAB), “Report from the Internet of Things (IoT) Software Update (IoTSU) Workshop 2016),”
<https://datatracker.ietf.org/doc/draft-iab-iotsu-workshop/>
- [11] “China electronics firm to recall some U.S. products after hacking attack,”
<http://www.reuters.com/article/us-cyber-attacks-manufacturers-idUSKCN1200MS>
- [12] “Framework for Improving Critical Infrastructure Cybersecurity,”
<https://www.nist.gov/sites/default/files/documents/cyberframework/cybersecurity-framework-021214.pdf>
- [13] US Department of Homeland Security, “Securing the Internet of Things,” <https://www.dhs.gov/securingtheIoT>

- [14] Bruce Schneier, “Security and Privacy Guidelines for the Internet of Things,”
https://www.schneier.com/blog/archives/2017/02/security_and_pr.html
- [15] William Stallings, “The Internet of Things: Network and Security Architecture,” *The Internet Protocol Journal*, Volume 18, No. 4, December 2015.
- [16] Lake, D., Rayes, A., and Morrow, M., “The Internet of Things,” *The Internet Protocol Journal*, Volume 15, No. 3, September 2012.
- [17] Stankovic, J., “Research Directions for the Internet of Things,” *Internet of Things Journal*, Volume 1, No. 1, 2014.
- [18] Frahim, J., et al., “Securing the Internet of Things: A Proposed Framework,” Cisco White Paper, March 2015.
- [19] Gareth Corfield, “Fix crap Internet of Things security, booms Internet daddy Cerf,” *The Register*, March 21, 2017,
https://www.theregister.co.uk/2017/03/21/vint_cerf_internet_things_security/

BOB HINDEN is a Check Point Fellow at Check Point Software, and co-chairs the IPv6 working group in the IETF. Hinden was the Chair of the Internet Society Board of Trustees from 2013 to 2016, and a member of the Board of Trustees from 2010–2016. Previously at Nokia, he was a Nokia Fellow, Chief Internet Technologist at Nokia Networks, and Chief Technical Officer (CTO) at the Nokia IP Routing Group. He was one of the early employees of Ipsilon Networks, Inc. Nokia acquired Ipsilon on December 31, 1997. Bob was previously employed at Sun Microsystems, where he was responsible for the Internet Engineering group that implemented internet protocols for Sun’s operating systems. Prior to this position he worked at Bolt, Beranek, and Newman, Inc. on a variety of internetwork-related projects, including the first operational Internet router and one of the first TCP/IP implementations. Hinden was co-recipient of the 2008 *IEEE Internet Award* for pioneering work in the development of the first Internet routers. He has been active in the IETF since 1985 and is the author of 39 RFCs, including two April 1 RFCs. He served as the chair of the *IETF Administrative Oversight Committee* (IAOC) from 2009 through 2013. Before that he served on the *Internet Architecture Board* (IAB), was Area Director for Routing in the Internet Engineering Steering group from 1987 to 1994, and chaired the IPv6, Virtual Router Redundancy Protocol, Simple Internet Protocol Plus, IPAE, the IP over ATM, and the Open Routing working groups. He is also a member of the RFC Editorial Board and the RFC Series Oversight Committee. Hinden holds an B.S.E.E., and a M.S. in Computer Science from Union College, Schenectady, New York. E-mail: bob.hinden@gmail.com

DNS Privacy

by Geoff Huston and Joao Luis Silva Dama, APNIC

The *Domain Name System* (DNS) is normally a relatively open protocol that smears its data (which is your data and mine too!) far and wide. Little wonder that the DNS is used in many ways, not just as a mundane name resolution protocol, but as a data channel for surveillance and as a common means of implementing various forms of content access control.

But this situation is poised to change. The material released by Edward Snowden has sensitized us to the level of such activities that we have now become acutely aware that many of our Internet tools are just way too trusting, way too chatty, and way too easily subverted. First and foremost, in this collection of vulnerable Internet tools is the DNS.

A query made to the DNS is a precursor to almost every Internet transaction. Whether it's performing a search, downloading a web page, sending a mail message, opening a chat session, or even receiving an online advertisement, the DNS is often invoked as the first step. As users, we work in a symbolic world of readable names, such as **facebook.com** or **netflix.com**, whereas the underlying fabric of the Internet can send and receive packets only by using binary IP addresses rather than these symbolic names. So, we use the DNS to map from a name to an IP address. Thus, if you were able to look at a log of the DNS queries I've made in the last day or so, you may well be able to reconstruct my recent web browsing history, for example.

Obtaining a log of the DNS queries I make is perhaps the equivalent in terms of information content to obtaining a telephone log of called numbers from a previous generation of communications. A DNS transaction log may not provide information about the precise nature of the network transactions I've made, but it does record which sites I've been using. This information is often just good enough, as it's exactly what you would need to build a highly accurate profile of what I do on the Internet. It's not just national security bodies that have an interest in assembling such data logs. These days we see many systems that target individual users, and build a comprehensive profile of their needs and desires. The difference between an annoying advertisement and a timely helpful suggestion is just information about the user, and many companies assemble such profiles as part of their own commercial activities.

It's also true that the DNS is incredibly chatty. For example, to resolve a new name, such as **www.example.com**, a DNS resolver first asks the root name servers for the IP address of **www.example.com**. The root name servers would not be able to provide the answer, but they will respond with the authoritative name servers for the **.com** domain.

The resolver will then repeat this query relating to the IP address of the name `www.example.com`. to a `.com` name server, and once more the answer is an indirect one, indicating that while it does not know the answer, the list of name servers for the domain `example.com` should be queried. At this point the resolver can repeat the same query to a server that is authoritative for the `example.com` domain and probably receive an answer that contains the address of `www.example.com`.

But let's think about these DNS queries for a second. In this case, a root server, a `.com` server, and an `example.com` server are all now aware that I am "interested" in `www.example.com`, and they probably have stored a log of these queries. I have no idea if these logs are private or public. I have no idea how they are analyzed, and what inferences are drawn from this data.

It's possible that the data leakage is a little worse than described, because the application I am using, such as a browser, normally does not perform DNS name resolution itself. It passes the query to the platform operating system via a `gethostbyname()` call. The operating system platform also has an opportunity to log this query. The platform normally does not operate a standalone DNS resolver, and often is configured by the local network provider with DNS resolvers to use. So, my service provider may also be privy to all my DNS activity. But it need not stop there. My service provider might farm out its queries to a recursive forwarder, so that it can avoid the overheads of running a full DNS resolver.

Normally such forms of query indirection imply a loss of attribution, as such forwarded queries do not have any of my identifying details. Unless of course the resolver uses the *EDNS0 Client Subnet Option*^[1], in which case the forwarded queries still contain some critical details of my network, and, by inference, me as well.

All of these DNS queries can represent a lot of information, even in these days of data intensity. Back in April 2015 Google reported that its public DNS servers deliver some 400 billion responses per day^[2], and it appears that Google resolves some 12% of the total DNS load^[3], so there were some 3 trillion DNS queries per day at that time. It can only be larger today.

Not only is the DNS a chatty protocol that gratuitously sprays out information about user behaviours, it does so in an entirely open manner. DNS queries and their responses are unencrypted, and are sitting on port 53 in *User Datagram Protocol* (UDP) and *Transmission Control Protocol* (TCP). Because they are open and unencrypted, DNS queries can be easy to intercept, and, if *Domain Name System Security Extensions* (DNSSEC) is not used, false responses can be inserted back into the data stream, and the client may be none the wiser. In some countries, DNS substitution appears to be relatively commonplace^[4].

Other countries have turned to DNS interception and blocking in response to problems associated with overloading IP addresses with virtual web hosting^[5]. Little wonder that many users have tried to get around these local efforts to block access by using a third-party DNS resolver, and the use of Google's Public DNS resolver is a common response to such local efforts to interfere in the DNS. Indeed, so common is this response that now the local DNS blocking measures appear to include intercepting access to Google's DNS service as well!^[6]

DNS privacy has been a matter of some interest to the *Internet Engineering Task Force* (IETF), and changes are being proposed to the DNS protocol that would make it far harder to be used as a snooper's and censor's tool of choice.

What is going on to improve this situation and introduce aspects of privacy into the DNS?

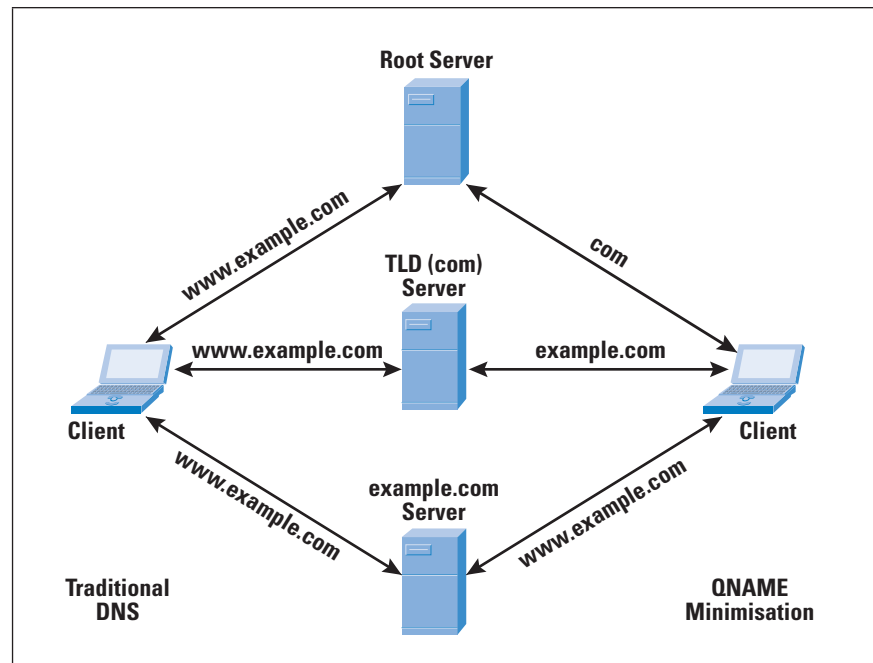
QNAME Minimisation

The IETF *DNS Operations* Working Group has been performing what has been called query name minimisation in the DNS, resulting in a specification for *QNAME Minimisation*.^[7] To quote from this document: "QNAME minimisation follows the principle [that] the less data you send out, the fewer privacy problems you have."

In the previous example, the query to the root servers for the A record for **www.example.com** has two elements of gratuitous information: the fully qualified domain name and the query type. A more targeted query that does not gratuitously leak extraneous information is a query directed to the root name servers for the name server records for the **.com** domain. Similarly, the **.com** name servers would be queried simply for the name servers of **example.com** and so on (see Figure 1).

In general, this approach is no less efficient than using a full query name at every point, and is equally capable of using cached information. The technique has exposed some inconsistencies with the handling of so-called *empty nonterminal* domain names, but the approach can be implemented in a robust manner, and it is a solid step in plugging a gratuitous information leak. It appears that the recently announced Knot DNS resolver from the Czech domain provider CZ.NIC is one of the first DNS resolvers to implement QNAME Minimisation (<https://www.knot-resolver.cz>).

Figure 1: The Intended Operation of QNAME Minimisation



DNS and TLS

However, QNAME Minimisation is only part of the privacy story. The open nature of DNS queries makes third-party monitoring, interception, and substitution incredibly easy, it appears. The *DNS PRIVate Exchange* (DPRIVE) Working Group of the IETF has been working on this topic, looking at ways for the DNS query and response interaction between a DNS client and a DNS resolver to be protected in some manner.

There are two parts to this work. Firstly, to ensure that the response you receive is a response from the DNS resolver that you intended to ask, and secondly, to ensure that the query you pass to your DNS resolver is not readily readable by anyone other than the addressed DNS resolver.

One issue here is whether to try to secure the current UDP-based resolution protocol, or head to a TCP-based approach where solutions already abound, typically based on *Transport Layer Security* (TLS). TLS encrypts the conversation between a client and a server using session keys that are generated based on random seed values coupled with public/private key pairs. If you have a way of associating a service name (such as its DNS name) with a public key (as would be the case with a conventional DNS Name public key certificate), and a way of validating that key (using conventional name certificate validation), then you can achieve both of these objectives. The remote server has demonstrated to you that it has knowledge of the private key associated with the DNS name of the service, which is theoretically known only to the server and no one else.

By encrypting your session with a session key that is based in part on this private key, the content of the data exchange should be protected from onlookers and potential interceptors in a man-in-the-middle attack.

But there is a problem here in terms of the transport protocol of choice. TLS conventionally requires a reliable transport channel, such as provided by TCP, and as such cannot be used directly to secure datagram traffic as used by UDP. However, the DNS has been heavily reliant on the use of UDP as a means of supporting speed and scalability in the DNS. So, from some perspectives, DNS-over-TLS-over-TCP is not seen as the optimal response to the problem. TCP attempts to ensure sequenced delivery, and in a message-oriented application, the loss of a message in TCP holds up the delivery of all subsequent messages until TCP can correct the data loss and deliver the lost message. This TCP “head-of-line blocking” can pose unacceptable overheads when using TCP to carry the datagram-like message payloads of the DNS. DNS over *Internet Protocol Security* (IPsec) could be seen as offering a cleaner fit when looking at securing a UDP-based application, but IPsec is a kernel function rather than an application module, and its semantics apply at the IP layer rather than as an attribute of the transport protocol. This reality makes it challenging to incorporate IPsec into an application and operate the cryptographic functions in user space, as happens with DNS name resolution.

One of the consequent investigations has been to see if the functionality of TLS could be mapped into a datagram transport environment. Out of this consideration has come a new protocol, *Datagram Transport Layer Security* (DTLS), which is an adaptation of the TLS function that can present to the application a datagram-like delivery function that does not require reliable transport services. DTLS is intended to be able to recover from packet loss and reordering, but it would be intolerant of UDP packet fragmentation^[8]. Its design is modelled upon TLS 1.2 and intends to use some explicit additional features that allow TLS to function over a datagram transport as distinct from a reliable stream transport. DTLS intends to minimise the impact of the use of TLS on the DNS experience, particularly when compared to DNS-over-TLS-over-TCP. The major change envisaged would be to require an initial DTLS handshake to set up a shared encryption state, and the use of cookies to reuse that state across multiple individual response/query interactions.

One of the main features of the current DNS protocol when used over UDP is how little shared state overhead each individual transaction incurs, resulting in a highly responsive and capable service. DNS over DTLS attempts, as far as possible, to preserve this simple query/response datagram exchange model but does so in a manner where the client is using an encryption based on the validated credentials offered by the server. The current state of play of this specification is described in an Internet Draft working document.^[9]

DTLS is intolerant of IP fragmentation, so the operation of DNS over DTLS is similar in design to the use of the *Truncated* bit in DNS over UDP as a signal to the client to repeat the query using TCP. Here the intended operation is that if a DNS-over-DTLS server has a response that is greater than the local Path *Maximum Transmission Unit* (MTU) estimate, then the server should set the Truncated bit in its response, and this response is to be interpreted by the client as a signal that the client should repeat the query using DNS-over-TLS-over-TCP, in a manner analogous to the current use of the Truncated bit to signal to a client to repeat the query using TCP rather than UDP. However, it needs to be noted that at this point DTLS remains a design exercise, and it may be some time before implementations of this specification are available for general use by end-user DNS libraries, recursive resolvers, and servers.

The other option here is to absorb the TCP overheads into the solution and just use conventional TLS, which is a TCP service. Much has been said on the use of TCP as a mainstream transport protocol for DNS, as distinct from its current intended role as a backup to UDP for large responses. It has been argued that the TCP connection state overheads of the servers seriously impair their ability to handle large query loads, and the additional overhead of the protocol handshake would negatively affect the user experience. On the other hand, it is argued that already the web is being used overwhelmingly as a short transaction service, and web servers appear to withstand the imposed load. It is also noted that the use of TCP is an effective measure against various forms of abuse that rely upon the ability to perform source address spoofing in UDP.

The specification for *DNS over Transport Layer Security*^[10] is a relatively straightforward description, in that the transport service offered by TLS is effectively the same as that offered by TCP, but running the listener of the server at TCP port 853, rather than port 443. There is perhaps one change here, and that is a suggestion for TLS session reuse: “In order to minimize latency, clients SHOULD pipeline multiple queries over a [single] TLS session.” For transactions between a client and a recursive resolver, the suggestion for session reuse makes some sense. For transactions between a client and authoritative name servers where the client is itself performing DNS resolution, this choice may not be so readily achievable. The choice of a distinguished TCP port is also interesting. If you wanted the secure channel DNS traffic to merge into all other traffic and pose a challenge to attempts to block this service, the temptation to use port 443 for DNS over TLS would be overwhelming (at least for me!). More information on the current state of clients and servers that support DNS over TLS can be found at:

<https://portal.sinodun.com/wiki/display/TDNS>.

Secure DNS over JSON

Last, but not least, there is the option to use an entirely different data encoding protocol, and here a recently announced service from Google is relevant. The server at <https://dns.google.com> performs a resolution function over TLS using port 443 with the results passed back as a *JavaScript Object Notation* (JSON) data structure. This function can readily be transformed into an alternative form of *gethostbyname()* by the application's substituting a web object retrieval for a conventional DNS query. This substitution offers the caller some level of privacy from third-party inspection and potential intrusion and censorship, although it's unclear precisely what "privacy" means when you are sharing your DNS activity with Google!

Example script:

```
#!/usr/bin/env python
import json, requests

url = "https://dns.google.com/resolve"
params = dict(
    name='www.potaroo.net',
    type='A',
    dnssec='true'
)

resp = requests.get(url=url, params=params)
data = json.loads(resp.text)
print data[u'Answer'][0][u'data']
```

Application-Level DNS

Concerns about data leakage is not limited to external forms of surveillance and interception. An appropriately paranoid application would not use the DNS resolution service of the underlying host platform, because that would release the name queries of the application into an uncontrolled environment where it may be logged and accessed by the platform and other applications.

One response to this potential for uncontrolled leakage of information is for the application to assume a greater role in performing DNS resolution. The simple approach is for the application to outsource this role to a trusted agent, and do so over a secured channel, which is where secure DNS over JSON comes in. In this case, the application is not performing DNS resolution and validation itself, but in creating a secured channel to Google's resolution service across a TLS connection the application can obtain some level of assurance that it is not performing a local leak of DNS information, and that with DNSSEC validation enabled, the responses it receives have some level of assurance that they are genuine, assuming that the name being resolved is itself DNSSEC-signed.

But perhaps even that is too much outsourced trust. Another approach is being constructed by the *GetDNS* project.^[11] In this case, it's not a secure channel to a recursive resolver that will resolve the application queries. GetDNS provides the application with a local validating resolver as a set of library calls. This project currently supports DNS over TLS. The GetDNS project operates as an open source project, and the GetDNS project page contains pointers to the code. A web application is built into the API, and a portal to a resolver implemented in this manner can be found at:

`https://getdnsapi.net/query`

This approach of pulling the DNS resolution function potentially all the way back into the application has some interesting trade-offs. The queries being made now have a source address of the local host, so the data that is leaked through the DNS queries can identify the local host. If an authoritative name server does not support a secure channel for queries using DNS over TLS, then the API will necessarily use an open unencrypted channel (at this stage the DNS over DTLS is not included in the GetDNS code base, but if someone wants to submit code ...).

On the other hand, the DNSSEC validation function can be performed locally as well by a GetDNS instance, so that the application is not forced to trust the authenticity of a bit flag in the response from a remote resolver. This way the application has direct control of the validation function, and direct knowledge of its outcome.

Between these two approaches there are further apparent trade-offs.

Making queries via a secure channel to a busy recursive resolver—and Google's Public DNS is about as busy as a DNS resolver can be—means that it is possible, to some extent, to hide behind the cache of such busy resolvers. As long as you are comfortable with sharing your DNS queries with Google, then to some extent you can use secured access to a DNS recursive resolver that intends to operate with integrity, accuracy, and completeness. The secure channel is far harder to subvert and more resistant to efforts to eavesdrop on the query stream.

If you are uncomfortable with this approach, then another option is to pull the name resolution function back into your platform and even back into the application itself using a framework such as GetDNS. The extent to which your queries may be readily visible to third parties, and the extent to which your query stream may be subverted in various ways, now depends on the capabilities of the authoritative name servers. Without name server support for DNS over TLS and possibly also potential future support for DNS over DTLS, and without DNSSEC signed zones, the local DNS resolver may still be misled in ways that may not be readily detected.

In this case, the local resolver is powerless to fix this problem, because the privacy and protection mechanisms are now in the hands of the authoritative name servers and the zone admins that are queried by the local resolver.

With DNS privacy, there is no such thing as a free lunch! But maybe in QNAME Minimisation there is the possibility of a much cheaper lunch. The queries are much the same as before, but the impact is that the query name is only progressively revealed to the name servers that are authoritative for the domain being queried.

However, an equal cause for concern is that the queries and responses are open and susceptible to eavesdropping. Here the lunch is definitely more expensive! The measures to introduce a secure channel for DNS queries and responses take a simple open query/response stateless protocol and introduce state by way of session establishment and crypto state establishment. The overheads of such additional state can be many packets and the imposition of additional Round Trip Time intervals. It is also true that no DNS privacy solution is absolute in these frameworks. While it is possible to set up a secure and encrypted channel to a recursive resolver, the implication is that the recursive resolver is privy to the query stream, even if the local network infrastructure cannot directly eavesdrop on the queries.

Alternatively, the user application can operate as a resolver itself, and attempt to direct queries to authoritative servers over a secured channel, but the authoritative servers now have visibility on your identity, and they also have to cope with a dramatic change in the query load. The authoritative name server would no longer be able to rely on recursive resolvers to absorb many of the queries, so they would presumably be exposed to a far larger query load. In addition, they would need to maintain a significant state overhead to support secured channels to these end application-based name resolvers. This really does not look like a likely scenario. As a result, it appears, at least for the moment, that this work on secured transport channels is likely to be a measure used between end-user applications and recursive resolvers.

What Does All This Mean?

While today the open nature of DNS queries makes third-party monitoring, interception, and substitution incredibly easy, there are now some grounds to be optimistic and start to contemplate a DNS environment that preserves privacy and integrity.

By performing QNAME Minimisation it is possible to radically reduce the level of leaked information coming from the DNS, and by wrapping up DNS queries and responses in a secured channel it is no longer trivial for third parties to monitor and intercept DNS queries and their responses on the wire.

If applications made use of services that push local DNS query traffic into encrypted TLS sessions with recursive resolvers, such as the service Google offers, the result would be that much of today's visible DNS would disappear from view. Not only that, but it would make the existing practices of selective local inspection and intervention in the DNS resolution process far more challenging, if not infeasible. It may be even better if authoritative name servers were to also support queries over TLS and DTLS, allowing a local host to take over the resolution function and still use encrypted query traffic services.

If this scenario were to be coupled with widespread use of DNSSEC, then it would be a somewhat different Internet from the one we have today. It's pretty obvious that national online censorship efforts will continue, and online monitoring and surveillance will also continue. But the ability to coopt the DNS into the role of an exceptionally cheap and simple means to achieve these ends will cease at some time if we collectively choose to head down this path for adding privacy and security to the DNS.

References

- [1] Wilmer van der Gaast, Carlo Contavalli, and Warren Kumari, "Client Subnet in DNS Queries," RFC 7871, May 2016.
- [2] Google Blog, "Google Public DNS and Location-Sensitive DNS Responses," December 2014,
<https://webmasters.googleblog.com/2014/12/google-public-dns-and-location.html>
- [3] APNIC Labs, "DNSSEC Validation Rate by Country,"
<https://stats.labs.apnic.net/dnssec>
- [4] Byron Ellacott and Geoff Huston, "Facebook and the GFW," August 2013,
<http://www.potaroo.net/presentations/2013-08-29-facebook.pdf>
- [5] Geoff Huston, "The Company You Keep," *The ISP Column*, June 2013.
<http://www.potaroo.net/ispcol/2013-06/company.html>
- [6] Babak Farrokhi, "Operator Level DNS Hijacking," RIPE Labs, July 2016.
https://labs.ripe.net/Members/babak_farrokhi/operator-level-dns-redirection
- [7] Stephane Bortzmeyer, "DNS Query Name Minimisation to Improve Privacy," RFC 7816, March 2016.
- [8] Eric Rescorla and Nagendra Modadugu, "Datagram Transport Layer Security Version 1.2," RFC 6347, January 2012.

- [9] Dan Wing, Tirumaleswar Reddy, and Prashanth Patil, “Specification for DNS over Datagram Transport Layer Security (DTLS),” Internet Draft, Work in Progress, December 2016, **draft-ietf-dprive-dnsodtls-15**.
- [10] Zi Hu, Liang Zhu, John Heidemann, Allison Mankin, Duane Wessels, and Paul Hoffman, “Specification for DNS over Transport Layer Security (TLS),” RFC 7858, May 2016.
- [11] GetDNS Project: **<https://getdnsapi.net>**

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: **gih@apnic.net**

JOAO LUIS SILVA DAMAS is currently Senior Researcher at APNIC. Joao was also co-founder of Hivecast Inc, a DNS services company later sold to Dyn. Previously he worked at ISC (Internet Systems Consortium) as CTO overseeing technical developments. Earlier, he served as CTO at RIPE NCC and later founded Bond Internet Systems as a consulting and research company. For around 7 years he organised the RIPE plenary program and launched the current RIPE program committee. In 2008, together with colleagues, Joao launched ESNOG to bring together Spanish ISPs to interact with each other, an activity that continues to this day. E-mail: **joao@bondis.org**

ICANN Launches Testing Platform for the KSK Rollover

The *Internet Corporation for Assigned Names and Numbers* (ICANN) is offering a testing platform for network operators and other interested parties to confirm that their systems can handle the automated update process for the upcoming *Root Zone Domain Name Systems Security Extensions* (DNSSEC) *Key Signing Key* (KSK) rollover^[1, 2]. The KSK rollover is currently scheduled for October 11, 2017.

“Currently, seven hundred and fifty million people are using DNSSEC-validating resolvers that could be affected by the KSK rollover,” said ICANN’s Vice President of Research, Matt Larson. “The testing platform is an easy way for operators to confirm that their infrastructure supports the ability to handle the rollover without manual intervention.”

Internet service providers, network operators and others who have enabled DNSSEC validation must update their systems with the new KSK. This can be done in one of two ways:

- An operator can configure a new trust anchor *manually* by obtaining the new root zone KSK from the *Internet Assigned Numbers Authority* (IANA) website at:
<https://www.iana.org/dnssec/files>
- An operator can enable a feature available in many validating resolvers that *automatically* detects and configures a new root zone KSK as a trust anchor, in which case they need take no action.

Check to see if your systems are ready by visiting:
go.icann.org/KSKtest

The KSK has been widely distributed and configured by every operator performing DNSSEC validation. If the validating resolvers using DNSSEC do not have the new key when the KSK is rolled, end users relying on those resolvers will encounter errors and be unable to access the Internet. A careful and coordinated effort is required to ensure that the update does not interfere with normal operations.

More information is available at www.icann.org/kskroll

[1] George Michaelson, Patrick Wallström, Roy Arends, and Geoff Huston, “Rolling Over DNSSEC Keys,” *The Internet Protocol Journal*, Volume 13, No 1, March 2010.

[2] ICANN Blog, “The Problem with ‘The Seven Keys,’” February 13, 2017. <https://www.icann.org/news/blog/the-problem-with-the-seven-keys>

Follow us on Twitter and Facebook



@protocoljournal



<https://www.facebook.com/newipj>

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Fabrizio Accatino
Scott Aitken
Antonio Cuñat Alario
Matteo D'Ambrosio
Jens Andersson
Danish Ansari
David Atkins
Jaime Badua
John Bigrow
Axel Boeger
Kevin Breit
Ilia Bromberg
Christophe Brun
Gareth Bryan
Stefan Buckmann
Scott Burleigh
Jon Harald Bøvre
Olivier Cahagne
Roberto Canonico
Lj Cemerar
Dave Chapman
Stefanos Charchalakakis
Greg Chisholm
Narelle Clark
Steve Corbató
Brian Courtney
Dave Crocker
Kevin Croes
John Curran
Morgan Davis
Freek Dijkstra
Geert Van Dijk
Ernesto Doelling
Karlheinz Dölger
Andrew Dul

Holger Durer
Peter Robert Egli
George Ehlers
Peter Eisses
Torbjörn Eklöv
ERNW GmbH
ESdatCo
Steve Esquivel
Mikhail Evstiounin
Paul Ferguson
Christopher Forsyth
Craig Fox
Tomislav Futivic
Edward Gallagher
Andrew Gallo
Chris Gamboni
Xosé Bravo Garcia
Kevin Gee
Serge Van Ginderachter
Greg Goddard
Octavio Alfageme Gorostiaga
Barry Greene
Martijn Groenleer
Geert Jan de Groot
Gulf Coast Shots
Sheryll de Guzman
Martin Hannigan
John Hardin
Edward Hauser
Headcrafts SRLS
Edward Hotard
Bill Huber
Hagen Hultzschn
Karsten Iwen
David Jaffe

Dennis Jennings
Jim Johnston
Jonatan Jonasson
Daniel Jones
Gary Jones
Amar Joshi
Merike Kaeo
David Kekar
Shan Ali Khan
Nabeel Khatri
Henry Kluge
Carsten Koempe
Alexander Kogan
Mathias Körber
John Kristoff
Terje Krogdahl
Bobby Krupczak
Warren Kumari
Darrell Lack
Yan Landriault
Markus Langenmair
Fred Langham
Richard Lamb
Tracy LaQuey Parker
Robert Lewis
Sergio Loreti
Guillermo a Loyola
Hannes Lubich
Dan Lynch
Miroslav Madic
Alexis Madriz
Carl Malamud
Michael Malik
Yogesh Mangar
Bill Manning

Harold March
David Martin
Timothy Martin
Gabriel Marroquin
Carles Mateu
Juan Jose Marin Martinez
Brian McCullough
Joe McEachern
Carsten Melberg
Kevin Menezes
Bart Jan Menkveld
William Mills
Thomas Mino
Mohammad Moghaddas
Charles Monson
Andrea Montefusco
Fernando Montenegro
Tariq Mustafa
Stuart Nadin
Mazdak Rajabi Nasab
Krishna Natarajan
Darryl Newman
Ovidiu Obersterescu
Mike O'Connor
Carlos Astor Araujo Palmeira
Alexis Panagopoulos
Manuel Uruena Pascual
Ricardo Patara
Dipesh Patel
Alex Parkinson
Craig Partridge
Dan Paynter
Leif-Eric Pedersen
Juan Pena
Chris Perkins

Rob Pirnie
Blahoslav Popela
Tim Pozar
David Raistrick
Priyan R Rajeevan
Paul Rathbone
Bill Reid
Justin Richards
Mark Risinger
Ron Rockrohr
Carlos Rodrigues
William Ross
Boudhayan Roychowdhury
Carlos Rubio
RustedMusic
Babak Saberi
George Sadowsky
Scott Sandefur
Arturas Satkovskis
Phil Scarr
Jeroen Van Ingen Schenau
Carsten Scherb
Roger Schwartz
SeenThere
Scott Seifel
Yury Shefer
Yaron Sheffer
Tj Shumway
Jeffrey Sicuranza
Thorsten Sideboard
Geoff Sisson
Helge Skrivervik
Darren Sleeth
Mark Smith
Job Snijders

Ignacio Soto Campos
Peter Spekrijse
Thayumanavan Sridhar
Matthew Stenberg
Adrian Stevens
Clinton Stevens
Viktor Sudakov
Edward-W. Suor
Vincent Surillo
Roman Tarasov
David Theese
Sandro Tumini
Phil Tweedie
Steve Ulrich
Unitek Engineering AG
John Urbanek
Martin Urwaleck
Betsy Vanderpool
Surendran Vangadasalam
Alejandro Vennera
Luca Ventura
Tom Vest
Dario Vitali
Randy Watts
Andrew Webster
Tim Weil
Jd Wegner
Rick Wesson
Peter Whimp
Jurrien Wijlhuizen
Pindar Wong
Bernd Zeimetz

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Diamond Sponsors



Ruby Sponsor



Sapphire Sponsors

Your logo here!

Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

ADDRESS SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

David Conrad, Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, Chairman
Internet Corporation for Assigned Names and Numbers

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

June 2017

Volume 20, Number 2

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
Automatic Certificate Management.....	2
The Root of the DNS.....	15
Fragments	26
Thank You.....	28
Call for Papers	30
Supporters and Sponsors	31

Every day we seem to read another news story about some form of *cyber attack*, be that Denial of Service incidents, ransom ware, malware, website intrusions, compromised databases, so-called *phishing*, leaked e-mails, election hacking, and much more. The underlying opportunities for such attacks are varied, ranging from human factors like easy-to-guess passwords to poorly designed and insecure technologies, as we have discussed many times in this journal. As you might expect, making the Internet more secure and robust involves numerous efforts at every layer of the protocol stack.

Encryption is a time-tested method for securing end-to-end communication as well as for storing information in a manner that prevents unauthorized access. Encryption is also used in the generation of trusted *certificates* for secure web communication. In our first article, Daniel McCarney presents an overview of the *Automatic Certificate Management Environment* (ACME).

The *Domain Name System* (DNS) is one of the core components of the Internet. We have covered many aspects of the DNS over the years, but not looked closely at the *root server system* until now. Geoff Huston describes the history and evolution of the DNS and its root servers.

As announced in the previous edition of IPJ, the *Latin America and Caribbean Network Information Centre* (LACNIC) has agreed to translate selective articles from IPJ and provide summaries in Spanish. This service is now available at: <http://lacnic.net/ipjournal>

If you have a print subscription to this journal, you will find an expiration date printed on the back cover. For the last couple of years, we have “auto-renewed” your subscription, but starting with this issue, we ask you to log in to our subscription system and perform this simple task yourself. You should receive an e-mail with instructions for how to access this system. When logged in, you can update your mail and e-mail address and change your delivery options. For any questions, e-mail us at ipj@protocoljournal.org

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

A Tour of the Automatic Certificate Management Environment (ACME)

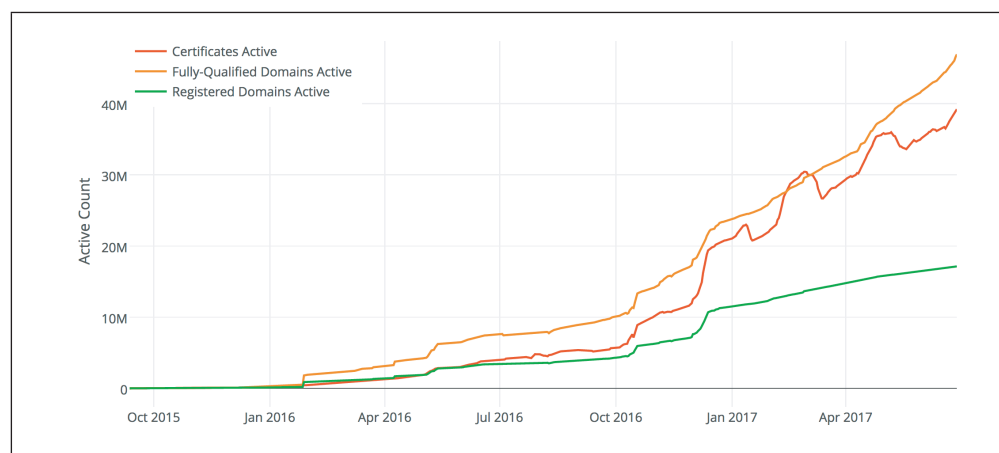
by Daniel McCarney, Internet Security Research Group

The introduction of *Let's Encrypt* has helped bolster *HyperText Transport Protocol Secure* (HTTPS) adoption by providing an easy-to-use and fully automated process for getting a trusted certificate for a domain name, all free of cost. Let's Encrypt is a service provided by the *Internet Security Research Group* (ISRG), a nonprofit organization with a mission to reduce financial, technological, and education barriers to secure communication over the Internet. To date, Let's Encrypt has issued certificates for more than 30 million websites.

Prior to Let's Encrypt, acquiring a certificate for a website was a difficult and error-prone process requiring frequent Google searches for obscure command-line incantations. Worse yet, this process typically had to be repeated manually with large periods of time elapsed between attempts—a recipe for disaster. One solution to both the usability and reliability issues created by placing this manual burden on a human is to augment the process with automation. A computer program will dutifully repeat complicated series of instructions at fixed intervals without missing a beat.

Frequently the methods of domain validation used by a certificate authority were difficult to automate at all (for example, e-mail-based validation) or required locking yourself in to a specific provider's *Application Program Interface* (API). Unlocking interoperable automation is critical to alleviating the burden on system administrators, promoting greater ease of use and helping create a fully encrypted web. Acknowledging this need for automation has been a core focus of Let's Encrypt from the very earliest days. The progress that Let's Encrypt has made in this area helps validate the premise that automation helps scale usage and has been reflected in the overall increase in HTTPS adoption observed by Mozilla Firefox telemetry, see Figure 1^[0].

Figure 1: Let's Encrypt Statistics



While much has been said about Let's Encrypt and the choice to provide the service free of cost, less attention has been directed at the work underpinning the automation aspect. Underneath the service provided by Let's Encrypt is hidden a larger effort to bring a secure, automated *Internet Engineering Task Force* (IETF) standardized protocol for certificate issuance. This protocol, called the *Automatic Certificate Management Environment* (ACME) protocol, is a multi-organizational standards effort currently in IETF Working Group *Last Call*. This article briefly describes the core aspects of ACME and some of the ecosystem that surrounds it.

ACME

In the context of IETF *Request for Comments* (RFC) documents, ACME is unfinished. To date the ACME Working Group has published six drafts as the standard has received external feedback, and lessons from Let's Encrypt have been folded back into the overall design. It is currently in a phase of development where no new major changes are expected and the primary work remaining is editorial and interoperability-related.

This article focuses on ACME as defined in **draft-07**^[1], the most current draft at the time of writing. As outlined in the next section, this focus means that our discussion of ACME will not perfectly match the Let's Encrypt service, which most closely resembles Version 3 (**draft-03**) of the draft specification for the moment while updates to the newer drafts are implemented.

The ACME Ecosystem

While ACME is a new protocol that is not yet standardized as an RFC, it already exists within an ecosystem of client and server implementations.

On the server side, Let's Encrypt has contributed an open source ACME implementation called *Boulder*^[2]. It's written in Go and as the production backend for the service provided by Let's Encrypt it is the most "battle-tested" ACME server codebase to date. At peak so far Let's Encrypt and the Boulder codebase have issued more than 1.2 million certificates in a single day. Internally Boulder is divided into subcomponents, each responsible for a portion of the responsibilities of a *Certificate Authority* (CA). Notable components include the CA, the *Validation Authority*, the *Registration Authority*, and the *Web Front End*. Components "talk" to each other using a universal *Remote Procedure Call* (RPC) framework called gRPC.

Boulder has evolved alongside the draft ACME protocol and so necessarily has some divergences where specification moved faster than implementation or vice versa. The Boulder developers document these divergences in the Boulder repository^[3] and adjust the codebase as the bounds of backwards compatibility with existing ACME clients allow. Since many ACME clients specifically target compatibility with Let's Encrypt and the Boulder implementation of ACME, most existing ACME clients will share these protocol divergences.

Our work to implement **draft-07**^[1] fully is not yet complete at the time of writing, but is expected to be available by the end of 2017 as a separate API endpoint to ease migration from older clients using the existing (largely **draft-03**-based) API.

One of the benefits of an open standard is the ease with which clients can be written to interact with Let's Encrypt and issue certificates. Better yet, if other CAs adopt ACME, these clients will be able to interact with those CAs with minimal modification. A plethora of clients have been written since the launch of Let's Encrypt targeting a variety of platforms, programming languages, and use cases. Whether you're looking to issue a single certificate from an embedded device or thousands of certificates on-demand for a large platform integration, an ACME client is available for your needs.

Most well-known of these clients is *Certbot*^[4]. Formerly known as the *letsencrypt* client, Certbot development has since been taken over by the *Electronic Frontier Foundation* (EFF) under its new name. As the first ACME client built for use with Let's Encrypt and the Boulder ACME server, it is still thought of as a reference client for Let's Encrypt. Certbot is an end-to-end solution capable of performing as much of the complicated administrative work as possible that is required to request, issue, and install a certificate for an HTTPS webserver.

Domain Validation

Let's Encrypt and the ACME protocol are both focused on *Domain Validated* (DV) certificates. *Organization Validated* (OV) and *Extended Validation* (EV) certificates are outside of the current scope of both the protocol and this article. OV and EV certificates require verification procedures that would be difficult if not impossible to automate programmatically. Unlike OV and EV certificates, a DV certificate does not attest identity but that the possessor of the private key corresponding to the public key in the certificate has demonstrated control of the domain names the certificate includes. The attestation from the CA that validated the domain control means you can be sure that your authenticated and encrypted *Transport Layer Security* (TLS) connection to the remote server is to a party in control of the domain name and not a *Man in the Middle* (MITM). It is important to know that this attestation does not vet anything about the trustworthiness of the domain owner.

Historically the technical method of domain validation that a CA employed was largely left to its own discretion and fairly ad-hoc. One method was to e-mail a token or activation link to an e-mail address believed to be authoritative for a given domain to prove ownership for an issuance request. The question of which addresses should be authoritative is the crux of this validation method, and mistakes in this decision process have led to certificates being issued to unauthorized parties in the past^[5].

Other popular methods involved generating a token to be placed in a well-known location in the HTTP webroot of the domain, or in a *Domain Name System* (DNS) record for the domain. Recent changes in the baseline requirements^[6] that CAs must meet have standardized the acceptable methods for domain validation and added some guardrails against mis-issuance.

Much of ACME directly addresses the domain validation process. This fact might be surprising if you were expecting to find a great deal of complicated cryptography related to the issuance of certificates themselves. Automating the act of issuing a certificate for a set of names is not the true challenge to scaling the web *Public Key Infrastructure* (PKI). The task of turning certificate signing requests into certificates with a software pipeline is well-understood. The larger challenge that must be addressed is how to scale the determination of whether the party requesting the certificate is authorized to act on behalf of all of the names the certificate includes. This contribution is one of the crucial ones of the ACME protocol: the introduction of clearly specified and peer-reviewed domain validation methods.

ACME Requests

All ACME requests are made over HTTPS. Protocol messages are primarily POSTed to HTTPS endpoints as *JavaScript Object Notation* (JSON) data. The JSON request data is authenticated and provided integrity through the application of *JSON Web Signatures* (JWS), as described in RFC 7515^[10]. GET requests do not have a JWS body and are not authenticated by the ACME account key; therefore only public resources are available via GET.

To provide anti-replay protection, all ACME server responses provide a *nonce* header. The value of this header must be provided in the next request to the server. A dedicated *new-nonce* endpoint also exists to request a fresh nonce without performing a throw-away request only to look at the nonce reply header.

Since JWS will not cover the *Uniform Resource Identifier* (URI) of an HTTPS request, the URI is also contained in all request bodies and must be verified by the server to ensure that an entity terminating the ACME HTTPS request (for example, a *Content Distribution Network* (CDN) or *Load Balancer*) did not modify the request URI from the one intended to be used by the client contained in the authenticated request body. The ACME draft threat model section covers these considerations with more detail.

Components of ACME

ACME was designed to be a *Representational State Transfer* (REST)ful protocol, so one way to approach understanding it is by examining the resources the protocol specifies. At its core ACME is made up of *Accounts*, *Orders*, *Authorizations*, and *Challenges*.

At a high level, issuing a certificate is a matter of creating an ACME account, submitting an order for a certificate containing a set of DNS identifiers, satisfying authorizations for each of the identifiers by solving challenges, and finally, polling the ACME server until a signed certificate satisfying the order is produced.

A special directory resource serves as the entry point for the account creation and certificate issuance flows of the ACME protocol. ACME clients identify servers by their directory URI and make an initial request to this resource in order to learn the URIs used for other resources and to get a first nonce value. The directory also contains metadata related to the ACME server (for example, terms of service requirements, *Certification Authority Authorization* (CAA) identifiers, etc.).

An example of Let's Encrypt's current `/directory` endpoint, as generated by Boulder, follows:

```
curl https://acme-v01.api.letsencrypt.org/directory
{
  "key-change": "https://acme-v01.api.letsencrypt.org/acme/key-change",
  "new-authz": "https://acme-v01.api.letsencrypt.org/acme/new-authz",
  "new-cert": "https://acme-v01.api.letsencrypt.org/acme/new-cert",
  "new-reg": "https://acme-v01.api.letsencrypt.org/acme/new-reg",
  "revoke-cert": "https://acme-v01.api.letsencrypt.org/acme/revoke-cert"
}
```

Accounts

The account resource is a container for information about a user and that user's account with the ACME server. Most importantly, an account contains a public key encoded as a *JSON Web Key* (JWK)^[11]. This public key is associated with the account at the time of account creation and is used to authenticate future requests. Like the other resources we'll see, an account is identified by its URI per usual REST practice. Account resources also contain additional metadata such as an e-mail address to contact for the account and whether a required Terms of Service Agreement has been acknowledged. In the earlier stages of the ACME draft accounts were called *registrations*, and you may still see references using this term in older material.

Accounts are created by POSTing an account resource to the new-account resource of the ACME server. Future updates (for example, to update contact information) are handled in a similar fashion. Notably you cannot view your current account information by sending a GET request to the account URI; instead you must use a POST request with an empty body. The rationale for this decision is rooted in the security model of the protocol. Only POST requests carry the required JWS to authenticate the request as coming from the account owner. If you use a JWS signed empty body in a POST request to retrieve account information, only the authorized account can view contact information.

An example POST body for a new account follows:

```
{
  "terms-of-service-agreed": true,
  "key": "...",
  "contact": [
    "mailto:example@example.com",
  ]
}
```

Orders

Orders encapsulate the request of an account for a certificate to be issued by the ACME server. The most important field of an order object is the *Certificate Signing Request* (CSR). You might be familiar with non-ACME CAs; the ACME CSR is a standard RFC 2986^[12] CSR, meaning existing tools (for example, *openssl*) can generate CSRs for use with ACME. For use within an Order the CSR is base64url-encoded, a practice used elsewhere in the protocol when binary data needs to be represented in a request.

The other important field of an Order object is the Authorizations field, containing an array of Authorization URIs. The ACME server is responsible for populating this field in the Order object returned to the client when a new order is created. Completing an order to obtain a certificate requires first completing each of the authorizations the order links to.

An example of a request body to create a new order for two DNS identifiers would resemble the following:

```
{
  "csr": "MIICmTCCAYECAQAw....cUc5i8XK-OBEMe",
}
```

Resulting in:

```
{
  "status": "pending",
  "expires": "2017-03-14T12:41:37-04:00",
  "csr": "MIICmTCCAYECAQAw....cUc5i8XK-OBEMe",
  "authorizations": [
    "/authZ/k4jO5648Y-qqrQ_F-bD6JLgtrfV4TJb6vef9GrlybvQ",
    "/authZ/c71yuTUHsuwIVeCk9B4DrsFA1MlCZMLtt4FDZ71KI20"
  ]
}
```

Presently, as described in the Boulder divergences document, Let's Encrypt does not implement the order resource. Instead clients must explicitly create authorization objects for each of the domains they wish to issue for themselves using the *new-authz* endpoint, as opposed to creating an order and receiving the URI of authorizations required from the server.

Authorizations

Authorizations are the core of the domain validation process in ACME. For an account to receive a certificate valid for an identifier, the CA needs to verify control of that identifier. If control is established, then the account is said to be authorized to request a certificate valid for the domain. In ACME an authorization starts its life in a *pending* status, indicating that the account has not yet completed the authorization process. In order to progress from the pending state to a *valid* state, the account holder must complete a set of required challenges. Authorizations also contain an expiry date, and both pending and valid authorizations fall out of usefulness after their expiry date. In the case of pending authorizations, this requirement keeps the challenges fresh. In the case of valid authorizations, it means that control must be reestablished through a fresh authorization and new challenges if the expiry has passed.

An example authorization follows:

```
{
  "status": "pending",
  "identifier": {
    "type": "dns",
    "value": "www.example.com"
  },
  "challenges": [
    {
      "type": "dns-01",
      "token": "T50nPYe3YNdKeqlqaelegDVftLpqG5D8klP_K7inCHY",
      "status": "pending",
      "error": {}
    }
  ],
  "expires": "2017-03-14T12:41:37-04:00"
}
```

Challenges

One or more challenges are embedded directly into authorizations and are identified by a type and a URI. Solving a challenge of an authorization will demonstrate the ACME account key holder's control over the identifier the authorization refers to, allowing issuance for that identifier. Each challenge type has its own method for demonstrating control, but all share the use of a random *token* and a key authorization.

The token is a random value used to identify the challenge. The token is always expressed in the *base64url* alphabet used throughout ACME, and to facilitate the usage in various challenge types it must not contain any padding characters.

The challenge key authorization is used to concretely link a specified ACME account key with the challenge for the purpose of validating an identifier. It is created by concatenating the random token present in the challenge and the Base64 URL encoding of the JWK thumbprint of the ACME account. The key authorization is provided in the subsequent JWS signed request from the ACME client to update a challenge, asking the server to attempt to verify control by performing the challenge verification process as required by the challenge type.

An example challenge follows:

```
{
  "type": "dns-01",
  "token": "T50nPYe3YNdKeqlqaelegDVftLpqG5D8klP_K7inCHY",
  "status": "pending",
  "error": {}
}
```

Getting a Certificate

After challenges have been completed successfully for each of the authorizations embedded in an order resource, the order is considered valid and the certificate can be issued. The ACME server proactively monitors order resources, and when an order is ready to be issued, it is responsible for issuing a certificate matching the domains from the CSR/authorizations. The order resource is then updated with a URI at which the client can download the issued certificate. After a client completes all of the authorizations the order requires, a polling state can be entered so the certificate URI can be added to the order to allow fetching the produced certificate.

Presently, as described in the Boulder divergences document, Let's Encrypt does not implement the order resource, so the issuance process is slightly different. Instead, clients must explicitly create authorization objects for each of the domains they wish to issue for themselves using the *new-authz* endpoint. After the authorizations are validated by completing challenges, the client can submit a CSR to the *new-cert* endpoint and will receive a certificate as a response provided the server is able to validate that the correct unexpired authorizations are in place.

Challenge Types

The ACME standard defines four distinct challenge types, each identified by the draft that it was introduced in: HTTP-01, DNS-01, TLS-SNI-01, and TLS-SNI-02. An additional *Out-of-Band* (OOB) challenge exists for integration with existing CAs. Let's Encrypt and Boulder presently do not implement TLS-SNI-02 or the OOB challenges.

HTTP-01 Challenges

The HTTP-01 challenge allows authenticating a domain by making an externally visible change to the domain website. The primary idea is that the ACME client must sign the requested key authorization and place the result in a pre-specified location in the domain webroot. The name of the file is the token value from the challenge, and the contents of the file will be the same computed key authorization that is included in the JWS signed POST body asking the server to validate the challenge.

For ACME, the pre-specified location for the challenge file is in `/.well-known/acme-challenge/`, a prefix registered with the *Internet Assigned Numbers Authority* (IANA) for the purpose of ACME domain validation. When the challenge is POSTed by the ACME client with the correct key authorization, the ACME server will make a GET request to this location on the domain referenced in the challenge authorization. Using the HTTP response, the server can validate the contents of the HTTP challenge file. If the correct key authorization was present at the correct location and signed by the correct ACME account key, then the challenge is completed and the account is considered to possess a valid authorization for this domain until the point at which it expires. Both the challenge and authorization objects are updated server side with a valid status and an expiry date.

The HTTP-01 challenge is a great fit if you are already running a world-accessible webserver on port 80 of your domains. Since the challenge requests are standard HTTP requests and will always be directed to a well-known path prefix, it is possible to implement more complex validation systems with ease. For instance, you could use the URL rewriting capabilities of a webserver to divert HTTP-01 challenge requests to a centralized server responsible for Let's Encrypt challenge validation. Many ACME clients (Certbot included) can start up a standalone HTTP server explicitly for the purpose of solving HTTP-01 challenges; this feature may be beneficial for issuing certificates for non-HTTPS services like the *Extensible Messaging and Presence Protocol* (XMPP) without needing to configure a heavyweight HTTP server.

DNS-01 Challenges

The DNS-01 challenge is conceptually similar to the HTTP-01 challenge but instead of provisioning a file at a well-known location the challenge responder provisions a TXT record at a well-known label.

For ACME, the required record is a TXT record for the label `_acme-challenge.` concatenated onto the domain being authorized. Rather than placing the entire key authorization as the value of this TXT record, the DNS-01 challenge asks that a SHA256 digest of the computed key authorization be used as the TXT record value.

When the ACME client POSTs the challenge with the JWS signed key authorization, the ACME server will verify the details of the key authorization and token match, and proceed to validate the TXT record by issuing a DNS query against one of the authoritative DNS servers for the domain being authorized. If the contents of that TXT record match the expectation of the server of the SHA256 of the challenge key authorization, then the account is considered to possess a valid authorization for this domain.

The DNS-01 challenge is often used in situations where ports 80 and 443 are not globally accessible (for example, because of corporate firewall policies), ruling out the use of HTTP-01 and TLS-SNI-02 challenges. Since the DNS-01 challenge requires only that a TXT record be updated, there's no requirement for a direct connection to the domain name that a certificate is to be issued for. Instead the challenge is validated through a query to the authoritative nameservers. The DNS-01 challenge is also well-suited to centralized management of certificate issuance. Many DNS providers support programmatic updates through an API, or with more traditional dynamic DNS updates through *nsupdate*-like tools. Using an ACME client that exposes hooks for adding and removing the required TXT records makes it easy to centrally issue certificates by automatically adjusting DNS as required by the challenges. Certbot presently supports DNS-01 in a manual-only mode, but some other ACME clients have fully automatic support with a variety of DNS providers.

TLS-SNI-02 Challenges

The TLS-SNI-02 challenge is perhaps the most unfamiliar of the ACME challenge types. For this challenge type the requester must configure a TLS server accessible at the domain to be authorized such that it will use a special self-signed certificate when processing TLS requests with a specific *Server Name Indication* (SNI)^[14] value.

The ACME client creates the self-signed certificate when it wishes to use a TLS-SNI-02 challenge to authorize a domain. The contents of the certificate are unimportant except for one crucial detail: the certificate must have two special *Subject Alternate Name* (SAN) values.

The first SAN value is a domain of the form `x.y.token.acme.invalid`, where `x` and `y` are computed as the SHA256 digest of the challenge token value, split into two labels. The second SAN value is a domain `x.y.ka.acme.invalid`, where `x` and `y` are computed as the SHA256 digest of the key authorization, split into two labels.

The TLS server used for responding to the TLS-SNI-02 challenge should be configured such that it returns the crafted challenge certificate whenever a TLS request arrives with the SNI value of the first SAN (for example, `x.y.token.acme.invalid`).

When the ACME client POSTs the challenge to begin the validation process, the ACME server will compute both SAN entries the same way the client did, and will send a TLS request to the domain using a SNI value of the first computed SAN. The ACME server can then validate that the challenge server presents a self-signed certificate with the two required SAN values verifying the challenge token and the key authorization.

The TLS-SNI-02 challenge is a good fit for environments where a webserver is already configured for HTTPS and you do not want to accept HTTP requests for HTTP-01 challenges or place files in the webroot of the domains. Similar to HTTP-01, Certbot and some other ACME clients can run a standalone TLS server for the purpose of solving TLS-SNI-02 challenges in place without requiring a heavier-weight server. The TLS-SNI-02 challenge uniquely employs the mechanics of certificates and TLS in order to provide authorization for the issuance of certificates; this symmetry of process and result is unique and satisfying from the perspective of an interested engineer.

The astute reader will note that unlike HTTP-01 and DNS-01, the TLS-SNI-02 challenge is on its second revision. Let's Encrypt and the Boulder codebase still use the original TLS-SNI-01 challenge from earlier drafts, but it suffers from one design flaw whereby all of the information required to complete the challenge was present in the request. This situation allows for a broken design where a TLS-SNI-01 challenge response server could be built that automatically replies to a challenge request without *a priori* knowledge of the challenges. To combat this design and its unintended security implications, the TLS-SNI-02 challenge requires that the key authorization, which isn't part of the challenge request, be returned as part of the challenge response.

What's Next?

We've covered the core of the ACME protocol, but the existing drafts have a great deal more information. Readers are encouraged to investigate the key rollover and revocation features of ACME from the standard since they were not covered in this article. Similarly the draft content offers more in-depth coverage of security considerations that may be of interest to readers with the hacker mindset.

The ACME standardization process is still underway. The IETF Working Group has most recently published **draft-07** and is undergoing a last-call process for interoperability testing. Sometime after this point we can expect the ACME draft will proceed to full RFC standard status.

Plenty of work remains to upgrade existing ACME servers and clients to support the latest iterations of the draft since most of the ecosystem is presently implementing ACME closer to the **draft-03** standard. Let's Encrypt intends to support the newer draft and final RFC version as independent directory endpoints alongside the current legacy **draft-03** era endpoint. This support will allow clients to gradually adopt support for the newest protocol features while continuing to renew legacy certificates produced with the **draft-03** endpoint.

The ACME protocol itself has left room for future improvements. Work is underway to develop a companion document^[7] describing additions to the CAA standard, RFC 6844^[13], that would allow domain owners to specify policy related to acceptable ACME account keys or challenge types. This work could allow for, as an example, adoption of a policy whereby only DNS-01 challenges could be used to issue certificates for a given domain name using ACME.

Standardization on challenges for non-DNS identifiers—such as IP addresses—is also an avenue for future ACME work. ACME was designed to handle additional identifier types and new challenges, and it will be interesting to see how the protocol evolves to handle use cases beyond domain validation of DNS identifiers.

Development of an open standard helps move the Web towards a world where HTTPS encryption is the norm. Certificates from Let's Encrypt are one avenue available to system administrators looking to increase the security of their websites. Adoption of ACME by other CAs and tools ensures that the decision to use HTTPS doesn't induce vendor lock-in and allows users the chance to change providers without abandoning automation. The future of ACME is still being written, and it's not too late to participate in the IETF Working Group^[8]. Readers are encouraged to subscribe to the mailing list^[9] and provide feedback as they envision integrating ACME into their own software and environments.

References

- [0] Let's Encrypt Statistics: <https://letsencrypt.org/stats>
- [1] Jacob Hoffman-Andrews, James Kasten, and Richard Barnes, "Automatic Certificate Management Environment (ACME)," Internet Draft, work in progress, **draft-ietf-acme-acme-07**, June 2017.
- [2] Github Repository for Boulder: <https://github.com/letsencrypt/boulder>
- [3] "Boulder divergences from ACME," <https://github.com/letsencrypt/boulder/blob/e81f7477a3169f77fd7247a6cdb8822fb29433aa/docs/acme-divergences.md>

- [4] “Automatically enable HTTPS on your website with EFF’s Certbot, deploying Let’s Encrypt certificates,”
<https://certbot.eff.org/>
- [5] Wayne Thayer, “Information about SSL Bug,” Godaddy Blog,
<https://www.godaddy.com/garage/godaddy/information-about-ssl-bug/>
- [6] CA Browser Forum, “Ballot 169, Revised Validation Requirements,” <https://cabforum.org/2016/08/05/ballot-169-revised-validation-requirements/>
- [7] Hugo Landau, “CAA Record Extensions for Account URI and ACME Method Binding,” Internet Draft, work in progress, [draft-ietf-acme-caa-01](#), February 2017
- [8] ACME Working Group Charter,
<https://datatracker.ietf.org/wg/acme/charter/>
- [9] ACME Mailing List Archive,
<https://mailarchive.ietf.org/arch/browse/acme/>
- [10] Nat Sakimura, Michael Jones, and John Bradley, “JSON Web Signature (JWS),” RFC 7515, May 2015.
- [11] Michael Jones, “JSON Web Key (JWK),” RFC 7517, May 2015.
- [12] Burt Kaliski, “PKCS #10: Certification Request Syntax Specification Version 1.7,” RFC 2986, November 2000.
- [13] Rob Stradling and Phillip Hallam-Baker, “DNS Certification Authority Authorization (CAA) Resource Record,” RFC 6844, January 2013.
- [14] Donald Eastlake 3rd, “Transport Layer Security (TLS) Extensions: Extension Definitions,” RFC 6066, January 2011.
- [15] Josh Aas, “Wildcard Certificates Coming January 2018,”
<https://letsencrypt.org/2017/07/06/wildcard-certificates-coming-jan-2018.html>

DANIEL MCCARNEY is a developer for the *Internet Security Research Group* (ISRG), where he works full-time on *Boulder*, the server-side software powering the Let’s Encrypt certificate authority. Prior to the ISRG Daniel was a security architect for a large content delivery network and focused on TLS and application security. He has a Masters in Computer Science from Carleton University, where his research touched both Android system security and password managers. Daniel resides in Montréal, Canada, where he enjoys long snowy walks with his dog Bart. Daniel can be reached at: cpu@letsencrypt.org

The Root of the Domain Name System

by Geoff Huston, APNIC

Few parts of the *Domain Name System* (DNS) are filled with such levels of mythology as its *root server system*. In this article I will explain what it is all about and ask the question whether the system we have is still adequate, or if it's time to think about some further changes.

The namespace of the DNS is a hierarchically structured label space. Each label can have an arbitrary number of immediately descendant labels and only one immediate parent label. Domain names are expressed as an ordered sequence of labels in left-to-right order starting at the terminal label and then enumerating each successive parent label until the root label is reached. In domain name expressions, the ASCII period character denotes a label delimiter. *Fully Qualified Domain Names* (FQDNs) are names that express a label sequence from the terminal label through to the apex (or *root*) label. In FQDNs this root is expressed as a trailing period character at the end of the label sequence. But there is a little more than that, and that's where the hierarchal structure comes in. The sequence of labels, as read from right to left, describes a series of name delegations in the DNS. If we take an example DNS name, such as `www.example.com`, then `com` is the label of a delegated zone in the root. Here we will call a *zone* the collection of all defined labels at a particular delegation point in the name hierarchy. The label "`example`" is the label of a delegated zone in the `com.` zone. And `www` is a terminal label in the `www.example.com.` zone.

But that is not all there is to the DNS. There are many more subtleties and possibilities for variation, but as we want to look specifically at the root zone, we're going to conveniently ignore all these other matters here. If you are interested, RFC 1034^[1] from November 1987 is still a good description of the way the DNS was intended to operate, and the recently published RFC 7719^[2] provides a good compendium of DNS jargon.

The most common operation performed on DNS names is to *resolve* the name; resolving is an operation to translate a DNS name to a different form that is related to the name. This most common form of name resolution is to translate a name to an associated *IP address*, although many other forms of resolution are also possible. The resolution function is performed by agents termed *resolvers*, and they function by passing queries to, and receiving results from, so-called *name servers*. In its simplest form, a name server can answer queries about a particular zone. The name itself defines a search algorithm that mirrors the same right-to-left delegation hierarchy.

Continuing with our simple example, to resolve the name `www.example.com.`, we may not know the IP addresses of the authoritative name servers for `example.com.`, or even `com.` for that matter. To resolve this name, a resolver would start by asking one of the *root zone name servers* to tell it the resolution outcome of the name `www.example.com.` The root name server will be unable to answer this query, but it will refer the resolver to the `com.` zone, and the root server will list the servers for this delegated zone, as this delegation information is part of the DNS root zone file for all delegated zones. The resolver will repeat this query to one of the servers for the `com.` zone, and the response is likely to be the list of servers for `example.com.` Assuming `www` is a terminal label in the `example.com.` zone, the third query, this time to a server for the `example.com.` zone, will provide the response we are seeking.

In theory, as per our example, every resolution function starts with a query to one of the servers for the root zone. But how does a resolver know where to start? What are the IP addresses of the servers for the root zone?

Common DNS resolver packages include a local configuration fragment that provides the DNS names and IP addresses of the authoritative name servers for the root zone. Another way is to pull down the current root hints file from <https://www.internic.net/domain/named.root>.

But it may have been some time between the generation of this list and the reality of the IP addresses of the authoritative root servers today, so the first actions of a resolver on startup will be to query one of these addresses for the name servers of the root zone, and use these name servers instead. This query is the so-called *priming query*^[3].

This priming implies that the set of root server functions includes supporting the initial bootstrap of recursive DNS resolvers into the DNS framework by responding to priming queries of the resolver, as well as anchoring the process of top-name name resolution by responding to specific name queries with the name server details of the next-level delegated zone. This role is critical in so far as if none of the root servers can respond to resolver queries, then at some point thereafter, as local caches of the resolvers expire, resolvers will be unable to respond to any DNS queries for public names. So, these root servers are important in that you may not know that they exist, or where they may be located in the net, but their absence, if that ever could occur, would definitely be noticed by all of us!

Moderating all considerations of the DNS is the issue of local caching of responses. For example, once a local resolver has queried a root server for the name `www.example.com.`, it will have received a response listing the delegated name servers for the `com.` zone.

If this resolver were to subsequently attempt to resolve a different name in the `com.` zone, then for as long as the `com.` name servers are still held in the resolver cache, the resolver will use the cached information and not query any root server. Given that the number of delegated zones in the root zone is relatively small (1,528 zones as of the start of 2017), then a busy recursive resolver is likely to assemble in its local cache the name servers of many of the top-level domain names. Then one would expect that it would have no further need to query the root name servers, except as required occasionally to refresh its local cache, assuming that it is answering queries about DNS names that exist in the DNS.

In that respect, the root servers would not appear to be that critically important in terms of the resolution of names, and certainly not so for large recursive name servers that have a large client population and therefore have a well-populated local cache. But this conclusion would not be a good one. If cached information of a recursive resolver for a zone has expired, it will need to refresh the cache with a query to a root server. At some point, all of the locally cached information will time out of the cache, and then the resolver will no longer be able to respond to any DNS query. To keep the DNS operating, recursive resolvers need to be able to query the root zone, so there is a requirement that collectively the root servers always need to be available.

In this respect, the root servers “anchor” the entire DNS system. They do not participate in every name resolution query, but without their feed of root zone information into the caches of recursive resolvers the DNS would stop. So these servers are important to the Internet, and it might be reasonable to expect a role of such importance to be performed by hundreds or thousands of such servers. But there are just 13 such root server systems.

Why 13?

The primary reason to have more than a single root server, and use multiple root servers, was diversity and availability. The root servers are intentionally located in different parts of the network, within different service provider networks. The intended objective is that in the case where a DNS resolver is incapable of contacting a root name server, then unless the resolver was itself completely isolated from the Internet, then the desired number of root servers was such that the likelihood that it could not reach any of the root name servers was considered to be acceptably low. By this reasoning, two is probably not enough, and three could well be insufficient as well, but perhaps hundreds or thousands of such root servers may well be a case of overkill!

This line of thought assumes that each named root server has a unique name, a unique IP address, and a single location. But perhaps we are assuming too much. There is a technique that places identically named and addressed servers at various locations across the Internet, called *anycast*^[5,6].

Using anycast, a user attempting to send an IP packet to an anycast service would be directed to the “closest” instance of the family of servers that share a common anycast IP address. Why not just use anycast for a collection of root servers and put as many root servers as we want behind a single IP address?

For a considerable time, anycast was viewed with some caution and trepidation, particularly in the days before *Domain Name System Security Extensions* (DNSSEC) of a signed root zone. What would stop a hostile actor from setting up a fake root server and publishing incorrect DNS information if the IP addresses the root servers used could be announced multiple times from any arbitrary location? There was also some doubt that the *Transmission Control Protocol* (TCP) would be adequately robust in such anycast scenarios. The original conservative line of thinking was that we needed multiple unitary DNS root zone servers, each with its own unique IP address announced from known points in the network fabric.

But needing “multiple” DNS root zone servers and coming up with the number 13 appears to be somewhat curious. It seems such an odd limitation in the number of root servers given that a common general rule in computer software design is Willem van der Poel’s *Zero, One, or Infinity Rule*, which states a principle in computer science that either an action or resource should not be permitted (zero), should happen uniquely (one), or should have no arbitrary limit at all (infinity). For root servers, it appears that we would like more than one root server. But why set the limit to 13?

The reason may not be immediately obvious these days, but when the DNS system was designed, the size limit of DNS responses using the *User Datagram Protocol* (UDP) was set to 512 bytes (Section 2.3.4 of RFC 1035). It seems a ludicrously small limit these days, but you have to also account for the fact that the requirement for IPv4 hosts was (and still is) that it accepts IPv4 packets up to 576 bytes long^[4]. Working backwards, that would imply that if you account for a 20-octet IPv4 packet header and an 8-byte UDP header, then the UDP payload could be up to 548 octets long, but no longer if you wanted some degree of assurance that the remote host would accept the packet. If you also allow for up to 40 bytes of IP options, then in order to ensure UDP packet acceptance under all circumstances the maximal UDP payload size should be 508 octets. The DNS use of a maximum payload of 512 bytes is not completely inconsistent with this assumption, but it is off by 4 bytes in this corner case!

This 512-byte size limit of DNS packets still holds, in that a query without any additional signal—that is, in today’s terms, a query that contains no DNS extension mechanisms that signal a capability to use a larger UDP response size—is supposed to be answered by a response with a DNS payload no greater than 512 octets long. If the actual response would be greater than 512 octets, then the DNS server is supposed to truncate the response to fit within 512 octets, and mark this partial response as *truncated*.

If a client receives a truncated response, then the client may repeat the query to the server, but use TCP instead of UDP, so that it could be assured of receiving the larger response.

The desire in the design of the DNS priming query and response was to provide the longest possible list of root name servers and addresses in the priming response, but at the same time ensure that the response was capable of being passed in the DNS using UDP, and not rely on the use of any form of optional DNS extension mechanism. The largest possible set of names that could be packed in a 512-octet DNS response in this manner was 13 such names and their IPv4 addresses—so there are at most 13 distinct root name servers in order to comply with this limit.

These days every root name server has an IPv6 address as well as an IPv4 address, so the DNS priming response that lists all these root servers and their IPv4 and IPv6 addresses is now 811 octets. If the resolver also requests that the response should include the DNSSEC signatures, then the size of the response would expand to 1,097 bytes. But if you pass a simple priming query to a root server without a UDP buffer size extension in the query, then you will still receive no more than 512 octets in response. The size-limited response will still list the names of all 13 root name servers, but will not list all of their IPv4 and IPv6 addresses in the additional section of the response.

The partial set of these additional records of root server names and their IPv4 and IPv6 addresses is passed back without any particular indication of what is missing. The decision as to which records to include and which to omit to meet the size restriction also varies between root name servers. Some root name servers provide the IPv6 addresses of root servers A through J in a 508-byte response, while others give all 13 IPv4 addresses and add the IPv6 addresses of A and B in a 492-byte response. The remainder provide the IPv4 and IPv6 addresses for A through F and the IPv4 address of G in a 508-byte response. I suppose that the message here is that recursive resolvers should support the *Extension Mechanisms for DNS* (EDNS(0)) as specified in RFC 6891^[14], and offer a UDP buffer size that is no less than 1,097 bytes if they want a complete DNSSEC-signed response to a root zone priming query.

However, even then the story is incomplete. These additional records are not DNSSEC-signed in the priming response, so if a resolver wants to assure itself that the IP addresses that are provided in this response are the actual IP addresses of the root servers, it needs to separately query these names and request DNSSEC credentials in the response. However, as of the time of writing of this article the zone **root-servers.net** is not DNSSEC-signed, so right at the heart of the DNS there is still a leap of faith that all resolvers need to make in order to link into the DNS through the priming process.

We are also entirely comfortable with anycast these days, and the root server system has enthusiastically adopted anycast, where most of the root servers are replicated in many locations. The overall result is that hundreds of locations host at least one instance of one of the root server anycast constellations, and often more. Part of the reason that our comfort level with anycast has increased is the use of a DNSSEC-signed zone, and recursive resolvers should be protecting their clients by validating the response they receive to ensure that they are using the genuine root zone data, to the extent that this data has been signed in the first place.

Should we do more?

It would certainly make some sense to sign the **root-servers.net** zone to further protect recursive resolvers from being led astray.

But what about the specification of 13 unique root server names and their associated anycast constellations? If we had more root servers, would it make everything else better? Should we contemplate further expanding these anycast constellations into thousands or even tens of thousands of root servers? Should we open up the root server letter set to more letters? Is there a limit to “more” or “many”? Where might that limit be, and why?

These days the response that recursive resolvers receive in 512 bytes or less is a partial view of the root name server system. From that perspective, 13 is not a practical protocol-derived ceiling on the number of distinct root server letters. Whether the partial response in 512 bytes reflects 6, 10, or 13 root name servers out of a pool of 13 or 14 or any larger number is largely no longer relevant. The topic has moved beyond a conversation about any numeric ceiling on the letter count into a consideration of whether more root server letters would offer any incremental benefit to the Internet, as distinct from the current practice of enlarging the root server anycast constellations. Indeed, rather than more root name servers, whether by adding more letters or enlarging anycast constellations, should we consider alternative approaches to the DNS that can scale and improve resilience under attack through answering root queries but not directly involving these root name servers at all? In other words, can we look at DNS structures that use the root servers as a distribution mechanism for the root zone data and use the existing recursive resolver infrastructure to directly answer all queries that relate to data in the root zone?

The reason to contemplate this question is that it is not clear that more root server letters or more root server anycast instances, or even both measures, make everything else better. Reducing the latency in querying a root name server has only a minimal impact for end users.

The design objective of the DNS system is to push the data as close to the user as possible in the first place, so that every effort is made to provide an answer from a local resolver cache.

It is only when there is a cache miss that the resolver query will head back into the authoritative DNS server infrastructure, a situation that would normally affect only a very small proportion of queries over time. The DNS derives its performance and efficiency through resolver caches, so the overall intention is to limit the extent to which resolvers query these root name servers to the minimal level possible.

Secondly, a local root name server may not necessarily provide any additional name resolution resilience in the case of local network isolation. Secondary root name servers also have an expiry time on the data they serve, and in the case of extended isolation the server will also time out a case to be able to respond. This timeout is as true for the root zone as it is for any other zone.

In many ways, the net effect of a local root name server on local users' Internet experience is minimal, and could well pass completely unnoticed in many cases.

In terms of the primary objectives of the root name server system, diversity and availability, there is little to be gained by adding additional root name letters. A significant expansion of the number of uniquely named root servers would ultimately make a complete priming response exceed 512 bytes, meaning either forcing all priming queries into TCP by signalling that the UDP response was truncated, or dropping some named root servers from a non-EDNS(0) priming query response.

But rather than resisting the hard limits imposed by protocol specifications in some early RFCs, perhaps we are asking the wrong question. Rather than trying to figure out how to field even more instances of root servers and keep them all current, there is perhaps a different question: Why do we need these special dedicated root zone servers at all?

If the only distinguishing feature of these root servers is the proposition that any response with a source address of any of these 26 distinguished IP addresses is by simple unfounded assertion the absolute truth, then it is laughably implausible. Anyone who has experienced DNS interceptors would have to agree that DNS lies are commonplace, and nation states as well as service providers across the entire Internet practice lying.

Enter DNSSEC

The DNSSEC-signing of the root zone of the DNS introduced further possibilities to the root zone service to resolvers. If a resolver has a validated local copy of the current *Key Signing Key* (KSK), then it can independently validate any response provided to it from any signed zone that has a chain of signing back to this KSK, including of course any signed response about the root zone itself.

A validating resolver no longer needs to obsess that it is querying a genuine root name server, and no longer needs to place a certain level of blind faith in the belief that its DNS queries are not being intercepted and that faked responses are not being substituted for the actual response. With DNSSEC it simply does not matter in the slightest how you get the response. What matters is that you can validate responses with your local copy of the root zone key. If you can perform this validation successfully, then the answer is much more likely to be genuine!

The ubiquitous use of DNSSEC casts the root server system in an entirely different light, and the relationship between recursive resolvers and the root servers can change significantly.

A relevant observation here is that some 75% of responses from the root zone are “no such domain” NXDOMAIN responses (for example,^[7]). Recursive resolvers could absorb much of the root server query load and answer these queries directly with NXDOMAIN responses if they used this form of response synthesis. The way resolvers could answer the queries is to use so called “aggressive NSEC caching^[11].” This approach uses the *Next Secure* (NSEC) records provided in the responses relating to the nonexistence of a name in the root zone to allow recursive resolvers to synthesise an authoritative NXDOMAIN response for queries relating to any name in the range specified in the NSEC data. Rather than caching a root zone NXDOMAIN answer for each individual nonexistent domain name, caching the NSEC response allows the recursive resolver to cache a common signed response for the entire span of query names as described in each NSEC response. With a cache of 1,528 defined top-level domains and another 1,528 NSEC records, a recursive resolver would be able to provide authoritative responses for any query that would otherwise be passed through to a root server.

Another approach is to use *local secondaries* for the root zone. This approach is not an architectural change to the DNS, or at least not intentionally so. For recursive resolvers that implement this approach, this change is a form of change in query behaviour in so far as a recursive resolver configured in this manner will no longer query the root servers for queries it would normally direct to an instance of the root. Instead, it directs these queries to a local instance of a slave server that is listening on the loopback address of the recursive resolver. This slave server is serving a locally held instance of the root zone, and the recursive resolver would perform DNSSEC validation of responses from this local slave to ensure the integrity of responses received in this manner. In effect, this technique loads a recursive resolver with the entire root zone into what is functionally similar to a local secondary root zone server cache. For users of this recursive resolver there is no apparent change to the DNS or to their local configurations. Obviously, there is no change to the root zone either.

This proposal provides integrity in the local root server through the mechanism of having the recursive resolver perform DNSSEC validation against the responses received from the local root slave. If the recursive resolver is configured as a DNSSEC-validating resolver, then this mechanism is configurable on current implementations of DNS recursive resolvers.

The advantage here is that the decision to set up a local slave root server or to use aggressive NSEC caching is a decision that is entirely local to the recursive resolver, and the impacts of this decision affect only the clients of this recursive resolver. No coordination with the root server operators is required, nor is any explicit notification. The outcomes are only indirectly visible to the clients of this recursive resolver, and no other.

Where does this leave the root server system?

In the light of increasing use of DNSSEC, the root server system is declining in relevance as a unique source of authoritative responses for the root zone, and we can forecast a time when their role in resolving queries would be largely anachronistic. A validated response can be considered a genuine response regarding the contents of the root zone, regardless of how the recursive resolver learned this response. It is no longer necessary to have a dedicated set of name servers running on a known set of IP addresses as the only means to protect the integrity of the root zone.

It is also true that the root servers are no longer being used as cache refresh for recursive resolvers for delegated domains. Today we see much of the time, effort and energy, and cost of root server operation being spent to ensure that NXDOMAIN answers are provided promptly and reliably. This use of time really does not make any sense these days. The use of local secondary root servers and the use of NSEC caching can remove all of these specific queries relating to undefined names to the root servers, and what would be left is the cache priming queries. If all recursive resolvers were able to use either of these measures, then the residual true role of the root server system would not be to respond to individual queries, but simply to distribute current root zone data into the resolver infrastructure.

If the functional intention of the root server system is to distribute signed root zone data to recursive resolvers, then perhaps we could look more widely for potential approaches. Regularising the times that changes are made to the root zone would help reduce opportunistic polling of the root servers to detect when a change might have occurred. Or using an approach based on *Incremental Zone Transfer* (IXFR) that would allow recursive resolvers to request incremental changes to the root zone based on differences between zone *Start of Authority* (SOA) numbers may be more efficient.

Maybe we can look further afield for additional ways to distribute the root zone contents. Social networks appear to be remarkably adept in their ability to distribute updates, and a thought is that the small set of incremental changes to the signed root zone would be highly amenable to similar techniques or even using the same social networks. One can readily imagine a feed of incremental root zone updates on media such as Twitter, for example!

I also can't help but wonder about the wisdom of the root zone servers being promiscuous with respect to whom they answer. Root zone query data points to some 75% of queries seen at the root zone servers generating NXDOMAIN responses, meaning that three-quarters of the responses from root servers are nonsensical questions in the context of the root zone. It's not clear to what extent the other 25% of queries reflect actual user activity. In an APNIC measurement exercise using synthetic domain names that included a time component, it was evident that more than 30% of the queries seen at the authoritative servers of the measurement reflected "old" queries, generated by query log replay or other DNS forms of stalking activities.

One way to respond to this situation is to farm out the query volume currently seen at the root servers into the existing recursive resolver infrastructure, so that all root zone responses are generated by these recursive resolvers, rather than passing queries onward to the root servers. If the root servers exclusively served some form of incremental zone transfer and did not answer any other query type directly, then we would see a shift in query traffic away from the root servers as a crucial DNS query attractor, leaving only a lower profile role as a server to recursive resolvers.

There is much to learn about the DNS, and there is still much we can do in trying to optimise the DNS infrastructure to continue to be robust, scalable, and accurate—all essential attributes to underpin the continued growth pressures of the Internet.

References and Further Reading

- [1] P.V. Mockapetris, "Domain names - concepts and facilities," RFC 1034, November 1987.
- [2] Kazunori Fujiwara, Paul Hoffman, and Andrew Sullivan, "DNS Terminology," RFC 7719, December 2015.
- [3] Peter Koch, Matt Larson, and Paul Hoffman, "Initializing a DNS Resolver with Priming Queries," RFC 8109, March 2017.
- [4] J. Postel, "Internet Protocol," RFC 791, September 1981.
- [5] Kurt Erik Lindqvist and Joe Abley, "Operation of Anycast Services," RFC 4786, December 2006.

- [6] David Oran, Dave Thaler, Eric Osterweil, and Danny McPherson, “Architectural Considerations of IP Anycast,” RFC 7094, January 2014.
- [7] <http://stats.dns.icann.org/rssac/2017/01/rcode-volume/1-root-20170130-rcode-volume.yaml>
- [8] “RSSAC023: History of the Root Server System—A Report from the ICANN Root Server System Advisory Committee (RSSAC),” November 4, 2016.
<https://www.icann.org/en/system/files/files/rssac-023-04nov16-en.pdf>
- [9] Marc Blanchet and Lars-Johann Liman, “DNS Root Name Service Protocol and Deployment Requirements,” RFC 7720, December 2015.
- [10] Paul Hoffman and Warren Kumari, “Decreasing Access Time to Root Servers by Running One on Loopback,” RFC 7706, November 2015.
- [11] Akira Kato, Warren Kumari, and Kazunori Fujiwara, “Aggressive use of DNSSEC-validated Cache,” May 2017, Internet Draft, work in progress, **draft-ietf-dnsop-nsec-aggressiveuse-10**
- [12] Geoff Huston, “Workshop on DNS Future Root Service,” December 2014.
<http://www.potaroo.net/ispcol/2014-12/futureroots.html>
- [13] DNS RFCs: <https://www.isc.org/community/rfc/dns/>
- [14] Paul Vixie, Joao Damas, and Michael Graff, “Extension Mechanisms for DNS (EDNS(0)),” RFC 6891, April 2013.

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

ISOC Issues Paper on Content Blocking

The *Internet Society* (ISOC) recently voiced its commitment to keeping the Internet on for everyone, in response to the increasing number of government orders to temporarily shut down or restrict access to Internet services. Speaking out at *RightsCon 2017*, the world's leading conference on Internet and human rights that took place in late March in Brussels, the organization underscored that any deliberate attempt to interrupt Internet communications or control the flow of information over the Internet puts society at risk.

Internet shutdowns, including those that impact social media sites or entire networks, occur when governments intentionally disrupt the Internet or mobile apps, often used in the context of elections, demonstrations or other tense social contexts. According to *Access Now*, there were 56 Internet shutdowns recorded worldwide in 2016, an upward trend from previous years.

A paper entitled “Internet Society Perspectives on Internet Content Blocking,”^[1] explores the most common Internet restriction techniques and highlights the shortcomings and collateral damage from the use of such measures. “From censorship to SMEs going out of business, the human, economic and technical costs of Internet shutdowns are just too high,” explains Nicolas Seidler, Senior Policy advisor at the Internet Society.

The paper describes and evaluates the most common content blocking techniques used by governments to restrict access to information (or related services) that is either illegal in a particular jurisdiction, is considered a threat to public order, or is objectionable for a particular audience.

According to Freedom House's *Freedom on the Net report 2016*, governments in 24 of the 65 countries assessed impeded access to social media and communication tools, up from 15 the previous year.

“Before they take action, we are calling policymakers to think twice: Internet shutdowns and content filtering are not the answer,” said Constance Bommelaer, Senior Director for Global Internet Policy at the Internet Society. “We are at a crossroads, and the actions we take today will determine whether the Internet will continue to be a driver of empowerment, or whether it will threaten personal freedoms and rights online,” added Bommelaer.

The Content Blocking paper can be downloaded in various forms and languages from ISOC's website^[1]. Quoting from the Foreword: “The use of Internet blocking by governments to prevent access to illegal content is a worldwide and growing trend. There are many reasons why policy makers choose to block access to some content, such as online gambling, intellectual property, child protection, and national security.

However, apart from issues relating to child pornography, there is little international consensus on what constitutes appropriate content from a public policy perspective.

The goal of this paper is to provide a technical assessment of different methods of blocking Internet content, including how well each method works and what are the pitfalls and problems associated with each. We make no attempt to assess the legality or policy motivations of blocking Internet content.

Our conclusion, based on technical analyses, is that using Internet blocking to address illegal content or activities is generally inefficient, often ineffective and generally causes unintended damages to Internet users.

From a technical point of view, we recommend that policy makers think twice when considering the use of Internet blocking tools to solve public policy issues. If they do and choose to pursue alternative approaches, this will be an important win for a global, open, interoperable and trusted Internet.”

[1] <https://www.internetsociety.org/doc/internet-content-blocking>

IAB Issues RFC on Protocol Adoption and Transition

The *Internet Architecture Board* (IAB) has recently published a *Request for Comments* (RFC) on Protocol Adoption and Transition^[1]. The abstract states: “Over the many years since the introduction of the Internet Protocol, we have seen a number of transitions throughout the protocol stack, such as deploying a new protocol, or updating or replacing an existing protocol. Many protocols and technologies were not designed to enable smooth transition to alternatives or to easily deploy extensions; thus, some transitions, such as the introduction of IPv6, have been difficult. This document attempts to summarize some basic principles to enable future transitions, and it also summarizes what makes for a good transition plan.”

[1] Thaler, D., Ed., “Planning for Protocol Adoption and Subsequent Transitions,” RFC 8170, May 2017.

Follow us on Twitter and Facebook



@protocoljournal



<https://www.facebook.com/newipj>

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Fabrizio Accatino
Scott Aitken
Antonio Cuñat Alario
Matteo D'Ambrosio
Jens Andersson
Danish Ansari
David Atkins
Jaime Badua
John Bigrow
Axel Boeger
Kevin Breit
Ilia Bromberg
Christophe Brun
Gareth Bryan
Stefan Buckmann
Scott Burleigh
Jon Harald Bøvre
Olivier Cahagne
Roberto Canonico
Lj Cemeraz
Dave Chapman
Stefanos Charchalakakis
Greg Chisholm
Narelle Clark
Steve Corbató
Brian Courtney
Dave Crocker
Kevin Croes
John Curran
Morgan Davis
Freek Dijkstra
Geert Van Dijk
Ernesto Doelling
Karlheinz Dölger
Andrew Dul
Holger Durer

Peter Robert Egli
George Ehlers
Peter Eisses
Torbjörn Eklöv
ERNW GmbH
ESdatCo
Steve Esquivel
Mikhail Evstiounin
Paul Ferguson
Christopher Forsyth
Craig Fox
Tomislav Futivic
Edward Gallagher
Andrew Gallo
Chris Gamboni
Xosé Bravo Garcia
Kevin Gee
Serge Van Ginderachter
Greg Goddard
Octavio Alfageme Gorostiaga
Barry Greene
Martijn Groenleer
Geert Jan de Groot
Gulf Coast Shots
Sheryll de Guzman
Martin Hannigan
John Hardin
Edward Hauser
Headcrafts SRLS
Robert Hinden
Edward Hotard
Bill Huber
Hagen Hultzschn
Karsten Iwen
David Jaffe
Dennis Jennings

Edward Jennings
Jim Johnston
Jonatan Jonasson
Daniel Jones
Gary Jones
Amar Joshi
Merike Kaeo
David Kekar
Shan Ali Khan
Nabeel Khatri
Henry Kluge
Carsten Koempe
Alexander Kogan
Mathias Körber
John Kristoff
Terje Krogdahl
Bobby Krupczak
Warren Kumari
Darrell Lack
Yan Landriault
Markus Langenmair
Fred Langham
Richard Lamb
Tracy LaQuey Parker
Robert Lewis
Sergio Loreti
Guillermo a Loyola
Hannes Lubich
Dan Lynch
Miroslav Madic
Alexis Madriz
Carl Malamud
Michael Malik
Yogesh Mangar
Bill Manning
Harold March

David Martin
Timothy Martin
Gabriel Marroquin
Carles Mateu
Juan Jose Marin Martinez
Brian McCullough
Joe McEachern
Carsten Melberg
Kevin Menezes
Bart Jan Menkveld
William Mills
Thomas Mino
Mohammad Moghaddas
Charles Monson
Andrea Montefusco
Fernando Montenegro
Soenke Mumm
Tariq Mustafa
Stuart Nadin
Mazdak Rajabi Nasab
Krishna Natarajan
Darryl Newman
Ovidiu Obersterescu
Mike O'Connor
Carlos Astor Araujo Palmeira
Alexis Panagopoulos
Manuel Uruena Pascual
Ricardo Patara
Dipesh Patel
Alex Parkinson
Craig Partridge
Dan Paynter
Leif-Eric Pedersen
Juan Pena
Chris Perkins
Rob Pirnie

Blahoslav Popela
Tim Pozar
David Raistrick
Priyan R Rajeevan
Paul Rathbone
Bill Reid
Rodrigo Ribeiro
Justin Richards
Mark Risinger
Ron Rockrohr
Carlos Rodrigues
Lex Van Roon
William Ross
Boudhayan Roychowdhury
Carlos Rubio
RustedMusic
Babak Saberi
George Sadowsky
Scott Sandefur
Sachin Sapkal
Arturas Satkovskis
Phil Scarr
Jeroen Van Ingen Schenau
Carsten Scherb
Roger Schwartz
SeenThere
Scott Seifel
Yury Shefer
Yaron Sheffer
Tj Shumway
Jeffrey Sicuranza
Thorsten Sideboard
Henry Sinnreich
Geoff Sisson
Helge Skrivervik
Darren Sleeth

Mark Smith
Job Snijders
Ignacio Soto Campos
Peter Spekrijse
Thayumanavan Sridhar
Matthew Stenberg
Adrian Stevens
Clinton Stevens
Viktor Sudakov
Edward-W. Suor
Vincent Surillo
Roman Tarasov
David Theese
Sandro Tumini
Phil Tweedie
Steve Ulrich
Unitek Engineering AG
John Urbanek
Martin Urwaleck
Betsy Vanderpool
Surendran Vangadasalam
Alejandro Vennera
Luca Ventura
Tom Vest
Dario Vitali
Randy Watts
Andrew Webster
Tim Weil
Jd Wegner
Rick Wesson
Peter Whimp
Jurrien Wijnhuizen
Pindar Wong
Bernd Zeimet

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Diamond Sponsors



Ruby Sponsor



Sapphire Sponsors

Your logo here!

Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

ADDRESS SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

David Conrad, Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, Chairman
Internet Corporation for Assigned Names and Numbers

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

November 2017

Volume 20, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

F R O M T H E E D I T O R

In This Issue

From the Editor	1
A Blockchain Tutorial.....	2
In Defence of NATs	25
Fragments	34
Thank You.....	36
Call for Papers.....	38
Supporters and Sponsors	39

Publication of this journal is made possible by numerous individuals and organizations. Every year in late August we initiate a sponsorship renewal campaign, and the total funding determines our publication frequency for the following year. This *November* edition will be the third and final issue in 2017, but we hope to return to our regular quarterly publication schedule in 2018. We still need more individual and corporate sponsors, so please make a donation at <http://tinyurl.com/IPJ-donate> or ask your company to sign up for a sponsorship.

In our series of articles on emerging technologies we turn to *Blockchain*, a term that is now found in mainstream news outlets. We asked Bill Stallings to give us an introduction to this technology and consider some of its applications, such as *Bitcoin*.

Network Address Translation (NAT) has been widely deployed in both home and corporate networks for many years. Since the IPv4 address space is largely depleted, NATs provide an easy option for creating local networks that use private address space as defined in RFC 1918, and communicate with the public Internet through a single IP address. There are many technical problems associated with NATs, some of which have been described in other articles in this journal. This time, Geoff Huston provides an opinion piece “In Defence of NATs.”

After many years on the ICANN Board, Steve Crocker has finished his term, but has agreed to continue serving on our Editorial Advisory Board. Thank you, Steve! In other news, the *Internet Society* recently celebrated its 25th anniversary, while the *Internet Engineering Task Force* (IETF) celebrated its 100th meeting. For more information on these events, visit <http://isoc.org> and <http://ietf.org>

As mentioned in our previous issue, if you have a print subscription to this journal, you will find an expiration date printed on the back cover. For the last couple of years, we have “auto-renewed” your subscription, but now we ask you to log in to our subscription system and perform this simple task yourself. This process will ensure that we have your current contact information, as well as delivery preference (print edition or PDF download). For any questions, e-mail us at: ipj@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

A Blockchain Tutorial

by William Stallings, Independent Consultant

Blockchain is a recently-developed distributed digital implementation of the hardcopy transaction *ledger* that has been used throughout the world for centuries. Businesses and other organizations use ledgers in a variety of applications, such as to determine ownership, establish valuations, and document liabilities. The most common ledger applications are for tracking and chronologically recording transactions that involve an exchange of value between parties. Another common use of ledgers is to record birth and death certificates.

Blockchain first came to public notice as the technology that supports the virtual currency *Bitcoin*. And while the interest in Bitcoin has tended to wax and wane, the interest in blockchain continues to grow^[1].

Distributed Ledgers

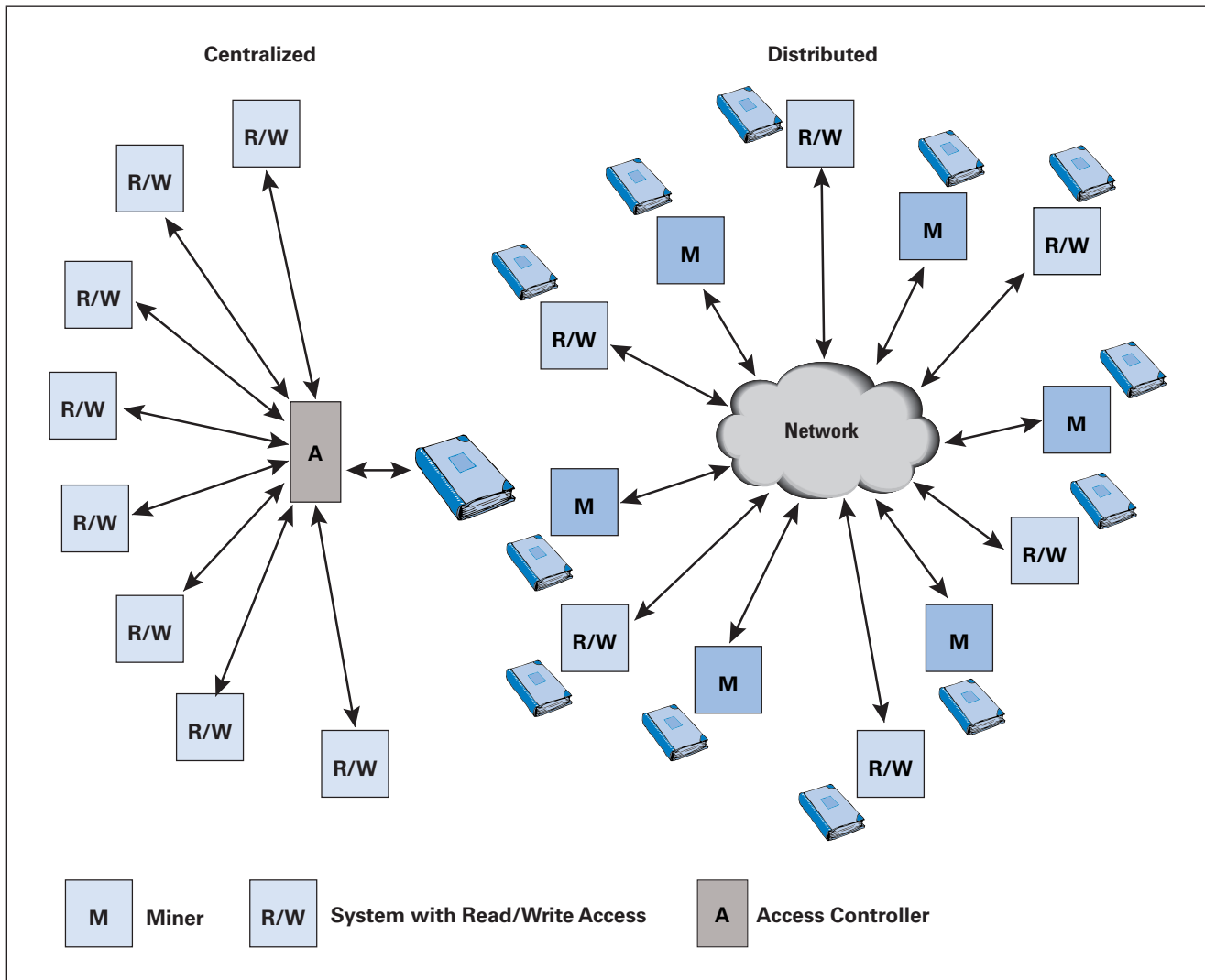
In the business context, a ledger, or *general ledger*, is defined as a central repository of the accounting information of an organization in which the summaries of all financial transactions during an accounting period are recorded. Also called the *book of final entry*, it provides all the data for preparing financial statements for the organization. As mentioned, the ledger can be used to record other types of transactions as well.

In the common business environment, a digital ledger is stored in a central server, and distributed access is provided with read and/or read/write privileges (Figure 1). To assure security, there is some sort of access control mechanism that authenticates users, enables secure access, and enforces access restrictions (for example, read-only). For a system with ongoing transactions and a heavy volume of read and write access, the central server model can be inefficient.

An alternative is a *secure distributed ledger*, which consists of an expandable list of cryptographically signed, irrevocable records of transactions that is shared by a distributed network of computers. Subject to network time delays, every participant has the same copy of the ledger. Each participant may propose a new transaction to be added to the ledger and when consensus that the transaction is valid is reached, it is added to the register.

Trust in a distributed ledger involves two concepts. First, security protocols and mechanisms, generally based on *Public-Key Cryptography*, ensure that the creator of each transaction is authenticated and validated. Transaction creators prove they are entitled to make a transaction by satisfying the particular conditions associated with this application. Meeting these conditions almost always involves the use of a secure digital signature.

Figure 1: Centralized and Distributed Ledgers



Second, a consensus mechanism is used in which computers on the network check each other to ensure records are consistent. In blockchain, this latter mechanism is implemented by systems called *miners*. Their job it is to determine that each new addition to the ledger is valid and consistent with previous entries. When the miners achieve consensus on a new entry, it is permanently added to the ledger.

Consider the use of a distributed ledger to record financial transactions or some other type of transaction that involves the exchange of value. Each transaction is a signed message that creates new outputs (transfer of value to another) while consuming old inputs (transfer of value from the transaction maker). For financial applications, each transaction is the digital equivalent of a paper check, and represents a promise by the payer to transfer control of a given amount of value to another party. The same funds or other value can be sent to only one party. An attempt at double spending, by creating two transactions that consume the same inputs, is prevented by the use of digital signatures and the trust mechanisms of the distributed ledger.

The Gartner Research document “What CIOs Should Tell the Board of Directors About Blockchain”^[2] lists the following benefits of using secure distributed ledgers:

- Civilians and computerized agents govern the economic and transaction infrastructure, which is global in scale, peer-to-peer, self-regulating, secure, and reliable.
- A decentralized, shared history of activity, obligations, rights, and records ensures transparency and certainty.
- Fine-grained and diverse (not just monetary) value exchange occurs directly between participants on a network, at lower cost and higher speed compared to legacy systems.
- The system is open to everyone, both public and private, but control and openness can be customized.
- Ownership and rights are recognized broadly. Value can be natively created and exchanged with no double spending or repudiation of transactions. The system guarantees proof of existence, process, and asset provenance.
- Embedded business logic enables dynamically self-executing smart contracts linked to diverse assets.
- Distributed autonomous organizations acting as full-fledged legal entities can execute transactions with no human intervention.

General Concept

In essence, blockchain is a data structure that makes it possible to create a digital ledger of transactions and share it among a distributed network of computers. After a block of data is recorded on the blockchain ledger, it is computationally infeasible to change or remove it. When someone wants to add to the ledger, participants in the network, all of which have copies of the existing blockchain, run algorithms to validate the proposed transaction. If a majority of nodes agree that the transaction looks valid—that is, identifying information matches the history of a blockchain—then the new transaction will be approved and a new block added to the chain. The transaction is fulfilled or executed only when it has been approved for addition to the blockchain. In contrast, in a typical computerized ledger scheme, transactions are submitted to a trusted central party that is responsible for validating the transactions and posting them in the ledger.

Blockchain provides a distributed public ledger containing transactions that are governed and maintained by specific protocols through consensus of the nodes participating in its network. The ledger consists of a linear time-sequenced chain of blocks, with each block containing one or more transactions. Each block is connected to the previous block via a hash (tamper-proof digital fingerprint). On the blockchain, users can observe transactions that have occurred, so they know which outputs are available for spending and which ones have been consumed.

Each block in the blockchain represents, in effect, the claim by someone on the network that the transactions contained inside the block are the first ones to spend the inputs involved, and therefore any transaction in the future that attempts to spend the same inputs should be rejected as invalid.

The term “blockchain” is used interchangeably to describe both the blockchain network (network of nodes) and the distributed ledger (chain of blocks). It offers a way for users who may not know or trust each other to create a record of *who* transacts *what* that will compel the assent of everyone concerned.

The blockchain ledger is not housed on a single privileged server. Rather, it is a shared data structure in which every node (user) on the network has the same copy of all other nodes (subject to propagation time delays) and can read any transaction in the ledger.

Blockchain Structure

A blockchain is a linear sequence of blocks used to store transactions. Each block contains one or more related transactions, and the blocks are ordered in increasing time sequence. Thus, each block represents a set of events that have occurred over a given time frame that is subsequent to the preceding block in the chain and prior to the following block in the chain. Users with application access to the chain can read any transaction in the sequence and can add a new block at the end of the sequence.

As shown in Figure 2, each block has a unique predecessor and successor. A block is added only at the newer, or higher end of the chain. As will be shown, there may temporarily be a branching structure as the chain grows. An essential element of blockchain is that each block is linked to its preceding block using a cryptographic algorithm. The scheme is designed such that it is computationally feasible to add a new block to the end of the chain but computationally infeasible to replace a block interior to the chain or to insert a new block between two existing blocks in the chain. After a block is added to the chain, it is read-only. Figure 3 shows the blockchain operation in general terms.

Figure 2: Block Chaining Concepts

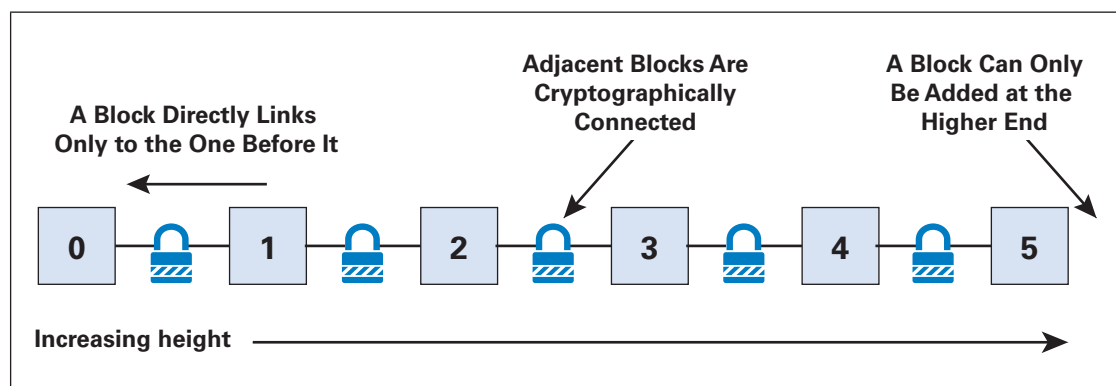
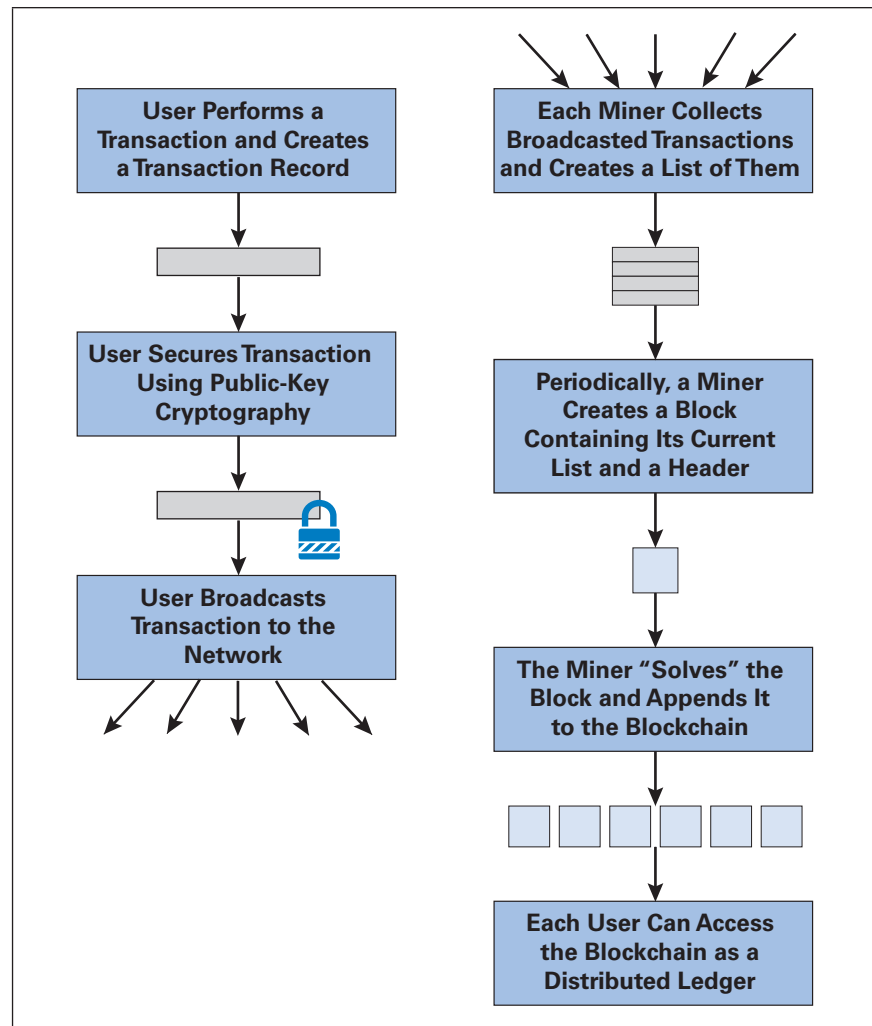


Figure 3: Basic Blockchain Logic



The exact structure of a block may vary from one application to another. Table 1 shows the typical block format. Each block begins with a “magic number” that uniquely identifies this chain. For Bitcoin, the magic number is 0xD9B4BEF9. This number is followed by a *blocksize* field that specifies the total number of bytes in the remainder of the block. Next comes the header, consisting of multiple fields. Finally, the block contains a transaction counter (≥ 1) followed by one or more transactions. The internal format of each transaction is application-dependent.

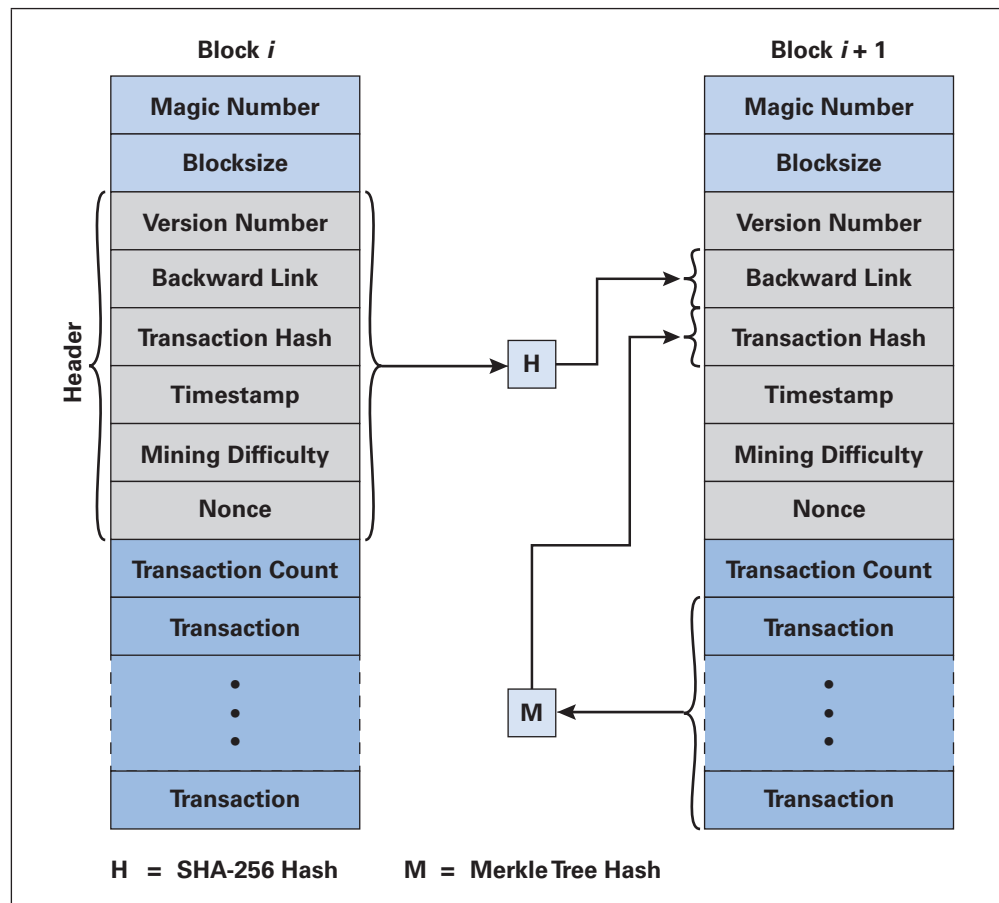
The header begins with a *Version Number*, to allow for future alterations to the block format. The blockchain application should be backward compatible so that older format versions can be processed. The foundation of the security of blockchain is found in the second field, which in effect provides a *Backward Link* to the preceding block. This backward link consists of the hash of all of the headers of the preceding block (Figure 4). By using a cryptographically strong hash function, such as SHA-256, this scheme secures the blockchain against an adversary’s altering a block or inserting a block.

In either case, the adversary would have to create a block with a header whose hash value equals a given value, and this creation is computationally infeasible for SHA-256.

Table 1: Contents of a Block

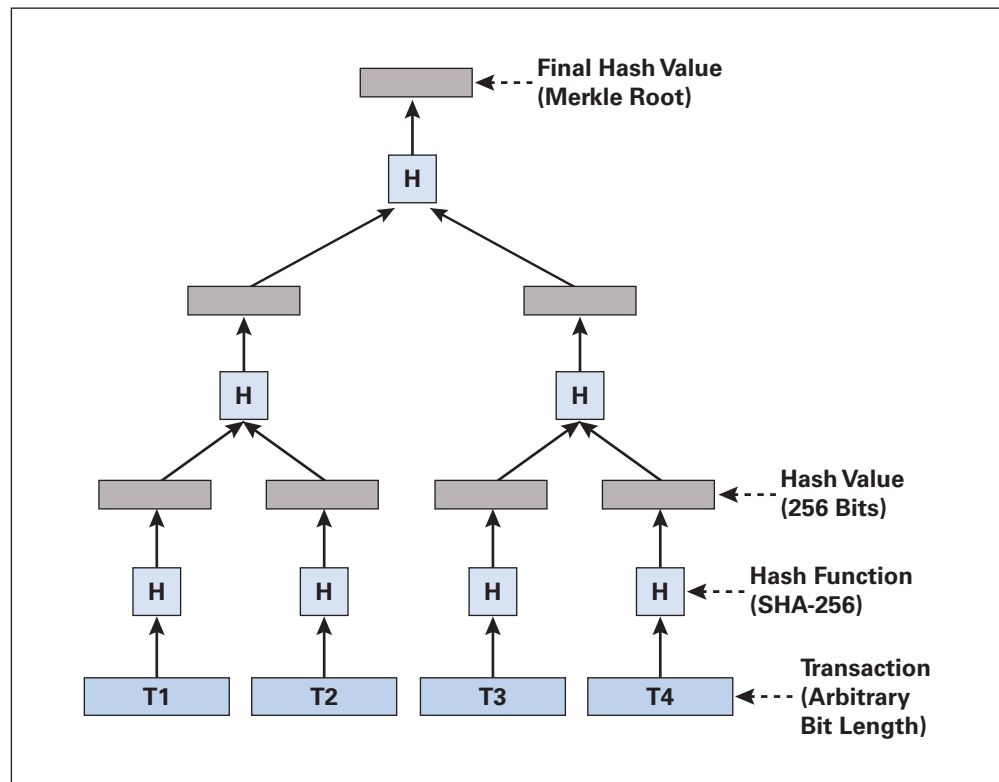
Item	Description
Magic Number	A unique identifier for the blockchain; remains constant for all subsequent blocks
Blocksize	Number of bytes following up to end of block
Version Number	Block format version
Link to Previous Block	Hash of preceding block header
Transaction Hash	The root node of a Merkle Tree, a descendant of all the hashed pairs in the tree. The root node is a 256-bit hash based on all of the transactions in the block.
Timestamp	When block was created
Mining Difficulty	A relative measure of how difficult it is to find a new block. The difficulty is adjusted periodically as a function of how much hashing power has been deployed by the network of miners.
Nonce	Used to calculate proof-of-work
Transaction Counter	Number of transactions in this block
Transactions	The (nonempty) list of transactions

Figure 4: Linkage Between Blocks



The next header field is the *Transaction Hash*. This hash value is computed from the set of data blocks that comprise the list of transactions. Rather than a single hash over this entire set, a *Merkle Tree* technique is used, illustrated in Figure 5. The transaction blocks are processed in pairs; if there is an odd number of transactions, the last transaction on the list is duplicated. Then, each pair of blocks is concatenated to form a binary block that is input to a hash function, typically SHA-256, which produces a 256-bit hash value. The resulting hash values are again paired, and each 512-bit pair is used as input to the hash function. This process continues until a single hash value results, known as the *Merkle Root*.

Figure 5: Example of a Merkle Tree



Following the transaction hash in the header is the *Timestamp* field, which indicates the relative time that this block was created, using a scheme specific to the application.

The next field, *Mining Difficulty*, is a measure of how difficult it is to find a new block. This procedure is explained subsequently. Finally, a one-time *Nonce* value is generated that is used for the proof-of-work concept, described subsequently.

Blockchain Mining

Consider an application that requires the storage of time-sequenced transactions for a distributed group of users. Numerous security issues arise, including authenticating users and ensuring the integrity of the sequence of stored transactions.

The latter includes the need for mechanisms to protect against malicious altering or insertion of transactions. Traditionally, these requirements are met by one or more trusted third parties that act as middleman. In a distributed environment with a large population of users, a peer-to-peer approach becomes more attractive as an efficient method for meeting these requirements. Such an approach is used in blockchain.

The distributed blockchain environment has the following characteristics:^[3]

- Each user has a copy of the blockchain.
- Each user running the blockchain client is part of the network.
- New blocks are broadcast to the network.
- Each user updates its local copy of the blockchain.
- If a user is behind the current height of the chain, it can ask other nodes for copies of the blocks needed to catch up.
- If every user has a copy of the blockchain, when the blockchain is queried, every user gets the same answer.

Within a given application, blocks are created periodically to be added to the chain. The linking of a new block to the end of the chain is most commonly done by a process called *mining*.

Each block in the blockchain is required to have evidence that a costly, nonreversible sacrifice of time and energy has been dedicated to that particular block and no others. This evidence is known as *Proof-of-Work*. The important characteristics of proof-of-work include that it represents a true sacrifice: the actions performed are absolutely useless for any purpose other than producing the proof; and that it is nonreversible: no matter what happens, the resources used to produce the proof cannot be recovered. When Bitcoin clients encounter two valid but different blockchains, they choose to accept the one that represents the highest total proof-of-work.

Some entities within a blockchain network act as miners. It is the task of the miners to add new blocks to the chain and, in effect, miners compete to do this task. Any user can create a set of transactions that are to be formed into a block and added to the end of the chain. The miners are a distributed, pooled resource that create the blocks and add them to the chain.

In a typical open, distributed blockchain application, there are no designated miners. The entity adding the next block to the chain is selected on a per-block basis based on whoever in the world chooses to produce the most proof-of-work. In effect, miners compete for the right to add the next block. The incentive for doing so is a reward based on the application, such as earning Bitcoins for adding to the Bitcoin blockchain.

Miners can enter the system without asking for or requiring anyone's permission, and the network will continue to operate seamlessly when any particular miner leaves the system. The system is kept stable by virtue of the nonrecoverable sacrifice and its ability to discourage non-cooperating miners.

The operation of the miners is governed by a *consensus protocol*. In general, a consensus protocol takes as an input the requests of the components and decides upon one of these requests^[4]. The blockchain consensus protocol ensures that among multiple conflicting proposed transactions, only one gets approved, preventing for example a double spending of the same coins.

A miner constructs a new block in the following fashion: Users broadcast transactions onto the network to be added to a new block. A miner collects these transactions to form a pool of transactions that are not yet part of a block.

Periodically, the miner constructs a new block with the pool of transactions it currently has. The miner validates all the transactions and decides on an ordering within the block. The miner then invests considerable computational effort to construct a new block; this process is called *solving* a block. This block is then broadcast to all the miners on the network and tentatively added to the end of the blockchain. The application requires that each block prove a significant amount of work was invested in its creation to ensure that untrustworthy peers who want to modify past blocks have to work harder than honest peers who only want to add new blocks to the blockchain.

In effect, the consensus mechanism for blockchain is a lottery race, in which the winner is rewarded in some fashion. The winner is the miner that is able to add a new block to the chain that is accepted by other miners.

The technique that is used for the proof-of-work may differ for different applications. Bitcoin uses a *cryptographic hash* technique that works as follows^[5]: The cryptographic hash value of the block header is calculated to form the backward link used by the next block in the chain. If any hash value is allowed, this operation is a simple one. To make the process more resource intensive, a *mining difficulty* is established, which defines how many leading zeros the header hash value must have. Thus, with a mining difficulty of 1, there must be one leading zero. The miner can vary the hash value of the header by varying the value of the nonce field. Typically, a miner will begin with the nonce equal to 1, calculate the hash, and see if it satisfies the difficulty requirement. If not, it increments the nonce and tries again, repeating the process until a hash value is produced that satisfies the difficulty measure.

This difficulty measure is simple to express and effective. For example, if a single leading zero is required, then half of the possible hash values meet the requirement; thus, on average every other hash attempt will result in a hit. If ten leading zeros are required, the level of effort is on the order of one thousand hash attempts. If twenty leading zeros are required, the level of effort is on the order of one million hash attempts. For the Bitcoin blockchain, the target time for solving a block is 10 minutes.

We can express the mining function as follows: For a difficulty level of *alpha*, the hash value *H* must satisfy the following inequality:

$$H(\text{version number, backward link, transaction hash, timestamp, alpha, nonce}) < \text{alpha}$$

The miner must choose a value of nonce that satisfies this inequality. For a secure hash function such as SHA-256, it is effectively impossible to guess a value of nonce that works. Instead, the miner must try out many different values of nonce (using much computing power) until the condition is satisfied. Moreover, the lower the value of *alpha*, the harder it is to satisfy the condition. A proposed solution, however, can easily be verified. That is, once the nonce value is fixed, it is easy to determine if $H(\text{version number, backward link, transaction hash, timestamp, alpha, nonce})$ is less than *alpha*.

A new block can be added to the Bitcoin blockchain only if its header hash is at least as challenging as a difficulty value expected by the consensus protocol. Every 2,016 blocks, the network uses timestamps stored in each block header to calculate the number of seconds elapsed between generation of the first and last of those last 2,016 blocks. The ideal value is 1,209,600 seconds (2 weeks). That is, if blocks are generated at a rate of once per 10 minutes (600 seconds), then 2,016 blocks should be generated in $2,016 \times 600 = 1,209,600$ seconds. If it took fewer than 2 weeks to generate the 2,016 blocks, the expected difficulty value is increased proportionally (by as much as 300%) so that the next 2,016 blocks should take exactly 2 weeks to generate if hashes are checked at the same rate. If it took more than 2 weeks to generate the blocks, the expected difficulty value is decreased proportionally (by as much as 75%) for the same reason.

Returning to a general discussion of blockchain, not specific to Bitcoin, we can now see how the interlocking values of the nonce, transaction hash, and header hash protect the blockchain. An adversary who wishes to successfully alter the transaction list is faced with two alternative challenges: (1) modify the transaction list in such a way that the transaction hash is unchanged, also leaving the header hash unchanged; this modification is computationally infeasible for a secure hash function such as SHA-256; or (2) allow the transaction hash to change but modify the nonce so that the header hash is unchanged; again, this modification is computationally infeasible. Similarly, to insert a new block interior to the chain, the adversary would have to find a nonce value so that the header hash of the inserted block equals that of the preceding block.

Note that the computational effort of the adversary is far greater than that of the miner. Using the Bitcoin difficulty measure, for example, a difficulty value of 30 means that 2^{226} possible hash values can satisfy the requirement and the level of effort is on the order of 2^{30} hash attempts. For an adversary, it is necessary to find a nonce value that will produce a given unique hash value out of the 2^{256} possible values. On average, this discovery will take about 2^{128} hash attempts.

Miner Selection

The *Proof-of-Work* mechanism discussed previously is a way of selecting which miner gets to append a block to the chain. As we have seen, in this scheme the miner is essentially chosen at random through the competition among miners to produce a proof-of-work that is costly to produce but easy to verify.

Other than proof-of-work, numerous alternative methods have been considered or implemented, including the following^[6]:

Proof-of-Stake grants mining rights to participants in proportion to their holding of the currency within the blockchain network. Miners must demonstrate that they hold more than a threshold amount of currency to be able to mine blocks. Proof-of-stake blockchains provide protection from a malicious attack because executing an attack would require the attackers to own a large amount of currency, which is very expensive. Besides, the miners owning a large stake most probably won't attack the system, for example, through double spending. Over time, such attacks will decrease the value of the cryptocurrency and the value of their stake.

The *Proof-of-Burn* process involves destroying Bitcoins by consuming them in a way that does not generate new Bitcoins^[7]. The idea is that miners should show proof that they burned some coins—that is, sent them to a verifiably unspendable address. This process is expensive from their individual point of view, just like proof-of-work; but it consumes no resources other than the burned underlying asset. To date, all proof-of-burn cryptocurrencies work by burning proof-of-work-mined cryptocurrencies, so the ultimate source of scarcity remains the proof-of-work-mined “fuel.”

Permacoin has proposed a modification to Bitcoin^[8], which uses *Proof-of-Retrievability* (POR) to re-purpose the mining resource of Bitcoin to distributed storage of archival data. A POR proves that a node is investing memory or storage resources to store a target file or file fragment. This approach provides additional incentives to contribute resources to the network.

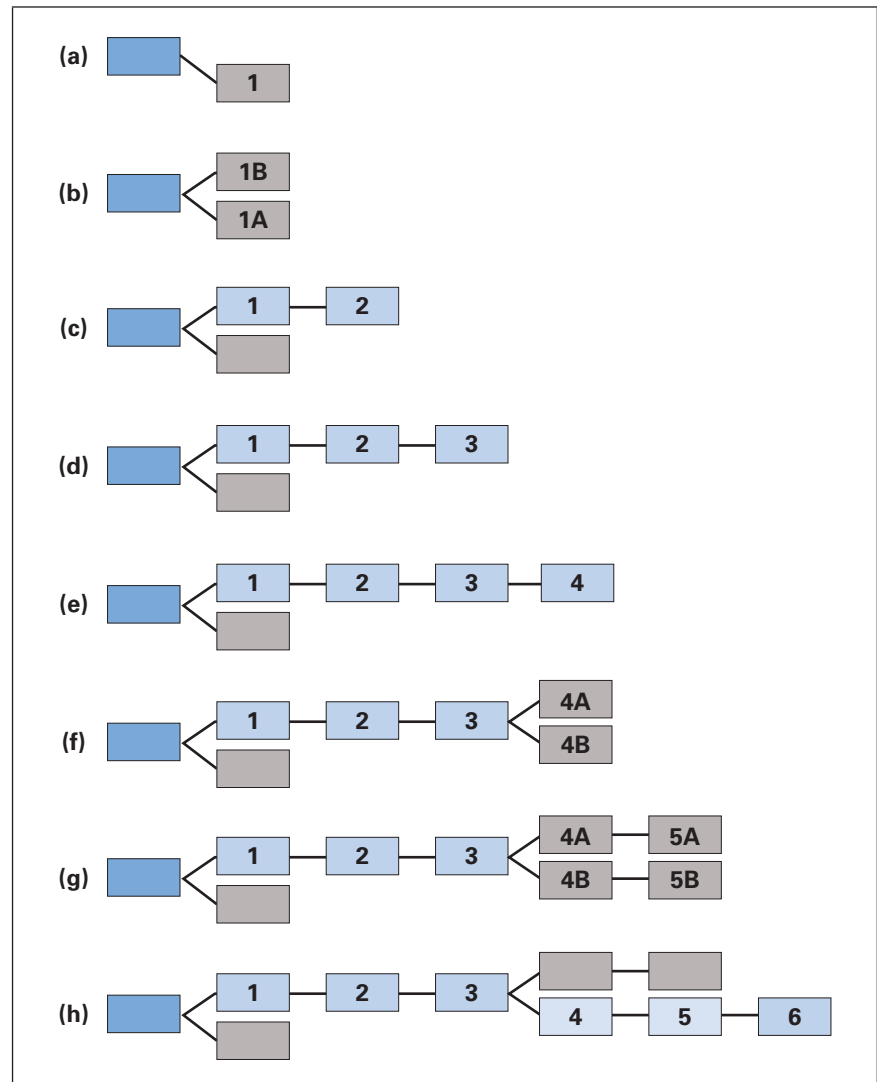
Building the Chain

After a miner has successfully hashed a block header, by finding a nonce value that satisfies the difficulty requirement, the miner can add the block to the end of the chain. It can happen that multiple blocks are created with the same height (distance from the origin block).

This situation occurs if more than one miner, working on different transactions, each produces a block at roughly the same time. This situation creates an apparent fork in the blockchain. Because the blockchain must be a linear, time-ordered sequence, the blocks in the branches following the fork are only provisional. The conflict is resolved when one of the branches exceeds the length of any competing branches.

Figure 6, based on [9], illustrates a typical sequence of events that resolves conflicts. There is an initial block at height 0. All miners try to solve the next block, and one miner solves a block at height 1 (see row *a* in the figure). But at almost the same time, another miner solves a block (row *b*). Block 1A may contain different transactions from 1B, and the users and miners in the network don't know which block should be the accepted one. So, both blocks are considered as provisional, and some miners work on adding to the chain at 1A, and some at 1B.

Figure 6: Adding Blocks to a Chain



Eventually, some miner creates a new block attached to 2B (row *c*). Because all miners must work at the highest height, those miners working on finding a successor to 1A stop that work. All miners are now working on creating a block to attach to the accepted block 2. At row *d* in Figure 6, one miner has successfully created and attached a block at height 3, and broadcasts this update to the network. All of the miners abandon their work at height two and now try to attach a block to the new block 4.

Next, a miner adds a block at height 4 and broadcasts it to the network (row *e*). Other miners, as soon as they receive this information, begin to work at block 4. However, at least one miner, before it receives this update, creates a block at height 4, creating a fork in the chain (row *f*), as happened back at row *c* in the figure. This situation creates a race condition that may continue. It is possible that both forks of the chain solve another block at about the same time (row *g*). In this example, the miners working on 5B solve a block, first adding a new block 6 (row *g*). This new chain, with block 6, is broadcast to all users and miners on the network. With block 6 in place, users are assured the blocks 4A and 5A are “locked in” to the blockchain. And miners who were working to add a block 6 realize they have lost the race. Now all miners begin working to try to append a new block after block 6. This activity continues as the chain grows, with occasional forks that are eventually discarded.

Confirming Transactions

As a miner is collecting transactions, it validates each one. The nature of the validation depends on the application but generally it depends on the use of public-key cryptography to authenticate the parties to the transaction and assure the integrity of the content of the transaction record^[10]. The miner then assembles its current pool of validated transactions into a block. When the block is established as the next block in the chain, it is referred to as a *confirmation*. If there is a fork in the chain, then this confirmation is only provisional until the fork is resolved.

The deeper a block is embedded in a chain (that is, farthest from the current height of the change) the more difficult it would be for an adversary to alter the block. Thus, in any given application, a user of the distributed ledger can decide how many confirmations to wait for before acting with full confidence in a particular ledger entry.

Scalability

A blockchain in active use grows over time and never shrinks, raising the question of the scalability of a blockchain application. For example, Bitcoin allows a maximum block size of 2 MB and as of November 6, 2017, Bitcoin blockchain activity had the following characteristics (<https://blockchain.info>):

- Blockchain size (total size of all block headers and transactions, not including database indexes): 140.295 GB
- Average block size: 1.03 MB

- Transactions per day (most recent day): 333,161
- Aggregate size of transactions waiting to be confirmed: 40.31 MB

To perform all the Bitcoin functions and store the entire blockchain requires considerable processing and memory resources. For many blockchain applications, however, it is not necessary for all users to perform all the blockchain tasks, which include mining management, *Peer-to-Peer* (P2P) network communication and blockchain management, key management, and virtual asset management. For many blockchain applications, systems can be configured that provide only a subset of the tasks of a full implementation, with the handling of public-private key pairs as the most common core feature.

The authors of [11] define five categories of configuration:

- *Basic Client*: A client that runs on a user-controlled device and can perform key management operations, but cannot perform any P2P network communication. It is not a stand-alone solution. Examples include some dedicated hardware clients/wallets and cold-storage (offline storage) clients that require a second online device for transaction processing.
- *Thin Client*: A client that runs on a user-controlled device and can perform key management operations. It performs some P2P tasks related to network communication and blockchain verification but does not keep a copy of the full blockchain.
- *Thick Client*: A client that runs on a user-controlled device and performs all P2P tasks related to network communication and blockchain verification, keeps a copy of the full blockchain, and can perform all key management-related operations. This type of client is sometimes referred to as a *full node*.
- *Fully Functional Basic Client*: A node that performs all of the functions of a thick client, and executes the mining algorithm.
- *Hosted Client*: A client that does not run on a user-controlled device and all tasks are performed by a trusted third party on behalf of the user. This type of client is sometimes referred to as *hosted* or *web client/wallet*. In this case, it is not relevant whether key management is handled in the browser (for example, via JavaScript) because this requirement would in turn require users to download and verify the script code from the website of the third party every time they want to use it.

Depending on the blockchain application, even a full client may not need to store the full chain going back to the genesis block. That task can be reserved for a few archival nodes. The archival nodes can be used to bootstrap fully validating nodes from the beginning but are otherwise not active.

An example of a thin client is the *Simple Payment Verification* (SVP) client used in the Bitcoin application^[12]. The majority of nodes on the Bitcoin network are SVP clients.

An SVP client stores only the portions of a blockchain needed to verify specific transactions of interest to this client. The node downloads the block headers and transactions that represent payments to its addresses. An SPV node doesn't have the security level of a fully validating node. Since the node has block headers, it can check that the blocks were difficult to mine, but it can't check to see that every transaction included in a block is actually valid because it doesn't have the transaction history and doesn't know the set of unspent transactions outputs. SPV nodes can validate only the transactions that actually affect them. The SPV nodes trust the fully validating nodes to have validated all the other transactions that are out there. The cost savings of being an SPV node are substantial. The block headers are only about 0.1% the size of the block chain. So instead of storing tens of gigabytes, the SPV node stores only a few tens of megabytes. Even a smartphone can easily act as an SPV node in the Bitcoin network.

Blockchain Types

Broadly speaking, there are three types of blockchains^[13]. A *Public Blockchain* can be accessed and mined by anyone with Internet access. Access includes not only reading but also posting transactions, and they will be included if they are valid. Nodes participating in a public blockchain network do not have to obtain permission to access the ledger or add transactions. These blockchains are generally considered to be fully decentralized. Public blockchains have the benefit of information transparency and auditability, but they sacrifice information privacy.

A *Consortium Blockchain* is used across multiple organizations. The consensus process is controlled by authorized nodes. For example, one might imagine a consortium of 15 financial institutions, each of which operates a node and of which 10 must sign every block in order for the block to be valid. The right to read the blockchain may be public or restricted to the participants, and there are also hybrid variations such as the root hashes of the blocks being public together with an *Application Programming Interface* (API) that allows members of the public to make a limited number of queries and get back cryptographic proofs of some parts of the blockchain state. These blockchains may be considered partially decentralized.

A *Fully Private Blockchain*, or *Permissioned Blockchain*, limits write permissions to within a single organization that owns the blockchain. Read permissions may be public or restricted to an arbitrary extent. Likely applications include database management and auditing internal to a single company, so public readability may not be necessary in many cases at all, though in other cases public auditability is desired.

Currently, and projected for the foreseeable future, the majority of blockchain applications by market share are public. Most of the remainder are fully private^[14].

Bitcoin

The original application of blockchain, for which it was invented, is *Bitcoin*. Bitcoin is perhaps the most widely used alternative (non-state-issued) currency in the world. Bitcoin is a digital currency scheme^[15]. The network of miners literally creates money out of computer processing cycles. Currency within the system is given value (as it is in any money system) by its scarcity; in this system, the scarcity is created by requiring that money be processed by computationally intensive procedures. Having a great deal of computational power enables a miner to create Bitcoin value more quickly.

Blockchain provides the ledger for recording all of the digital currency transactions. Each transaction is potential only until it is recorded in a block that is accepted as valid and added to the chain. Recording requires the cooperative effort of the miners to achieve. As an incentive, miners are paid in Bitcoins for successfully adding a block to the blockchain.

Other Blockchain Applications

Great interest is being shown in applying blockchain technology to a wide variety of commercial and government applications. The following applications are listed in^[16]:

- *Nasdaq* is using its Linq blockchain technology to complete and record private securities transactions.
- *Depository Trust & Clearing Corporation*, working with market participants and technology firm Axoni, is managing post-trade events for credit default swaps.
- *Factom* is providing blockchain technology for the Honduran land registry project. The focus is data security.
- *Everledger's* focus is on the identity and legitimacy of objects. Blockchain works well here because its history cannot be changed and it enables trust by consensus. The company's initial work provides a distributed ledger of diamond ownership and transaction history verification for owners, insurance companies, claimants, and law enforcement agencies. The system assists with prevention of fraud in the supply chain, but also helps consumers decide whether to buy particular diamonds. The ultimate goal is to track diamonds from mine to market, so that consumers can see if correct duties and taxes have been paid and whether a diamond is a "blood diamond" that has been mined and traded in a war zone and contributed to human atrocity.

Blockchain technology has also caught the interest of numerous government agencies dealing with national security and homeland security, including *Defense Advanced Research Projects Agency* (DARPA) and the U.S. Air Force. The *Department of Homeland Security* (DHS) has awarded contracts for five projects that will use distributed ledger technology to develop new solutions for identity management and privacy protection^[17]:

- *Digital Bazaar* is developing a Linked Data ledger format and architecture to demonstrate how to publish identity credentials.
- *Respect Network Corporation* is developing a decentralized registry and discovery service to integrate with the public blockchain.
- *Narf Industries* is developing an identity management solution built on a permission-less blockchain, with a focus on confidentiality (with selective information disclosure), integrity, availability, non-DHS repudiation, provenance, and pseudo-anonymity.
- *Xcelerate Solutions* is researching blockchain solutions to enable users to establish and maintain trusted identity transactions with public and private organizations.
- *Factom* is studying possible blockchain-based advancements for the security of digital identities for the *Internet of Things* (IoT). The project will create an identity log that captures the identification of a device, who manufactured it, lists of available updates, known security issues, and granted authorities while adding the dimension of time for added security. The goal is to limit would-be hackers' abilities to corrupt the past records for a device, making it more difficult to spoof.

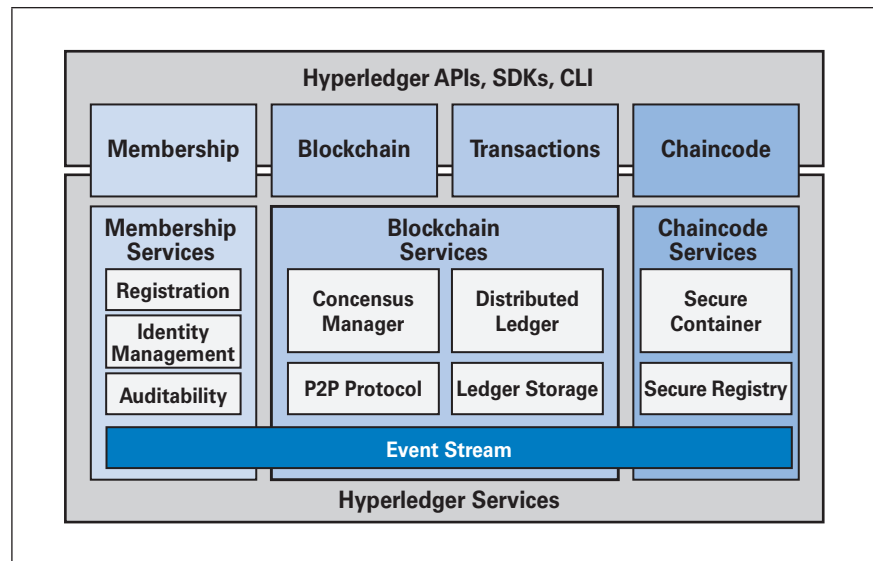
Another interesting development, one that indicates the growing and widespread popularity of the blockchain technology, is the *Initial Coin Offering* (ICO). With the ICO model, instead of selling ownership shares to investors to finance the start of the company, the startup sells digital coins, or tokens, that have value within the application or service the company offers. Sales of ICO tokens exceeded US\$250 million in 2016 and are estimated to exceed US\$1 billion in 2017^[18].

Open-Source Blockchain

In 2016, the Linux Foundation, a nonprofit that champions open-source technologies, announced the *Hyperledger* project, an effort to create an enterprise-grade distributed blockchain ledger framework (<https://www.hyperledger.org>)^[19, 20]. Participants in the group include R3, Cisco Systems, IBM, Intel, and VMware, among others. The objective of this project is to develop a standardized, production-grade digital ledger fabric. The project focuses on identifying and addressing important features for an enterprise-class, cross-industry open standard for distributed ledgers that can transform the way business transactions are conducted globally.

Figure 7 illustrates the current Hyperledger reference architecture within which open-source code is being developed.

Figure 7: Hyperledger Reference Architecture



Four main elements make up a Hyperledger-based application:

- **Membership:** Deals with registering, identifying, and auditing the activity of the peers who will use this particular ledger. The system distinguishes between two kinds of peers. A *validating peer* is a node on the network responsible for running consensus, validating transactions, and maintaining the ledger. A *non-validating peer* is a node that functions as a proxy to connect clients (issuing transactions) to validating peers.
- **Blockchain:** Consists of all the functions associated with building, storing, and providing access to a blockchain ledger.
- **Chaincode:** Implemented in Go, Chaincode is the realization of a smart contract. Each chaincode is encapsulated in a Docker container.
- **Transactions:** Examples of transaction types include the following: A *deploy* transaction takes a chaincode as a parameter; the chaincode is installed on the peers and is ready to be invoked. An *invoke* transaction invokes a transaction of a particular chaincode that has been installed earlier through a deploy transaction; the arguments are specific to the type of transaction. A *query* transaction returns an entry of the state directly from reading the peer's persistent state.

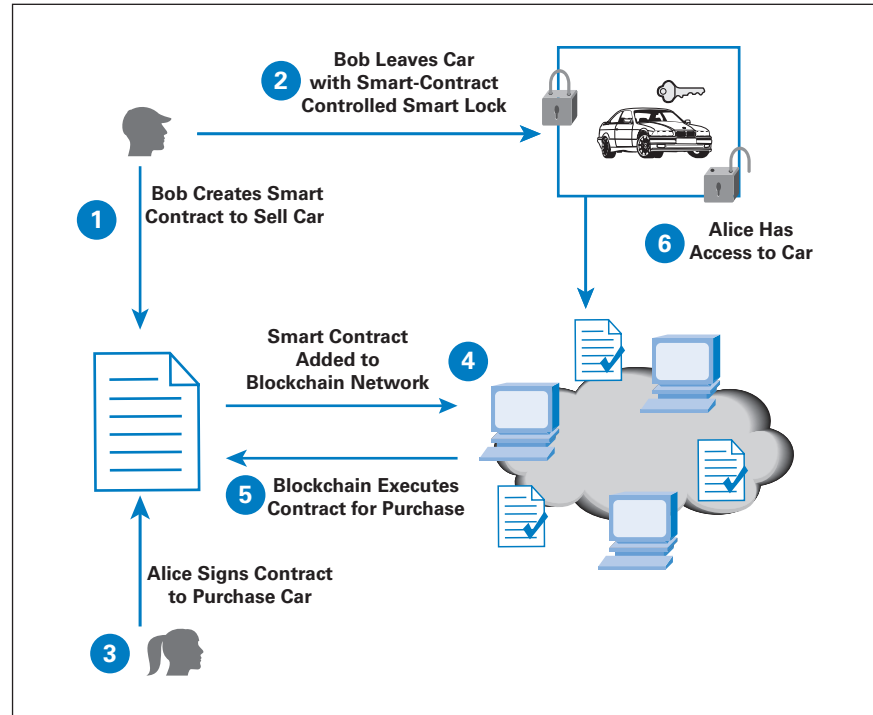
Smart Contracts

Smart Contracts (also called *Self-Executing Contracts*, *Blockchain Contracts*, or *Digital Contracts*) are computer programs that act as agreements, in which the terms of the agreement can be preprogrammed with the ability to self-execute and self-enforce itself. The main goal of a smart contract is to enable two anonymous parties to trade and do business with each other, usually over the Internet, without the need for a middleman.

The concept of smart contracts predates blockchain technology, but it is blockchain that has enabled the sudden growth in the use of smart contracts. One of the most prominent platforms for blockchain-based smart contracts is the open-source platform *Ethereum*^[21].

An example of the use of smart contracts is shown in Figure 8^[22].

Figure 8: Smart Contract Example



The steps involved follow:

1. Bob creates a digital contract to sell his car. He identifies himself with his blockchain address (757382), which is his public key, uses a smart contract to define the terms of the sale, and signs the contract with his private key. The terms might read as follows:

```

IF $20,000 is sent to my account number 757382
THEN Transfer car ID 73849Z to account number that transferred the money
Grant smart lock access to account number that transferred the money
  
```

That English language description corresponds to code embedded in the digital contract. When the contract is added to the blockchain ledger, the code is automatically executed.

2. Bob leaves his car and car key in a garage locked with a smart lock controlled by the smart contract. The car has its own blockchain address 73849Z, which is a public key stored on the blockchain.

3. Alice wants to buy the car and searches for a suitable contract via a web browser. She finds Bob's car listed and signs the contract with her private key, which triggers an automatic transfer of \$20,000 from her blockchain address (389157), which is her public key, to Bob's blockchain address 757382.
4. The signed smart contract is verified by each node in the blockchain network to verify that Bob is the owner of the car and that Alice has sufficient funds for the purchase.
5. If the network verifies the conditions, Alice automatically gets the access code to the smart lock for the garage (encrypted with Alice's public key), Alice is registered as the new owner, and \$20,000 is transferred to Bob.
6. Alice can obtain the access code using her private key and then use the access code to pick up her car.

The whole process is distributed, automated, and does not require a central authority. In general, any blockchain-based smart contract employs the following steps:

1. A contract is defined. For some applications, there is a pre-defined contract specifying the terms of the contract and conditions for execution.
2. An event triggers the execution of the contract. An event could be the initiation of a transaction or the receipt of information.
3. When the contract is added to the blockchain ledger, the contract is executed, a process that typically involves movement of some value based on conditions met.
- 4a. For digital assets on the blockchain, such as cryptocurrency, accounts are automatically settled.
- 4b. For assets that are not part of the blockchain, such as stocks, changes to accounts in the ledger will match settlement instructions off the blockchain.

The smart contract model is very flexible and can be used in a wide variety of applications; some of them are listed in Table 2^[23] on the following page.

Table 2: Blockchain Use Cases

Use Case	Description
Trade Clearing and Settlement	Manages approval workflows between counterparties, calculates trade settlement amounts, and transfers funds automatically
Coupon Payments	Automatically calculates and pays periodic coupon payments and returns principle upon bond expiration
Insurance Claims Processing	Performs error checking, routing, and approval workflows, and calculates payout based on the type of claim and underlying policy
Micro-insurance	Calculates and transfers micropayments based on usage data from an IoT-enabled device (for example, pay-as-you-go automotive insurance)
Electronic Medical Records	Provides transfer and/or access to medical records upon multi-signature approvals between patients and providers
Population Health Data Access	Grants health researchers access to certain health information; micropayments are automatically transferred to the patient for participation
Personal Health Tracking	Tracks patients' health-related actions through IoT devices and automatically generates rewards based on specific milestones
Royalty Distribution	Calculates and distributes royalty payments to artists and other associated parties according to the contract
Autonomous Electric Vehicle Charging Station	Processes a deposit, enables the charging station, and returns remaining funds when complete
Record Keeping	Updates private company share registries and capitalization table records, and distributes shareholder communications
Supply Chain and Trade Finance Documentation	Transfers payments upon multi-signature approval for letters of credit and issues port payments upon custody change for bills of lading
Product Provenance and History	Facilitates chain-of-custody process for products in the supply chain where the party in custody is able to log evidence about the product
Peer-to-Peer Transacting	Matches parties and transfer payments automatically for various peer-to-peer applications: lending, insurance, energy credits, etc.
Voting	Validates voter criteria, logs vote to the blockchain, and initiates specific actions as a result of the majority vote

Summary

Four elements characterize blockchain:

Blockchain is a Replicated Ledger. The ledger provides a history of all transactions that is immutable and is changed only by appending at the newest end of the linear chain of blocks that implements the ledger. The ledger is distributed and readable by all participants.

Blockchain operates by Consensus. Consensus is implemented by means of a shared, decentralized, peer-to-peer protocol. This shared control of the blockchain tolerates disruption, in that from time to time there may be temporary forks in the otherwise linear chain. The consensus mechanism is a means for validating transactions.

Fundamental to the operation of blockchain is *Cryptography*. Various cryptographic algorithms ensure the integrity of the ledger, the authenticity of transactions, the privacy of transactions, and the identity of participants.

Finally, blockchain is a versatile framework for implementing *Business Logic* that is embedded in the ledger. This logic is application-dependent and is reflected in the format and content of the transactions.

References

- [1] “Blockchain: The next big thing,” *The Economist*, May 9, 2015.
- [2] “What CIOs Should Tell the Board of Directors About Blockchain,” Gartner Research, February 14, 2017.
- [3] Giles, A., “Blockchains 101,” Slide presentation from the #StartingBlock2015 tour by @blockstrap, available at: slideshare.net
- [4] Lamport, L. “The Part-Time Parliament,” *ACM Transactions on Computer Systems* (TOCS), May 1998.
- [5] Bitcoin Developer Guide, bitcoin.org/en/developer-guide
- [6] Xu, X. et al., “The Blockchain as a Software Connector,” 2016 13th Working IEEE/IFIP Conference on Software Architecture, April 2016.
- [7] Kaye, M., “How to Secure a Blockchain with Zero Energy,” *Bitcoin Magazine*, January 16, 2014.
- [8] Miller, A. et al., “Permacoin: Repurposing Bitcoin Work for Data Preservation,” IEEE Symposium on Security and Privacy, May 2014.
- [9] Giles, A., “Blockchains 102,” Slide presentation from the #StartingBlock2015 tour by @blockstrap, available at: slideshare.net
- [10] Stallings, W., *Cryptography and Network Security*, Seventh Edition, ISBN-13: 978-0134444284, Pearson, 2017.
- [11] Judmayer, A., Stifter, N., Krombholz, K., and Weippl, E., *Blocks and Chains: Introduction to Bitcoin, Cryptocurrencies, and Their Consensus. Synthesis Lectures on Information Security, Privacy, and Trust*, ISBN-13: 978-1627057165, Morgan & Claypool, 2017.
- [12] Narayanan, A., Bonneau, J., Felton, E., Miller, A., and Goldfeder, S., *Bitcoin and Cryptocurrency Technologies: A Comprehensive Introduction*, ISBN-13: 978-0691171692, Princeton University Press, 2016.

- [13] Buterin, M.V., “On Public and Private Blockchains,” Ethereum Blog, August 7, 2015, <https://blog.ethereum.org/2015/08/07/on-public-and-private-blockchains/>
- [14] Transparency Market Research, “Blockchain Technology Market (Type - Public Blockchain, Private Blockchain, and Consortium Blockchain; Application - Financial Services and Non-financial Sector) - Global Industry Analysis, Size, Share, Growth, Trends, and Forecast 2016–2024,” 2017.
- [15] Velde, F., “Bitcoin: A primer,” Chicago Fed Letter, December 2013, http://www.chicagofed.org/digital_assets/publications/chicago_fed_letter/2013/cfldecember2013_317.pdf
- [16] Underwood, S., “Blockchain Beyond Bitcoin,” *Communications of the ACM*, November 2016.
- [17] Prisco, G., “Department of Homeland Security Awards Blockchain Tech Development Grants for Identity Management and Privacy Protection,” *Bitcoin Magazine*, August 18, 2016.
- [18] Regnier, P., “ICO Is the New IPO,” *Bloomberg BusinessWeek*, June 19, 2017.
- [19] Androulaki, E., “Cryptography and Protocols in Hyperledger Fabric,” Real-World Cryptography Conference, 2017.
- [20] Cachin, C., “Architecture of the Hyperledger Blockchain Fabric,” Workshop on Distributed Cryptocurrencies and Consensus Ledgers, July 2016.
- [21] Bogatyy, I., “A Next-Generation Smart Contract and Decentralized Application Platform,” Ethereum White Paper, <https://github.com/ethereum/wiki/wiki/White-Paper>
- [22] BlockchainHub, “Smart Contracts,” <https://blockchainhub.net/smart-contracts/>
- [23] Odini, M., “Understanding Blockchain,” Slideshare presentation, February 16, 2017, [slideshare.net](https://www.slideshare.net)

WILLIAM STALLINGS is an independent consultant and author of numerous books on security, computer networking, and computer architecture. His latest book is the forthcoming *Effective Cybersecurity: A Practical Guide to Standards and Best Practices* (Pearson, 2018). He maintains a computer science resource site for computer science students and professionals at ComputerScienceStudent.com and is on the editorial board of *Cryptologia*. He has a Ph.D. in computer science from M.I.T. He can be reached at ws@shore.net

In Defence of NATs

by Geoff Huston, APNIC

Network Address Translation (NAT) has often been described as an unfortunate aberration in the evolution of the Internet, and one that will be expunged with the completion of the transition to *Internet Protocol Version 6* (IPv6). I think that this view, which appears to form part of today's conventional wisdom about the Internet, unnecessarily vilifies NATs. In my opinion, NATs are far from being an aberration; instead I see them as an informative and positive step in the evolution of the Internet, particularly as they relate to possibilities in the evolution of name-based networking. Here's why.

Background

It was in 1989, some months after the US National Science Foundation-funded IP backbone network had been commissioned, and at a time when there was a visible momentum behind the adoption of IP as a communications protocol of choice, that the first inklings of the inherent finite nature of the IPv4 address became apparent in the *Internet Engineering Task Force* (IETF)^[1].

Progressive iterations over the IP address consumption numbers reached the same general conclusion: that the momentum of IP deployment meant that the critical parts of the 32-bit address space would be fully committed within 6 or so years. It was predicted that by 1996 we would have fully committed the pool of Class B networks, which encompassed one quarter of the available *Internet Protocol Version 4* (IPv4) address space. At the same time, we were concerned about the pace of growth of the routing system, so stop-gap measures that involved assigning multiple Class C networks to sites could have staved off exhaustion for a while, but perhaps at the expense of the viability of the routing system^[2].

The IETF considered other forms of temporary measures, and the stop-gap measure that was adopted in early 1994 was the dropping of the implicit network/host partitioning of the address in classful addressing in favour of the use of an explicit network mask, or *classless* addressing. This change directly addressed the pressing nature problem of the exhaustion of the Class B address pool, as the observation at the time was that while a Class C network was too small for many sites given the recent introduction of the personal computer, Class B networks were too large, and many sites were unable to realise reasonable levels of address use with Class B addresses. This move to classless addressing (and classless routing, of course) gained some years of breathing space before the major impacts of address exhaustion, and the time gained was considered enough to complete the specification and deployment of a successor IP protocol^[3].

In the search for a successor IP protocol, several ideas were promulgated. The decisions around the design of IPv6 related to a desire to make minimal changes to the IPv4 specification, while changing the size of the address fields and changing some of the encoding of control functions by using the extension header concept, and changing the fragmentation behaviour to stop routers from performing fragmentation in real time^[4].

The common belief at the time was that the adoption of classless addressing in IPv4 bought sufficient time to allow the deployment of IPv6 to proceed. It was anticipated that IPv6 would be deployed across the entire Internet well before the remaining pools of IPv4 addresses were fully committed. This assumption, together with a deliberate approach for hosts to prefer to use IPv6 for communication when both IPv4 and IPv6 was available for use, would imply that the use of IPv4 would naturally dwindle away as more IPv6 was deployed, and that no “flag day” or other means of coordinated action would be needed to complete this Internet-wide protocol transition^[5].

In the flurry of documents that discussed a successor protocol was work that explored the concepts behind “address realms” where one single unique address realm could be replaced by a number of distinct address realms, where the addresses in a packet header could be rewritten when the packet passed across a realm boundary^[6]. One paper at that time described the concept of source address sharing^[7]. If a processing unit was placed on the wire, it was possible to intercept all outbound *Transmission Control Protocol* (TCP) and *User Datagram Protocol* (UDP) packets and replace the source IP address with a different address and change the packet header checksum, and then forward the packet on towards its intended destination. As long as this unit used one of its own addresses as the new address, then any response from the destination would be passed back to this unit. The unit could then use the other fields of the incoming IP packet header, namely the source address and the source and destination port addresses, to match this packet with the previous outgoing packet and perform the reverse address substitution, this time replacing the destination address with the original source address of the corresponding outgoing packet. This scenario allowed multiple internal end systems to use a “public” address, provided that they were not all communicating simultaneously. More generally, a pool of public addresses could be shared across a larger pool of internal systems.

It may not have been the original intent of the inventors of this address-sharing concept, but the approach was enthusiastically adopted by the emerging *Internet Service Provider* (ISP) industry in the 1990s. ISPs were seeing the emergence of the home network and were unprepared to respond to it. With the previous deployment model, using dial-up modems, each active customer was assigned a single IP address as part of the session start process.

A NAT in the gateway to the home network could extend this “single IP address per customer” model to include households with home networks and multiple attached devices. To do so efficiently a further refinement was added, namely that the source port was part of the translation. That way up to 65,535 simultaneous TCP sessions could theoretically share a single external address, provided that the NAT could rewrite the source port along with the source address^[8].

For the ensuing decade NATs were deployed at the edge of the network, and ISPs have used them as a means of externalising the need to conserve IP addresses. The address-sharing technology was essentially deployed by, and operated by, the end customer, and within the ISP network each connected customer still required just a single IP address.

But perhaps that role is underselling the value of NATs in the evolution of the Internet. NATs provided a *firewall* between the end customer and the carrier. The telephony model shared the same end-to-end service philosophy, but it achieved this protection by exercising overarching control over all components of the service. For many decades, the telephone network was a controlled monopoly that was intolerant of any form of competitive interest in the customer. The Internet did not go down this path, and one of the reasons why is that NATs allowed end customers to populate their home network with whatever equipment they chose, and via a NAT, present to the ISP carrier as a single “termination” with a single IP address. This effective segmentation of the network created a parallel segmentation in the market, which allowed the consumer services segment to flourish without carrier-imposed constraint. And at the time that was critically important. The Internet wasn’t the next generation of the telephone service. It was an entirely different utility service operating in an entirely different manner.

More recently, NATs have appeared within the access networks themselves, performing the address-sharing function across a larger set of customers. This function was first associated with mobile access networks but has been used in almost all recent deployments of access networks, as a response to the visible scarcity in the supply of available IPv4 addresses.

NATs have not been universally applauded. Indeed, in many circles within the IETF, NATs were deplored.

It was observed that NATs introduced active middleware into an end-to-end architecture, and divided the pool of attached devices into clients and servers. Clients (behind NATs) had no constant IP address and could not be the target of connection requests. Clients could communicate only with servers, not with each other. It appeared to some to be a step in a regressive direction that imposed a reliance on network middleware with its attendant fragility, and imposed an asymmetry on communication^[9].

For many years, the IETF did not produce standard specifications for the behaviour of NATs, particularly in the case of handling of UDP sessions. Because UDP has no specific session controls, such as session opening and closing signals, how was a NAT meant to maintain its translation state? In the absence of a specific standard specification, different implementations of this function made different assumptions and implemented different behaviour, introducing another detrimental aspect of NATs: *variability*.

How could an application operate through a NAT if the application used UDP? The result was the use of various NAT discovery protocols that attempted to provide the application with some understanding of the particular form of NAT behaviour that it encountered^[10].

NATs in Today's Internet

Let's now look at the situation today in the Internet of late 2017. The major hiatus in the supply of additional IPv4 addresses commenced in 2011 when the central *Internet Assigned Numbers Authority* (IANA) pool of unallocated IPv4 addresses was exhausted. Progressively the *Regional Internet Registries* (RIRs) ran down their general allocation address pools: *Asia Pacific Network Information Centre* (APNIC) in April 2011, *Réseaux IP Européens Network Coordination Centre* (RIPE NCC) in September 2012, *Latin America and Caribbean Network Information Centre* (LACNIC) in 2014, and *American Registry for Internet Numbers* (ARIN) in 2015. The intention from the early 1990s was that the impending threat of imminent exhaustion of further addresses would be the overwhelming impetus to deploy the successor protocol. By that thinking then the Internet would have switched to use IPv6 exclusively before 2011. Yet, that has not happened.

Today a minimum of 90% of the connected device population of the Internet still uses IPv4 exclusively, while the remainder use IPv4 and IPv6^[11]. This network is an all-IPv4 network with a minority proportion also using IPv6. Estimates vary of the device population of today's Internet, but they tend to fall within a band of 15 to 25 billion connected devices^[12]. Yet only some 2.8 billion IPv4 addresses are visible in the Internet routing system. This reality implies that on average each announced public IPv4 address serves from 3 to 8 hidden internal devices.

Part of the reason why estimates of the total population of connected devices are so uncertain is that NATs occlude these internal devices so effectively that no conventional Internet census can expose these hidden internal device pools with any degree of accuracy.

And part of the reason why the level of IPv6 deployment is still so low is that users, and the applications that they value, appear to operate perfectly well in a NATed environment. The costs of NAT deployment are offset by preserving the value of existing investment, both as a tangible investment in equipment and as an investment in knowledge and operational practices in IPv4.

NATS can be deployed incrementally, and they do not rely on some ill-defined measure of coordination with others to operate effectively. They are perhaps one of the best examples of a piecemeal incremental deployment technology where the incremental costs of deployment directly benefit the entity who deployed the technology. This situation is in direct contrast to IPv6 deployment, where the ultimate objective of the deployment, namely the comprehensive replacement of IPv4 in the Internet, can be achieved only after a significant majority of the population of the Internet are operating in a mode that supports both protocols. Until then the deployments of IPv6 are essentially forced to operate in a dual-stack mode, and also support IPv4 connectivity. In other words, the incremental costs of deployment of IPv6 generate incremental benefit only when others also take the same decision to deploy this technology. Viewed from the perspective of an actor in this space, the pressures and costs to stretch the IPv4 address space to encompass an ever-growing Internet are a constant factor. The decision to complement that factor with a deployment of IPv6 means an additional cost that in the short term does not offset any of the IPv4 costs.

So, for many actors the question is not “Should I deploy IPv6 now?” but “How far can I go with NATs?” By squeezing some 25 billion devices into 2 billion active IPv4 addresses, we have used a compression ratio of around 14:1, or the equivalent of adding 4 additional bits of address space. These bits have been effectively “borrowed” from the TCP and UDP port address space. In other words, today’s Internet uses a 36-bit address space in aggregate to allow these 25 billion devices to communicate.

Each additional bit doubles this pool, so the theoretical maximum space of a comprehensively NATed IPv4 environment is 48 bits, fully accounting for the 32-bit address space and the 16-bit port address space. This number is certainly far less than the IPv6 128 bits of address space, but the current division of IPv6 into a 64-bit network prefix and a 64-bit interface identifier drops the available IPv6 address space to 64 bits. The prevalent use of a /48 as a site prefix introduces further address use inefficiencies that effectively drop the IPv6 address space to span the equivalent of some 56 bits.

NATs can be pushed harder. The “binding space” for a NAT is a 5-tuple consisting of the source and destination IP address, a source and destination port address, and a protocol identifier. This 96-bit NAT address space is a highly theoretic ceiling, but the pragmatic question is how much of this space can be exploited cost-effectively such that the marginal cost of exploitation is lower than the cost of an IPv6 deployment.

NATs as Architecture

NATs appear to have pushed applications to a further level of refinement and abstraction that were at one point considered to be desirable objectives rather than onerous limitations.

The maintenance of both a unique fixed-endpoint address space and a uniquely assigned name space for the Internet could be regarded as an expensive luxury when it appears that only one of these spaces is strictly a necessity in terms of ensuring integrity of communication.

The IPv4 architecture made several simplifying assumptions—one of which was that an IPv4 address was overloaded with both the unique identity of an endpoint and its network location. In an age where computers were bolted to the floor of a machine room, this assumption seemed very minor. However, in today's world it appears that the overwhelming number of connected devices are portable devices that change their location constantly, both in a physical sense and in terms of network-based location. This paradigm places stress on the IP architecture, and the result is that IP is variously tunnelled or switched in the final-hop access infrastructure in order to preserve the overloaded semantics of IP addresses.

NATs deliberately disrupt this relationship, and the presented client-side address and port have a particular interpretation and context only for the duration of a session.

In the same way that clients now share IP addresses, services now also share addresses. Applications cannot assume that the association of a name to an IP address is a unique 1:1 relationship. Many service-identifying names may be associated with the same IP address, and in the case of multihomed services the name could be associated with several IP addresses.

With this change comes the observation that IP addresses are no longer the essential “glue” of the Internet. They have changed to a role of ephemeral session tokens that have no lasting semantics. NATs are pushing us to a different network architecture that is far more flexible—a network that uses names as the essential glue that binds it together.

We are now in the phase of the Internet evolution where the address space is no longer unique, and we rely on the name space to offer coherence to the network.

From that perspective, what does IPv6 really offer?

More address bits? Well perhaps not all that much. The space created by NATs operates from within a 96-bit vector of address and port components, and the usable space may well approach the equivalent of a 50-bit conventional address architecture. On the other hand, the IPv6 address architecture has stripped off some 64 bits for an interface identifier and conventionally uses a further 16 bits as a site identifier. The resulting space is of the order of 52 bits. It's not clear that the two pools of address tokens are all that much different in size.

More flexibility? IPv6 is a return to the overloaded semantics of IP addresses as being unique endpoint tokens that provide a connected device with a static location and a static identity. This situation appears to be somewhat ironic in view of the observation that increasingly the Internet is largely composed of battery-powered mobile devices of various forms.

Cheaper? Possibly, in the long term, but not in the short term. Until we get to the “tipping point” that would allow a network to operate solely using IPv6 without any visible impact on the user population of the network, then every network still must provide a service using IPv4.

Permanent address-to-endpoint association? Well, not really. Not since we realised that having a fixed interface identifier represented an unacceptable privacy leak. These days IPv6 clients use so-called *privacy addresses* as their interface identifier, and regularly change this local identifier value.

Perhaps we should appreciate the role of NATs in supporting the name-based connectivity environment that is today’s Internet. It was not a deliberately designed outcome, but a product of incremental evolution that has responded to the various pressures of scarcity and desires for greater flexibility and capability. Rather than eschewing NATs in the architecture as an aberrant deviation in response to a short-term situation, we may want to contemplate an Internet architecture that embraces a higher level of flexibility of addressing. If the name space is truly the binding glue of the Internet, then perhaps we might embrace a view that addresses are simply needed to distinguish one packet flow from another in the network, and nothing more.

Appreciating NATs

When NATs were first introduced to the Internet, they were widely condemned as an aberration in the Internet architecture. And in some ways NATs have directly confronted the model of a stateless packet switching network core and capable attached edge devices.

But that model has been a myth for decades. The Internet as it is deployed is replete with various forms of network “middleware,” and the concept of a simple stateless packet switching network infrastructure has been relegated to the status of an historical, but now somewhat abstract, concept.

In many ways, this condemnation of NATs was unwarranted, as we can reasonably expect that network middleware is here to stay, irrespective of whether the IP packets are formatted as IPv4 or IPv6 and irrespective of whether the outer IP address fields in the packets are translated or not.

Rather than being condemned, perhaps we should appreciate the role that NATs play in the evolution of the architecture of the Internet.

We have been contemplating what it means to have a name-based data network, where instead of using a fixed relationship between names and IP addresses, we eschew this mapping and perform network transactions by specifying the name of the desired service or resource^[13]. NATs are an interesting step in this direction, where IP addresses have lost their fixed association with particular endpoints, and are used more as ephemeral session tokens than endpoint locators. This step certainly appears to be an interesting one in the direction of named data networking.

The conventional wisdom is that the endpoint of this current transitioning Internet is an IPv6 network that has no further use for NATs. But it may not be true. We may find that NATs continue to offer an essential level of indirection and dynamic binding capability in networking that we would rather not casually discard. It may be that NATs are a useful component of network middleware and that they continue to have a role in the Internet well after this transition to IPv6 has been completed, whenever that may be!

References

- [1] Frank Solensky, "Continued Internet Growth," Proceedings of the 18th Internet Engineering Task Force Meeting, August 1990.
- [2] Hans Werner Braun, Peter Ford, and Yakov Rekhter, "CIDR and the Evolution of the Internet," SDSC Report GA-A21364, Proceedings of INET'93, Republished in *ConneXions—The Interoperability Report*, Volume 7, No. 9, September 1993.
- [3] Vince Fuller, Tony Li, Jessica Yu, and Kannan Varadhan, "Classless Inter-Domain Routing (CIDR): An Address Assignment and Aggregation Strategy," RFC 1519, September 1993.
- [4] Scott Bradner and Allison Mankin, "The Recommendation for the IP Next Generation Protocol," RFC 1752, January 1995.
- [5] Dan Wing and Andrew Yourtchenko, "Happy Eyeballs: Success with Dual-Stack Hosts," RFC 6555, April 2012.
- [6] Lixia Zhang, "A Retrospective View of Network Address Translation," *IEEE Network*, September/October 2008.
- [7] Paul Tsuchiya and Tony Eng, "Extending the IP Internet Through Address Reuse," ACM SIGCOMM *Computer Communications Review*, Volume 23, No.1, January 1993.

- [8] Pyda Srisuresh and Der-hwa Gan, “Load Sharing Using IP Network Address Translation (LSNAT),” RFC 2391, August 1998.
- [9] Tony Hain, “Architectural Implications of NAT,” RFC 2993, November 2000.
- [10] Geoff Huston, “Anatomy: A Look Inside Network Address Translators,” *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [11] IPv6 Deployment Measurement,
<https://stats.labs.apnic.net/ipv6>
- [12] Internet of Things connected devices 2015–2025:
<https://www.statista.com/statistics/471264/iot-number-of-connected-devices-worldwide/>
- [13] L. Zhang, et Al., “Named Data Networking,” ACM SIGCOMM *Computer Communication Review*, Volume 44, No. 3, July 2014.
- [14] Daniel Karrenberg, Yakov Rekhter, Elliot Lear, and Geert Jan de Groot, “Address Allocation for Private Internets,” RFC 1918, February 1996.

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

Fragments

Kimberly C. Claffy Honored with Postel Award

The Internet Society, a global non-profit dedicated to ensuring the open development, evolution and use of the Internet, recently announced that Dr. Kimberly C. Claffy, founder and director of the *Center for Applied Internet Data Analysis* (CAIDA) is this year's recipient of the prestigious *Jonathan B. Postel Service Award*.



© Stonehouse
Photographic/Internet Society

Dr. Claffy is a pioneer in the field of measuring and understanding the Internet, not only through her research contributions, but her commitment to establishing and operating infrastructure to support large-scale data collection, curation, and sharing with the scientific research community.

The Postel Award was established by the Internet Society to honor individuals or organizations that, like Jon Postel, have made outstanding contributions to the data communications community. The award is focused on sustained and substantial technical contributions, service to the community, and leadership.

Dr. Claffy was selected by an international award committee comprised of former Postel Award winners. The committee placed particular emphasis on candidates who have supported and enabled others in addition to their own contributions. The committee noted that the award is being presented to Dr. Claffy in recognition for: “her pioneering work on Internet measurement through the development of infrastructure and methodologies for data collection, analysis, and sharing around the world.”

The first of Dr. Claffy’s many papers on Internet traffic measurement and analysis was published in 1992, years before the Internet transitioned to the global, private sector led network it is today. Since then, she has published dozens of papers and received numerous grants and awards for her work.

In 1997 Dr. Claffy founded CAIDA, based at the University of California’s San Diego Super-computer Center, as a center which conducts network research and builds research infrastructure to support large-scale data collection, curation, and data distribution to the scientific research community.

“Simply put, Dr. Claffy’s long-standing and pioneering work has helped the global community better understand the Internet and how it is used,” explained Kathy Brown, President and CEO of the Internet Society, who presented the award.

“In addition to leading the way in the field of Internet measurement and analysis itself, her dedication of resources to ensure widespread access to measurement data has allowed a range of disciplines—from network science and network operations to political science and public policy—to benefit from her efforts.”

KSK Rollover Postponed

The *Internet Corporation for Assigned Names and Numbers* (ICANN) recently announced that the plan to change the cryptographic key that helps protect the *Domain Name System* (DNS) is being postponed. Changing the key involves generating a new cryptographic key pair and distributing the new public component to the *Domain Name System Security Extensions* (DNSSEC)-validating resolvers. Based on the estimated number of Internet users who use DNSSEC validating resolvers, an estimated one-in-four global Internet users, or 750 million people, could be affected by the KSK rollover.

The changing or “rolling” of the *Key Signing Keys* (KSK) was originally scheduled to occur on October 11, 2017, but it is being delayed because some recently obtained data shows that a significant number of resolvers used by *Internet Service Providers* (ISPs) and Network Operators are not yet ready for the rollover. The availability of this new data is due to a very recent DNS protocol feature that adds the ability for a resolver to report back to the root servers which keys it has configured. There may be multiple reasons why operators do not have the new key installed in their systems: some may not have their resolver software properly configured and a recently discovered issue in one widely used resolver program appears to not be automatically updating the key as it should, for reasons that are still being explored.

ICANN is reaching out to its community, including its Security and Stability Advisory Committee, the *Regional Internet Registries*, Network Operator Groups and others to help explore and resolve the issues. In the meantime, ICANN believes it prudent to follow its process and to delay the changing of the key rather than run the risk of a significant number of Internet users being adversely affected. ICANN is committed to continuing its education, communication and engagement with the relevant technical organizations to ensure readiness for the key change.

“The security, stability and resiliency of the domain name system is our core mission. We would rather proceed cautiously and reasonably, than continue with the roll on the announced date of 11 October,” said Göran Marby, ICANN CEO. “It would be irresponsible to proceed with the roll after we have identified these new issues that could adversely affect its success and could adversely affect the ability of a significant number of end users.”

A new date for the Key Roll has not yet been determined. ICANN’s Office of the Chief Technology Officer says it is tentatively hoping to reschedule the Key Roll for the first quarter of 2018, but that it will be dependent on more fully understanding the new information and mitigating as many potential failures as possible. ICANN will provide additional information as it becomes available and the new Key Roll date will be announced as appropriate.

For more information, visit:

<https://www.icann.org/resources/pages/ksk-rollover>

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Fabrizio Accatino	George Ehlers	Aart Jochem
Scott Aitken	Peter Eisses	Jim Johnston
Antonio Cuñat Alario	Torbjörn Eklöv	Jonatan Jonasson
Matteo D'Ambrosio	ERNW GmbH	Daniel Jones
Jens Andersson	ESdatCo	Gary Jones
Danish Ansari	Steve Esquivel	Amar Joshi
David Atkins	Mikhail Evstiounin	Merike Kaeo
Jaime Badua	Paul Ferguson	David Kekar
John Bigrow	Gary Ford	Shan Ali Khan
Axel Boeger	Christopher Forsyth	Nabeel Khatri
Kevin Breit	Craig Fox	Anthony Klopp
Ilia Bromberg	Tomislav Futivic	Henry Kluge
Christophe Brun	Edward Gallagher	Andrew Koch
Gareth Bryan	Andrew Gallo	Carsten Koempe
Stefan Buckmann	Chris Gamboni	Alexander Kogan
Scott Burleigh	Xosé Bravo Garcia	Antonin Kral
Jon Harald Bøvre	Kevin Gee	Mathias Körber
Olivier Cahagne	Serge Van Ginderachter	John Kristoff
Roberto Canonico	Greg Goddard	Terje Krogdahl
John Cavanaugh	Octavio Alfageme Gorostiaga	Bobby Krupczak
Lj Cemeraz	Barry Greene	Warren Kumari
Dave Chapman	Martijn Groenleer	Darrell Lack
Stefanos Charchalakakis	Geert Jan de Groot	Yan Landriault
Greg Chisholm	Gulf Coast Shots	Markus Langenmair
Narelle Clark	Sheryll de Guzman	Fred Langham
Steve Corbató	Martin Hannigan	Richard Lamb
Brian Courtney	John Hardin	Tracy LaQuey Parker
Dave Crocker	Edward Hauser	Simon Leinen
Kevin Croes	David Hauweele	Robert Lewis
John Curran	Headcrafts SRLS	Sergio Loreti
Morgan Davis	Robert Hinden	Guillermo a Loyola
Freek Dijkstra	Edward Hotard	Hannes Lubich
Geert Van Dijk	Bill Huber	Dan Lynch
Richard Dodsworth	Hagen Hultzsich	Miroslav Madi
Ernesto Doelling	Karsten Iwen	Alexis Madriz
Karlheinz Dölger	Ashford Jaggernauth	Carl Malamud
Andrew Dul	David Jaffe	Michael Malik
Holger Durer	Dennis Jennings	Yogesh Mangar
Peter Robert Egli	Edward Jennings	Bill Manning

Harold March
David Martin
Timothy Martin
Gabriel Marroquin
Carles Mateu
Juan Jose Marin Martinez
Ioan Maxim
Miles McCredie
Brian McCullough
Joe McEachern
Carsten Melberg
Kevin Menezes
Bart Jan Menkveld
William Mills
Desiree Miloshevic
Thomas Mino
Mohammad Moghaddas
Charles Monson
Andrea Montefusco
Fernando Montenegro
Soenke Mumm
Tariq Mustafa
Stuart Nadin
Mazdak Rajabi Nasab
Krishna Natarajan
Darryl Newman
Marijana Novakovic
Ovidiu Obersterescu
Mike O'Connor
Carlos Astor Araujo Palmeira
Alexis Panagopoulos
Manuel Uruena Pascual
Ricardo Patara
Dipesh Patel
Alex Parkinson
Craig Partridge
Dan Paynter
Leif-Eric Pedersen
Juan Pena

Chris Perkins
Derrell Piper
Rob Pirnie
Jorge Ivan Pincay Ponce
Blahoslav Popela
Tim Pozar
David Raistrick
Priyan R Rajeevan
Paul Rathbone
Bill Reid
Rodrigo Ribeiro
Justin Richards
Mark Risinger
Ron Rockrohr
Carlos Rodrigues
Lex Van Roon
William Ross
Boudhayan Roychowdhury
Carlos Rubio
RustedMusic
Babak Saberi
George Sadowsky
Scott Sandefur
Sachin Sapkal
Arturas Satkovskis
Phil Scarr
Jeroen Van Ingen Schenau
Carsten Scherb
Roger Schwartz
SeenThere
Scott Seifel
Yury Shefer
Yaron Sheffer
Tj Shumway
Jeffrey Sicuranza
Thorsten Sideboard
Henry Sinnreich
Geoff Sisson
Helge Skrivervik

Darren Sleeth
Bob Smith
Mark Smith
Job Snijders
Ignacio Soto Campos
Peter Spekrijse
Thayumanavan Sridhar
Matthew Stenberg
Adrian Stevens
Clinton Stevens
Viktor Sudakov
Edward-W. Suor
Vincent Surillo
Roman Tarasov
David Theese
Sandro Tumini
Phil Tweedie
Steve Ulrich
Unitek Engineering AG
John Urbanek
Martin Urwaleck
Betsy Vanderpool
Surendran Vangadasalam
Alejandro Vennera
Luca Ventura
Tom Vest
Dario Vitali
Randy Watts
Andrew Webster
Tim Weil
Jd Wegner
Rick Wesson
Peter Whimp
Jurrien Wijlhuizen
Pindar Wong
Bernd Zeimetz

Follow us on Twitter and Facebook



@protocoljournal



<https://www.facebook.com/newipj>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Diamond Sponsors



Ruby Sponsor



Sapphire Sponsors

Your logo here!

Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

ADDRESS SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

David Conrad, Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

April 2018

Volume 21, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

F R O M T H E E D I T O R

In This Issue

From the Editor	1
Considerations in Network Complexity.....	2
IPv6 Fragmentation	13
Letters to the Editor.....	24
Thank You.....	28
Call for Papers.....	30
Supporters and Sponsors	31

We live in a complex and increasingly interconnected world. With this complexity comes a desire by network engineers to design systems that can cope with increasing demands while still offering predictable performance, manageability, and maintainability. In our first article, Russ White discusses ways to analyze network designs from a complexity point of view.

The *Internet Protocol* (IP) was designed to operate over a variety of underlying network technologies, such as Ethernet, X.25, FDDI, Frame Relay, WiFi, and even mobile telephone networks. Applications that use IP must deal with the fact that datagrams may be split into *fragments* as they travel across the network with subsequent *reassembly* at the receiving end. Previous articles in this journal have discussed fragmentation, largely in the context of IPv4. This time Geoff Huston describes fragmentation in IPv6 and the particular challenges that arise with this protocol in conjunction with applications such as the *Domain Name System* (DNS).

We usually provide a section of announcements entitled “Fragments,” but this time it has been replaced by a selection of Letters to the Editor—all in response to articles in our November 2017 issue. We are very happy to receive feedback on any aspect of this journal, and we would also point you to our website, which contains additional articles and material as well as all of our back issues in PDF format.

As mentioned in previous issues, if you have a print subscription to this journal, you will find an expiration date printed on the back cover. For the last couple of years, we have “auto-renewed” your subscription, but now we ask you to log in to our subscription system and perform this simple task yourself. The subscription portal is here: <https://www.ipjsubscription.org/> This process will ensure that we have your current contact information as well as delivery preference (print edition or download). For any questions, contact us by e-mail at: ipj@protocoljournal.org

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

Considerations in Network Complexity

by Russ White

Computer networks are complex—and getting more complex by the day. At one time, knowing the *Internet Protocol* (IP) was enough; today there are underlays, overlays, virtualized services, service chains, and a host of other technologies engineers need to plan around and for. With complexity on the rise, maybe it's time to ask some fundamental questions, such as—what does complexity mean? Can complexity be solved? How can engineers manage complexity?

Why So Complex?

While the most obvious place to begin might be with a definition of complexity, it's actually more useful to consider why complexity is required in a more general sense. To put it more succinctly, is it possible to “solve” complexity? Why not just design networks and protocols that are simpler? Why does every attempt to make anything simpler in the networking world end up apparently making things more complex in the long run? For instance, tunneling on top of (or through) IP reduces the complexity of the control plane and makes the network simpler overall. Why is it, then, that tunneled overlays end up containing so much complexity?

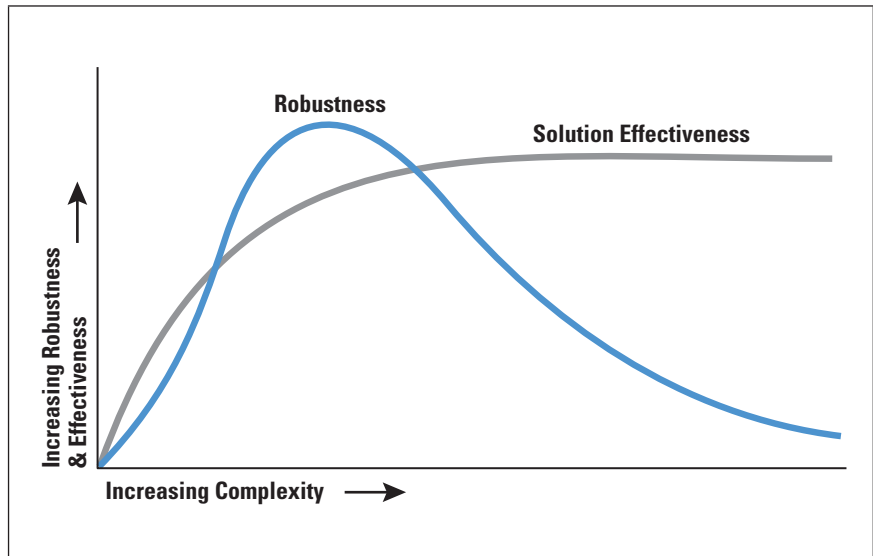
This question has two answers: The first is that human nature being what it is, engineers will always invent 10 different ways to solve the same problem. This reality is especially true in the virtual world, where new solutions are (relatively) easy to deploy, it's (relatively) easy to find a problem with the last set of proposed solutions, and it's (relatively) easy to move some bits around to create a new solution that is “better than the old one.” The virtual space, in other words, is partially so messy because it's so easy to build something new there.

- Abstract the complexity away, to build a black box around each part of the system, so each piece and the interactions among these pieces are more immediately understandable.
- Toss the complexity over the cubicle wall—to move the problem out of the networking realm into the realm of applications, or coding, or a protocol. As RFC 1925^[1] says, “It is easier to move a problem around (for example, by moving the problem to a different part of the overall network architecture) than it is to solve it.”
- Add another layer on top, to treat all the complexity as a black box by putting another protocol or tunnel on top of what's already there. Returning to RFC 1925, “It is always possible to add another level of indirection.”
- Become overwhelmed with the complexity, label what exists as “legacy,” and chase some new shiny thing that will solve all the problems in what is perceived as a much less complex way.

- Ignore the problem and hope it will go away. Argue for an exception “just this once,” to meet a particular business goal, or fix some problem, within a very tight schedule, with the promise that the complexity issue will be dealt with “later,” is a good example.

The second answer, however, lies in a more fundamental problem: complexity is necessary to deal with the uncertainty involved in problems that are difficult to solve (Figure 1).

Figure 1: Complexity, Effectiveness, and Robustness



Adding complexity, then, allows a network to handle future requirements and unexpected events more easily, as well as providing more services over a smaller set of base functions. If this condition is the case, why not simply build a single protocol running on a single network that can handle all the requirements potentially thrown at it, and can handle any sequence of events you can imagine? A single network running a single protocol would certainly reduce the number of moving parts network engineers need to deal with, making all our lives simpler, right?

Maybe not. At some point, any complex system becomes brittle—*robust yet fragile* is one phrase you can use to describe this condition. A system is robust yet fragile when it is able to react resiliently to an expected set of circumstances, but an unexpected set of circumstances will cause it to fail. As an example from the real world—knife blades are required to have a somewhat unique combination of characteristics. They must be hard enough to hold an edge and cut, and yet flexible enough to bend slightly in use, returning to their original shape without any evidence of damage, and they must not shatter when dropped. It has taken years of research and experience to find the right metal to make a knife blade, and there are still long and deeply technical discussions about which material is right for specific properties, under what conditions, etc.

Complexity is necessary, then: it cannot be “solved.”

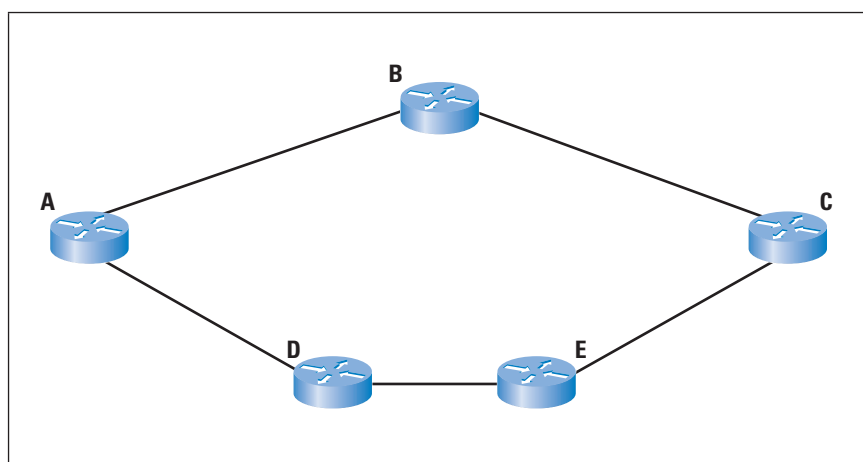
Defining Complexity

Given complexity is necessary, engineers are going to need to learn to manage it in some way, by finding or building a model or framework. The best place to begin to build such a model is with the most fundamental question: what does complexity mean in terms of networks? Can you put a network on a scale and have the needle point to “complex?” Is there a mathematical model into which you can plug the configurations and topology of a set of network devices that will, in turn, produce a “complexity index?” How do the concepts of scale, resilience, brittleness, and elegance relate to complexity? The best place to begin in building a model is with an example.

Control-Plane State versus Stretch

What is network *stretch*? In the simplest terms possible, it is the difference between the shortest path in a network and the path traffic between two points actually takes. Figure 2 illustrates this concept.

Figure 2: A Small Network to Illustrate State and Stretch



Assuming the cost of each link in this network is the same, the shortest physical path between Routers A and C will also be the shortest logical path: [A,B,C]. What happens, however, if we change the metric on the [A,B] link to 3? The shortest physical path is still [A,B,C], but the shortest logical path is now [A,D,E,C]. The differential between the shortest physical path and the shortest logical path is the distance a packet being forwarded between Routers A and C must travel—in this case, the stretch can be calculated as $(4 [A,D,E,C]) - (3 [A,B,C])$, for a stretch of 1.

How Is Stretch Measured?

In terms of hop count, is stretch measured by the summary of the metrics, the delay through the network, or some other way? It depends on what is most important in any given situation, but the most common way is by comparing hop counts through the network, and this method is used in the examples here for simplicity. In some cases, it might be more important to consider the metric along two paths, the delay along two paths, or some other metric, but the important point is to measure it consistently across every possible path to allow for accurate comparison between paths.

It's sometimes difficult to differentiate between the physical topology and the logical topology. In this case, was the [A,B] link metric increased because the link is actually a slower link? If so, whether this is an example of stretch or an example of simply bringing the logical topology in line with the physical topology is debatable.

In line with this observation, it's much easier to define policy in terms of stretch than almost any other way. Policy is any configuration that increases the stretch of a network. Using *Policy-Based Routing* or *Traffic Engineering* to push traffic off the shortest physical path and onto a longer logical path to reduce congestion on specific links, for instance, is a policy—it increases stretch.

Increasing stretch is not always a bad thing. Understanding the concept of stretch simply helps us understand various other concepts, and put a framework around complexity tradeoffs. The shortest path, physically speaking, isn't always the best path.

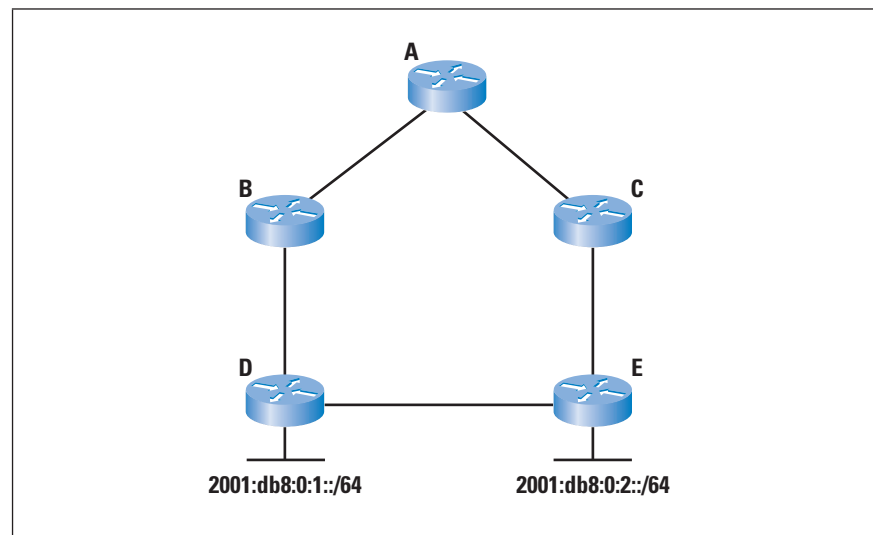
Stretch, in this illustration, is very simple—it affects every destination, and every packet flowing through the network. In the real world, things aren't so simple. Stretch is actually per source/destination pair, making it very difficult to measure on a networkwide basis.

With all of this information in mind, let's look at two specific examples of the tradeoff between stretch and optimization.

Aggregation versus Stretch

Aggregation is a technique used to reduce not only the amount of information carried in the control plane, but also the rate of state change in the control plane. Aggregation is built into IP (both IPv4 and IPv6)—even a single subnet contains multiple host addresses. By connecting a single broadcast segment to a set of hosts, the IP routing protocol doesn't need to manage Layer 2 reachability, nor individual host addresses. Aggregation within the control plane can also cause stretch, as Figure 3 shows.

Figure 3: Aggregation and Stretch



Two different situations illustrate increasing stretch through route aggregation:

1. Assume the [A,B] link has a cost of 2, and all the other links in this network have a cost of 1. If Routers B and C both aggregate to `2001:db8::/61`, then the path through [A,C] would be preferred for everything within the aggregate. Traffic destined to `2001:db8:0:1::/64` will pass along the path [A,C,E,D] to reach its destination, even though the shortest (physical) path is [A,B,D]. The stretch for `2001:db8:0:2::/64` isn't changed, but the stretch for `2001:db8:0:1::/64` is increased by 1.
2. Assume all the links in the network have a cost of 1. If Routers B and C both aggregate to `2001:db8::/61`, then Router A will somehow load share traffic toward the two subnets behind Routers D and E across the two equal-cost paths it has available. Given perfect load sharing, 50% of the traffic destined to `2001:db8:0:1::/64` will flow along [A,C,E,D], with a stretch of 1, and 50% of the traffic destined to `2001:db8:0:2::/64` will flow along [A,B,D,E], with a stretch of 1.

Implementing aggregation removes specific reachability information about the two /64 prefixes behind Routers D and E from the state of Router A. Implementing aggregation also disconnects the state of the individual /64's behind Routers D and E from the state at Router A. Aggregation, then, decreases complexity from the perspective of Router A by reducing the *amount* and *speed* of state in the routing table of Router A.

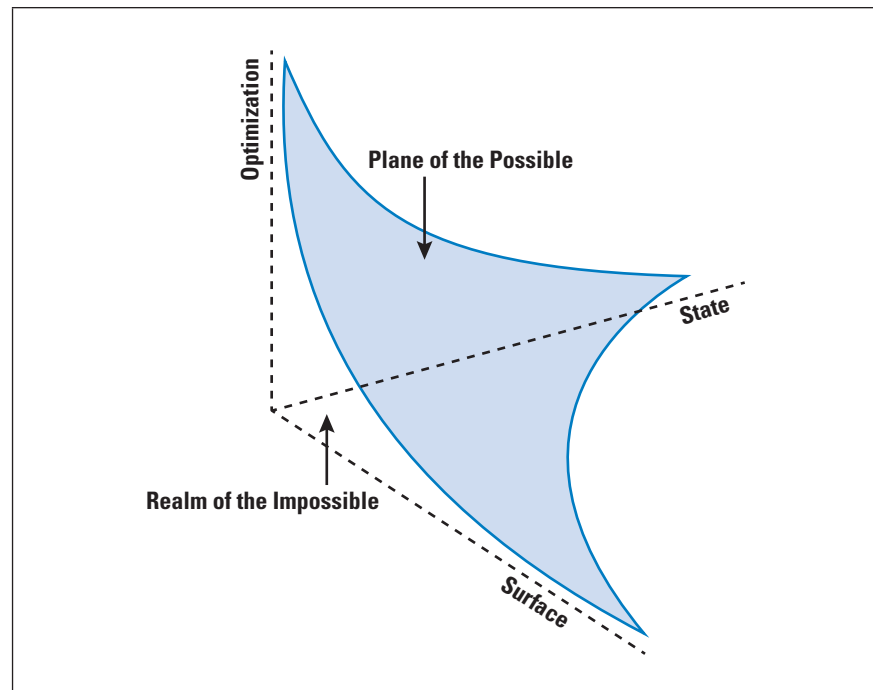
Increasing stretch increases the overall use of the network without any actual increase in the amount of traffic being carried through the network. In the example given in Figure 2, traffic that would normally take a two-hop path is directed along a three-hop path, meaning one more link and router are involved in forwarding and switching the packets in the flow(s) across the network. In purely mathematical terms, increasing stretch decreases the overall efficiency of the network by increasing the number of devices and links used to forward any particular flow.

Finally, to implement aggregation Routers B and C must be configured to summarize the two longer prefixes into a single shorter one. This additional configuration introduces an additional bit of interaction between the human operator (or at least the configuration system) and the control plane. This situation can be described as an increase in *surface* in the network.

Defining Complexity: A Model

These three components—*state*, *optimization*, and *surface*—are common in virtually every network or protocol design decision. They can be seen as a set of tradeoffs, as illustrated in Figure 4.

Figure 4: The Plane of the Possible



Increasing optimization always moves towards more state or more interaction surfaces. Decreasing state always moves towards less optimization or more interaction surfaces. Decreasing interaction surfaces always moves towards less optimization or more state. These rules aren't ironclad, of course; they are contingent on the specific network, protocols, and requirements, but they are generally true often enough to make this model useful for understanding tradeoffs in complexity.

Interaction Surfaces

While state and optimization are fairly intuitive, it's worthwhile to spend just a moment more on interaction surfaces. The concept of interaction surfaces is difficult to grasp primarily because it covers such a wide array of ideas. Perhaps an example would be helpful; assume a function that:

- Accepts two numbers as input
- Adds them
- Multiplies the resulting sum by 100
- Returns the result

This single function can be considered a subsystem in some larger system. Now assume you break this single function into two functions, one of which does the addition, and the other of which does the multiplication. You've created two simpler functions (each one does only one thing), but you've created an interaction surface between the two functions—you've created two interacting subsystems within the system where there used to be only one. This example is really simple, I know, but consider a few more that might help.

The routing information carried in *Open Shortest Path First* (OSPF) is split into *external* routes being carried in *Border Gateway Protocol* (BGP) and *internal* routes being carried in OSPF. You've gone from one system with more state to two systems with less state, but you've created an interaction surface between the two protocols—they must now work together to build a complete forwarding table.

A single set of hosts with different access policies is split onto multiple virtual topologies on the same physical network. You've simplified the amount of state in filtering, but you've created an interaction surface between the two virtual topologies and between the two topologies and the control plane. In addition, you've exposed new shared risk groups where a single physical failure can cause multiple logical ones. Hence you've traded state in one control plane for interaction surfaces between multiple control planes.

Even two routers communicating within a single control plane can be considered an interaction surface. This breadth of definition is what makes it so very difficult to define what an interaction surface is.

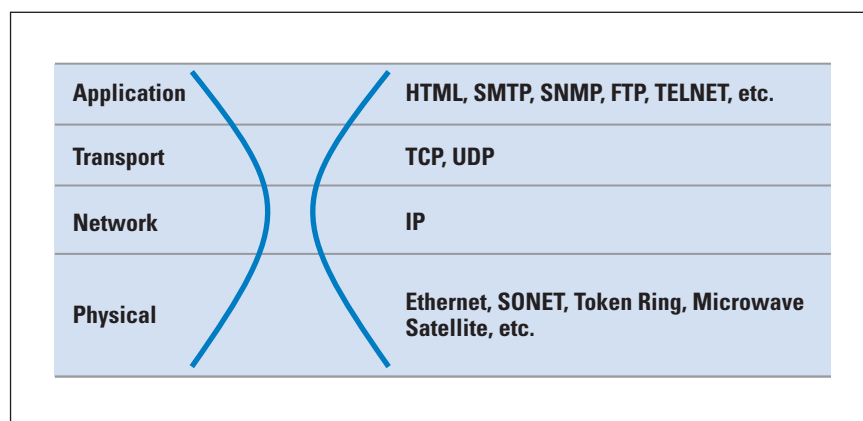
Interaction surfaces aren't a bad thing; they help engineers and designers divide and conquer in any given problem space, from modeling to implementation. At the same time, interaction surfaces are all too easy to introduce without thought.

Managing Complexity through the Wasp Waist

There is a simple model that is ubiquitous throughout the natural world, and is widely mimicked in the engineering world. While engineers don't often consciously apply this model, it's actually used all the time. What is this model?

Figure 5 illustrates the hourglass model in the context of the four-layer *Department of Defense* (DoD) model that gave rise to the *Internet Protocol Suite*.

Figure 5: The DoD Model and the "Wasp Waist"



At the bottom layer, the physical transport system, there is a wide array of protocols, from Ethernet to Satellite. At the top layer, where information is marshaled and presented to applications, there is a wide array of protocols, from *Hypertext Transfer Protocol* (HTTP) to TELNET (and thousands of others besides). A funny thing happens when you move towards the middle of the stack, however: the number of protocols decreases, creating an hourglass. Why does this work to control complexity?

Going back through the three components of complexity—*state*, *optimization*, and *surface*—exposes the relationship between the hourglass and complexity.

- *State* is divided by the hourglass into two distinct types of state: information about the network, and information about the data being transported across the network. While the upper layers are concerned with marshaling and presenting information in a usable way, the lower layers are concerned with discovering what connectivity exists and what the properties of that connectivity actually are. The lower layers don't need to know how to format a *File Transfer Protocol* (FTP) frame, and the upper layers don't need to know how to carry a packet over Ethernet—state is reduced at both ends of the model.
- *Optimization* is traded off by allowing one layer to reach into another layer, and by hiding the state of the network from the applications. For instance, the *Transmission Control Protocol*, (TCP) doesn't really know the state of the network other than what it can gather from local information. TCP could potentially be much more efficient in its use of network resources, but only at the cost of a layer violation, which opens up interaction surfaces that are difficult to control.
- *Surfaces* are controlled by reducing the number of interaction points between the various components to precisely one—IP. This single interaction point can be well defined through a standards process, with changes in the one interaction point closely regulated to prevent massive rapid changes that will reflect up and down the protocol stack.

The layering of a stacked network model is, then, a direct attempt to control the complexity of the various interacting components of a network.

Managing Complexity as an Engineer

Managing complexity in design in a general sense is just one application of the state/optimization/surface model. Another use is in learning how to understand complex systems quickly—a skill every engineer could use in everyday life. Here this three-sided model is used as part of a process, or a way of thinking about systems.

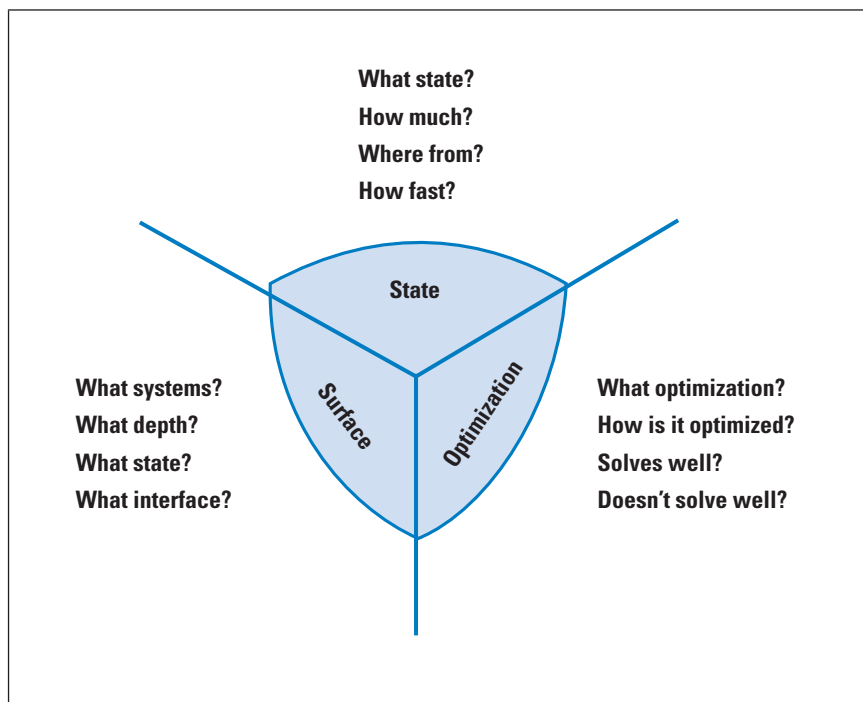
There are three steps in this model, or rather three questions:

- Why is this being done this way?
- What is being accomplished?
- What is this like?

Why acts as an abstractor, focusing on the problem at hand by excluding others. This question uncovers the purpose, or the goal, of the system. Asking why also connects the system to the business driver. As an example: “Why does OSPF elect a designated router?” might be a perfectly valid question in some settings, but not necessarily while you’re troubleshooting slow network convergence. Asking why can focus energy and uncover motives that drive configuration, policy, and other choices.

What is the question engineers normally engage with first, but they often do so with little structure. The three-pronged complexity model provides a solid model with which to find the right questions to ask. Figure 6 illustrates this method.

Figure 6: A Structure for Asking What



By focusing on questions about each of the three prongs of the complexity model, you can quickly uncover the tradeoffs made in system design—even if those tradeoffs weren’t made intentionally.

Finally, *what is this like* helps relate the problem space, the potential solutions, and even potential problems with potential solutions to the system being considered. Assume a new working group is formed to solve a particular problem in the routing space. The working group quickly decides on using Dijkstra's *Shortest Path First* (SPF) algorithm to find the solution to one particular problem the protocol is set out to solve.

Asking *what is this like* should immediately uncover the relationship of new solution to existing link-state protocols, such as OSPF. From there, given experience with OSPF, the engineer knows what sort of convergence characteristics the newly proposed solution is likely to have, and where to look for potential problems. As link-state protocols are subject to microloops in some situations, so the newly proposed solution is likely to be subject to microloops. As link-state protocols can suffer from overwhelming amounts of state during large-scale convergence events, so the new solution might suffer from the same sort of problem.

Concluding Thoughts

You can't run, and you can't hide from complexity; there's no point in even trying. You're going to encounter it; ignoring it doesn't make the problem go away, it just allows the problem to fester under some "rug" in some corner of your network. The complexity problems you create today will return as bigger, more complex problems in just a few years. To quote someone who's spent years looking at complexity:

"Trying to make a network proof against predictable problems tends to make it fragile in dealing with unpredictable problems (through an ossification effect as you mentioned). Giving the same network the strongest possible ability to defend itself against unpredictable problems, it necessarily follows, means that it MUST NOT be too terribly robust against predictable problems—not being too robust against predictable problems is necessary to avoid the ossification issue, but not necessarily sufficient to provide for a robust ability to handle unpredictable network problems."

—Tony Przygienda

That call at 2 a.m. might not be pleasant, but solving it the wrong way might cause a much worse call at 2 a.m. sometime later. Hardening the network against all failures eventually means to make it fail spectacularly when a failure that you didn't predict occurs—there's just no way around this reality.

When dealing with engineering problems, then, a little humility around what can, and cannot, be solved is in order. Don't ignore complexity, but don't think you can solve it, either. Instead, remember to treat every situation as a set of tradeoffs—and if you don't see the tradeoffs, you're not looking hard enough.

References

- [1] Ross Callon, “The Twelve Networking Truths,” RFC 1925, April 1996.
- [2] Russ White and Shawn Zandi, “Cloudy-Eyed: Complexity and Reality with Software-Defined Networks,” *The Internet Protocol Journal*, Volume 19, No. 3, September 2016.
- [3] Russ White and Jeff Tantsura, *Navigating Network Complexity: Next-Generation Routing with SDN, Service Virtualization, and Service Chaining*, Addison-Wesley Professional, 2015, ISBN-13: 978-0133989359.
- [4] Russ White and Ethan Banks, *Computer Networking Problems and Solutions: An Innovative Approach to Building Resilient, Modern Networks*, Addison-Wesley Professional, 2018, ISBN-13: 978-1587145049.

RUSS WHITE began working with computers in the mid-1980s, and computer networks in 1990. He has experience in designing, deploying, breaking, and troubleshooting large-scale networks, and is a strong communicator from the white board to the board room. Across that time, he has co-authored more than 40 software patents, participated in the development of several Internet standards, helped develop the *Cisco Certified Design Expert* (CCDE) and the *Cisco Certified Architect* (CCAr) programs, and worked in Internet governance with the Internet Society. Russ has a background covering a broad spectrum of topics, including radio frequency engineering and graphic design, and is an active student of philosophy and culture. Russ is a co-host at the *Network Collective*, serves on the Routing Area Directorate at the IETF, co-chairs the BABEL working group, serves on the *Technical Services Council* as a maintainer on the open source *FR Routing* project, and serves on the *Linux Foundation* (Networking) board. His most recent works are *Computer Networking Problems and Solutions*, *The Art of Network Architecture*, *Navigating Network Complexity*, and the Intermediate System-to-Intermediate System LiveLesson. He holds a *Master of Science in Information Technology* (MSIT) from Capella University, and a *Master of Arts in Christian Ministries* (MACM) from Shepherds Theological Seminary and is currently working on a PhD at Southeastern Baptist Theological Seminary. E-mail: russ@riw.us

IPv6 and Packet Fragmentation

by Geoff Huston, APNIC

Version 6 of the Internet Protocol (IPv6) introduced very few changes to its Version 4 predecessor (IPv4). The major change was of course the expansion of the size of the IP source and destination address fields in the IP packet header from 32 bits to 128 bits. Some other changes, however, apparently were intended to subtly alter IP behaviour. One of them was the change in treatment of IP *packet fragmentation*.

It appears that rather than effecting a slight improvement from IPv4, the manner of fragmentation handling in IPv6 appears to be significantly more difficult. In light of these problems, it is perhaps unsurprising that calls have been made from time to time to dispense completely with packet fragmentation in IPv6^[1], as the current situation with IPv6 appears to be worse than both no fragmentation and the IPv4-style of fragmentation.

Packet Fragmentation

One of the more difficult design exercises in packet-switched network architectures is the design of packet fragmentation.

In time-switched networks, developed to support a common bearer model for telephony, each “unit” of information passed through the network occurred within a fixed timeframe (an analogue voice stream was digitized into 8,000 sample points per second, so the basic time unit for switching these digital samples was 1/8,000 of a second), which resulted in fixed-size packets, all clocked off a common time base.

Packet-switched data networks could dispense with a constant common time base, in turn allowing individual data packets to be sized according to the needs of the application as well as the needs and limitations of the network substrate.

For example, smaller packets have a higher packet header-to-packet payload ratio and are consequently less efficient in data carriage and impose a higher processing load as a function of effective data throughput. On the other hand, within a packet-switching system the smaller packet can be dispatched faster, reducing head-of-line blocking in the internal queues within a packet switch and potentially reducing network-imposed jitter as a result. This reduction can make it easier to use the network for real-time applications such as voice or video. Larger packets allow larger data payloads, in turn allowing greater carriage efficiency. Larger payloads per packet also allows a higher internal switch capacity when measured in terms of data throughput, in turn facilitating higher carriage capacity and higher-speed network systems.

Various network designs adopted various parameters for packet size. The original Ethernet specification, invented in the early 1970s, adopted a variable packet size, with supported packet sizes of between 64 and 1,500 octets. *Fiber Distributed Data Interface* (FDDI), a fibre ring local network, used a variable packet size of up to 4,478 octets. *Frame Relay* used a variable packet size of between 46 and 4,470 octets. The choice of variable-size packets allows applications to refine their behaviour. Jitter and delay-sensitive applications, such as digitised voice, may prefer to use a stream of smaller packets in an attempt to minimise jitter, while reliable bulk data transfer may choose a larger packet size to increase the carriage efficiency. The nature of the medium may also have a bearing on this choice. If there is a high *Bit Error Rate* (BER) probability, then reducing the packet size minimises the impact of sporadic errors within the data stream, possibly increasing throughput in such environments.

In designing a network protocol that is intended to operate over a wide variety of substrate networking media and support as wide a variety of applications as possible, the designers of IP could not rely on a single packet size for all transmissions. Instead, the designers of IPv4 provided a *packet length field* in the packet header. This field was a 16-bit octet count, allowing for an IP packet to be anywhere from the minimum size of 20 octets (corresponding to an IP header without any payload) to a maximum of 65,535 octets.

Obviously not all packets can fit into all underlying media. If the packet is too small for the minimum payload size, then it can be readily padded. But if it's too big for the maximum packet size of the media, then the problem is a little more challenging. IPv4 solved this problem using “forward fragmentation.” The basic approach is that any IPv4 router that is unable to forward an IPv4 packet into the next hop because the packet is too large for the next-hop network may split the packet into a set of smaller “fragments,” copying the original IP header fields into each of these fragments, and then forwarding each of these fragments instead. The fragments continue along the network path as autonomous IP packets, and the destination host is responsible for re-assembling them back into the original IP packet and pass the result, namely the packet as it was originally sent, back up to the local instance of the end-to-end transport protocol.

It is a clever approach, as it hides the entire network-level fragmentation issue from the upper-level protocols, including the *Transmission Control Protocol* (TCP) and *User Datagram Protocol* (UDP). The transport protocols and the upper-level application protocols can, in theory, treat the underlying network as capable of supporting any IP packet size, and the IP layer performs the necessary adaptation to allow the IP packet to traverse any media layer. However, even with this intended transparency of operation, this approach has managed to earn a very poor reputation.

Packet fragmentation was seen as being a source of performance inefficiency, a security vulnerability, and it even posed a cap on maximal delay-bandwidth product on data flows across networks. IP fragmentation was considered harmful^[2,3].

The IPv6 designers removed the fragmentation controls from the common IPv4 packet header and placed them into an 8-octet IPv6 *Extension Header*. This additional packet header, placed between the IPv6 packet header and the end-to-end transport packet header, was present *only* in fragmented packets (whereas the IPv4 fragmentation control fields are present in *all* IPv4 packet headers). Secondly, IPv6 did not permit fragmentation to be performed when the packet was in transit within the network: all fragmentation was to be performed by the packet source prior to transmission. This stipulation, too, has resulted in an uncomfortable compromise, where an unforeseen need for fragmentation relies on *Internet Control Message Protocol for IPv6* (ICMPv6) signalling from the interior of the network back to the original packet source and retransmission.

In the case of TCP, a small amount of layer violation goes a long way. If the sending host is permitted to pass an IPv6 *Packet Too Big* (PTB) ICMPv6 diagnostic message up to the TCP session that generated the original packet, then it's possible for the TCP driver to adjust its sending *Maximum Segment Size* (MSS) to the new, smaller value and carry on. In this case, no fragmentation is required.

UDP is different. In UDP a functional response to path message size issues inevitably relies on interaction with the upper-level application protocol.

It appears that when we consider fragmentation in IPv6 we have to consider the treatment of IPv6 Extension Headers and UDP. And that story is not a robust one^[4].

The DNS and IPv6 Packet Fragmentation

The *Domain Name System* (DNS) is the major user of UDP. As a consequence of the increasing use of *Domain Name System Security Extensions* (DNSSEC) as a security mechanism, coupled with the increasing use of IPv6 as the IP protocol transition gathers momentum, it is time to look once more at the interaction of larger DNS payloads over IPv6.

To illustrate this situation, here are two DNS queries, both made by a recursive resolver to an authoritative name server, both using UDP over IPv6.

Query 1:

```
$ dig +bufsize=4096 +dnssec 000-4a4-000a-000a-  
0000-b9ec853b-241-1498607999-2a72134a.ap2.  
dotnxdomain.net. @8.8.8.8  
139.162.21.135
```

```
(MSG SIZE rcvd: 1190)
```

Query 2:

```
$ dig +bufsize=4096 +dnssec 000-510-000a-000a-0000-
b9ec853b-241-1498607999-2a72134a.ap2.
dotnxdomain.net. @8.8.8.8
status: SERVFAIL

(MSG SIZE rcvd: 104)
```

What we see here are two almost identical DNS queries that have been passed to Google's Public DNS service to resolve.

In the first case, the DNS response is 1,190 octets long, and in the second case the response is 1,346 octets long. The DNS server is an IPv6-only server, and the underlying host of this name server is configured with a local maximum packet size of 1,280 octets. Therefore, in the first case the response being sent to the Google resolver is a single, unfragmented IPv6 UDP packet, and in the second case the response is broken into two fragmented IPv6 UDP packets. And it is this single change that triggers the Google Public DNS Server to provide the intended answer in the first case, but to return a SERVFAIL failure notice in response to the fragmented IPv6 response. When the local *Maximum Transmission Unit* (MTU) on the server is lifted from 1,280 octets to 1,500 octets, the Google resolver returns the server DNS response in both cases.

The only difference in these two responses is IPv6 fragmentation, but there is perhaps more to it than that.

IP fragmentation in both IPv4 and IPv6 “raises the eyebrows” of firewalls. Firewalls typically use the information provided in the transport protocol header of the IP packet to decide whether to admit or deny the packet. For example, you may see firewall rules admitting packets using TCP ports 80 and 443 as a way of allowing web traffic through the firewall filter. For this process to work, the inspected packet needs to contain a TCP header and use the fields in the header to match against the filter set. Fragmentation in IP duplicates the IP portion of the packet header, but the inner IP payload, including the transport protocol header, is not duplicated in every ensuing packet fragment. Thus trailing fragments pose a conundrum to the firewall. Either all trailing fragments are admitted, a situation that has its own set of consequent risks, or all trailing fragments are discarded, a situation that also poses connection issues^[5].

IPv6 adds a further factor to the picture. In IPv4 every IP packet, fragmented or not, contains IP fragmentation control fields. In IPv6 these same fragmentation *control fields* are included in an IPv6 *Extension Header* that is attached only to packets that are fragmented.

This 8-octet Extension Header is placed immediately after the IPv6 packet header in all fragmented packets, meaning that a fragmented IPv6 packet does not contain the *Upper Level Protocol Header* starting at octet offset 40 from the start of the IP packet header. But in the first packet of this set of fragmented packets, the Upper Level Protocol Header is chained off the fragmentation header, at byte offset 48, assuming that there is only a *Fragmentation Extension Header* in the packet. The implications of this fact are quite significant. Instead of always looking at a fixed point in a packet to determine its upper-level protocol, the packet-handling device needs to unravel the Extension Header chain, raising two rather tough questions. First, how long is the device prepared to spend unravelling this chain? And second, would the device be prepared to pass on a packet with an Extension Header that it did not recognise?

In some cases, implementers of IPv6 equipment have found it simpler to just drop all IPv6 packets that contain Extension Headers. Some measurements of this behaviour are reported in RFC 7872^[6]. This document reports a 38% packet-drop rate when sending fragmented IPv6 query packets to DNS Name servers. But the example provided previously is in fact the opposite case to that reported in RFC 7872, and the example illustrates a more conventional case. It's not the queries in the DNS that can readily grow to sizes that require packet fragmentation, but the responses. The relevant question concerns the anticipated probability of packet drop when sending fragmented UDP IPv6 packets as responses to DNS queries. To rephrase the question slightly, how do DNS recursive resolvers fare when the IPv6 response from the server is fragmented?

For a start, it appears from the example cited here that Google's Public DNS resolvers experienced some packet-drop problem when they passed a fragmented IPv6 response (this problem was noted in mid-2017, and Google has subsequently corrected it). But was this problem limited to just one or two DNS resolvers, or do many other DNS resolvers experience a similar packet-drop issue? How widespread is this problem?

We used an experiment that tested resolver capabilities in handling DNS responses that entailed the use of fragmented UDP IPv6 packets. The experiment used a measurement script embedded in an online ad to enlist a large number of endpoints to perform resolution of a domain name^[7]. For this measurement, we altered the DNS resolution system to fragment certain DNS responses.

The approach we took in this experiment was to use a user-level packet-processing system that listens on UDP port 53 and passes all incoming DNS queries to a back-end DNS server. When it receives a response from this back-end server it generates a sequence of IPv6 packets that fragments the DNS payload and uses a raw device socket to pass these packets directly to the device interface.

We are relying on the observation that IPv6 packet fragmentation occurs at the IP level in the protocol stack, so the IPv6 driver at the remote end will reassemble the fragments and pass the UDP payload to the DNS application, and if the resolver receives the payload packets, there will be no trace that the IPv6 packets were fragmented.

The results of this experiment follow:

- 10,851,323 experiments used IPv6 queries for the name server address.
- 6,786,967 experiments queried for the terminal DNS name.
- Fragmented response: $6,786,967 / 10,851,323 = 62.54\% = 37.45\%$ drop.

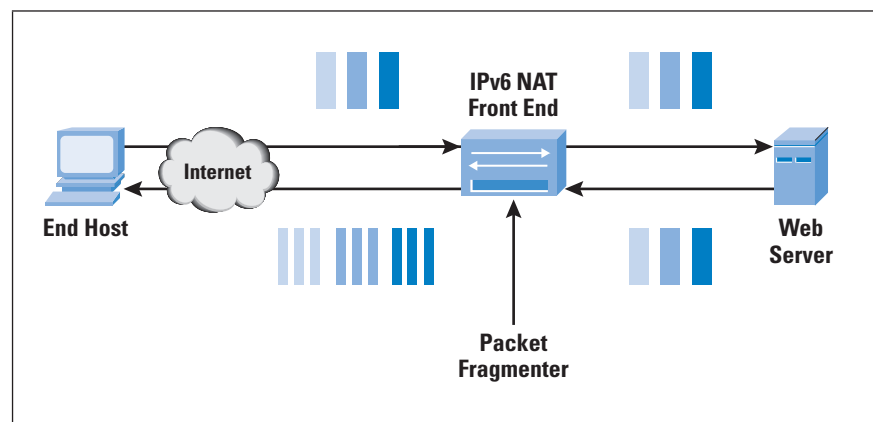
Some 37% of client endpoints used IPv6-capable DNS resolvers that were incapable of receiving a fragmented IPv6 response.

TCP and IPv6 Packet Fragmentation

The use of IPv6 Extension Headers implies that any transport protocol-sensitive functions within network equipment must follow the Extension Header chain of the packet header. This process takes a variable number of cycles for the device. It also requires that the device should recognise all the Extension Headers encountered on the header chain as passing through Extension Headers that the device either does not understand or is not prepared to check to determine whether or not it is a security risk. It's easier to drop all packets with Extension Headers! And that is what a lot of deployed equipment evidently does.

To measure the extent to which equipment drops fragmented IPv6 packets in TCP, we used a front-end unit to a web server and configured this front end to perform packet fragmentation on outbound packets as required. All TCP packets passed across the unit from the back end towards the Internet that contain a TCP payload larger than 15 octets are fragmented (Figure 1).

Figure 1: Experiment Configuration



The subsequent data-analysis phase can detect if the end host has received and successfully reassembled the set of fragments by looking at a log of packets. If an incoming TCP *Acknowledgement* (ACK) has a sequence number that encompasses the sending sequence number of outbound fragments within the same TCP session, that is evidence that the remote end has successfully reassembled the fragmented packet.

How “Real” Is This Experiment?

Before looking at the results, it may be useful to ask whether this experiment represents a “real” scenario that is commonly encountered on the Internet.

It’s certainly the case that in TCP over IPv6 we do not expect to see packet fragmentation in the normal course of events.

A TCP sender should ensure that all outbound TCP segments fit within the local interface MSS size, so in the absence of network path MTU issues, a sender should not be fragmenting outbound TCP packets before sending them.

What about the case where the path MTU is smaller than the local interface MTU? When a packet encounters a network path next hop where the packet is larger than the next-hop MTU, the IPv6 router constructs an ICMPv6 PTB message, noting the size of the next hop, and also including the original packet headers as the payload of this ICMPv6 message. It sends this ICMPv6 diagnostic message back to the original sender and discards the original packet.

When a sending host receives this ICMPv6 PTB message, it also has the TCP packet header as part of the inner payload. This information can be used to find the local TCP control entry for this session, and the outbound MSS value of this TCP session can be updated with the new value. In addition to the updated size information, the TCP header in the ICMPv6 PTB message payload also contains the sequence number of the lost packet. The sending TCP process can interpret the ICMPv6 PTB message as an implicit *Negative Acknowledgement* (NACK) of the lost data, and resend the discarded data, using the updated MSS size. Again, no packet fragmentation is required.

All this sounds like a blatant case of “layer violation” and we should call in the Protocol Police. But before we do so, maybe we should think about the hypothetical situation where the host did not pass the ICMPv6 PTB message to the TCP control block. This situation is analogous to the case where the ICMPv6 PTB message is not passed to the host at all, where, for example, some unhelpful piece of network filtering middleware is filtering out all ICMPv6 messages.

In this case, the sending TCP session has sent a TCP segment and is waiting to receive an ACK. The receiver will not get this packet, so it cannot ACK it. The sender might have a retransmission timer and it might try to resend the offending large packet, but that too will get lost, so it will never get the ACK.

This situation results in a wedged TCP state, or a *Path MTU Black Hole* condition. Hiding ICMPv6 PTB messages from the TCP controller, either because of local processing rules within the host or because some network element has decided to drop them, is invariably harmful.

In that sense, we have constructed a somewhat “unreal” experiment, and we should not expect to see applications that critically depend on the correct working of packet fragmentation in TCP experiencing the same network conditions as those we’ve set up here.

On the other hand, fragmentation is an IP function, not a function performed by an end-to-end transport protocol. Therefore, the question of whether a host can receive a fragmented UDP packet is essentially the same question as whether a host can receive a fragmented TCP packet—at least from the perspective of the host itself. In both cases the real question is whether the IPv6 process on the host can receive fragmented IPv6 packets.

While the experiment itself uses conditions that are essentially an artifice, the result, namely the extent to which IPv6 Extension Header drop occurs when passing fragmented IPv6 packets towards end hosts, is nevertheless a useful and informative result.

Results

Over a period in August 2017, this experiment presented fragmented TCP packets to 1,702,949 unique IPv6 addresses. The results are summarized in Table 1.

Table 1: Results of Fragmentation Test

	Count	% of Total
Sent Fragmented TCP Packets	1,675,898	
Acknowledged Fragmented TCP Packets	1,324,834	79.03%
Failed to Acknowledge Fragmented TCP Packets	351,514	20.97%

Compared to the earlier DNS packet fragmentation result, namely that some 37% of endpoints who used IPv6-capable DNS resolvers used resolvers that were incapable of receiving IPv6 Fragmentation Extension Headers, the overall failure rate observed here of some 21% looks somewhat better. However, “better” is a relative term, as it is still the case that one-fifth of IPv6-capable endpoints are unable to receive a fragmented IPv6 packet.

Having one-fifth of the end-user population incapable of receiving fragmented large responses over IPv6 is indeed a serious problem.

Let's look briefly at the IPv6-over-IPv4 auto-tunnelling techniques Teredo and 6to4 (Table 2), as these two auto-tunnelled IPv6 bridging technologies just don't seem to want to die!

Table 2: Results of IPv6 Fragmentation Test for Teredo and 6to4 Prefixes

	Teredo	%	6to4	%
Sent Fragmented TCP Packets	53,780		24,384	
Acknowledged Fragmented TCP Packets	263	0.5%	1,486	6.1%
Failed to Acknowledge Fragmented TCP Packets	53,517	98.5%	22,898	93.9%

Both of these auto-tunnelling services are atrocious in this respect! Almost no Teredo endpoints can handle IPv6 fragmentation, and the 6to4 failure rate is not much better. Having no IPv6 at all is far better than having such a terrible service, and I can think of few better justifications for turning off the remaining Teredo and 6to4 gateways than these figures! What is even more depressing is that these two auto-tunnelling technologies represent one-quarter of the count of unique /64 prefixes seen in this experiment.

There is a considerable level of variation in the extent to which networks support the delivery of IPv6 Fragmentation Extension Headers to hosts. In some cases, it appears that the choice of customer premises equipment, or the configuration of IPv6 firewalls, may be a factor. Where the failure rate is very high it would point to the drop point being part of the behaviour of the provider network rather than the behaviour of the customer premises equipment.

Conclusions

Whatever the reasons, the conclusion here is unavoidable: IPv6 fragmentation is just not a viable component of the IPv6 Internet.

We need to adjust our protocols to avoid fragmentation.

For TCP, this adjustment should not be a major issue. Of course, this assertion relies on ICMPv6 PTB messages getting back to the sender's TCP process, but that is a major topic in its own right, so we won't delve deeper into this subject right now.

However, for UDP, this conclusion should be cause for some major rethinking of the way the DNS works, as the combination of DNSSEC, UDP, and IPv6 is really not going to work very well. It has implications for other UDP-based protocols as well, particularly where the protocol can generate large payloads.

The *Quick UDP Internet Connections* (QUIC) protocol, which uses a TCP-like control protocol embedded within a UDP encapsulation^[8], has taken the pragmatic position of using a maximum packet size of 1,350 octets as a universal base and does not expect to encounter fragmentation issues given this somewhat conservative choice of packet size. If the DNS over IPv6 used a similar upper limit of UDP packet size and always sent back truncated responses for larger answers, we could probably avoid many of the packet-loss problems that we encounter today. Of course, the consequent larger use of TCP has its own implications in terms of query processing capacity for DNS resolvers and servers, so there are no free points here.

However, one conclusion looks starkly clear to me from these results. We can't just assume that the DNS as we know it today will just work in an all IPv6 future Internet. We must make some changes in some parts of the protocol design to get around this current widespread problem of IPv6 Extension Header packet loss in the DNS, assuming that we want to have a DNS at all in this all-IPv6 future Internet.

References and Further Reading

- [0] J. Postel, "Internet Protocol," RFC 791, September 1981.
- [1] Ron Bonica, Warren Kumari, Randy Bush, and Hagen Pfeifer, "IPv6 Fragment Header Deprecated," July 2013, Internet Draft, Work in Progress, **draft-bonica-6man-frag-deprecate-02**.
- [2] Christopher A. Kent and Jeffrey C. Mogul, "Fragmentation Considered Harmful," Proceedings of Frontiers in Computer Communications Technology, ACM SIGCOMM '87, August 1987.
- [3] Matt Mathis, Ben Chandler, and John W. Heffner, "IPv4 Reassembly Errors at High Data Rates," RFC 4963, July 2007.
- [4] Ron Bonica, Fred Baker, Geoff Huston, Robert Hinden, Ole Troan, and Fernando Gont, "IP Fragmentation Considered Fragile," March 2018, Internet Draft, Work in Progress, **draft-bonica-intarea-frag-fragile-01**.
- [5] Joel Jaeggli, Lorenzo Colitti, Warren Kumari, Eric Vyncke, Merike Kaeo, and Tom Taylor, "Why Operators Filter Fragments and What It Implies," December 2013, Internet Draft, Work in Progress, **draft-taylor-v6ops-fragdrop-02**.
- [6] Fernando Gont, J. Linkova, and Tim Chown, "Observations on the Dropping of Packets with IPv6 Extension Headers in the Real World," RFC 7872, June 2016.

- [7] Geoff Huston, Joao Damas, and George Michaelson, “How we Measure IPv6,” Presentation to APNIC 44 Conference, September 2017.
- [8] Jana Iyengar and Martin Thomson, “QUIC: A UDP-Based Multiplexed and Secure Transport,” April 2018, Internet Draft, Work in Progress, **draft-ietf-quic-transport-11**.
- [9] Pekka Savola, “MTU and Fragmentation Issues with In-the-Network Tunneling,” RFC 4459, April 2006.
- [10] G. Ziemba, D. Reed, and P. Traina, “Security Considerations for IP Fragment Filtering,” RFC 1858, October 1995.
- [11] Geoff Huston, “Fragmentation,” *The Internet Protocol Journal*, Volume 19, No. 2, June 2016.

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. He is an active contributor to the Internet Engineering Task Force. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

Letters to the Editor

Ole and William,

As a loyal IPJ reader for decades, I think you did a terrific job on your blockchain piece in the November 2017 issue of IPJ. I don't think I have seen anything yet as comprehensive, understandable, and interesting to read on this subject. Thanks for writing and publishing the piece. There is far too much fluff and misinformation on blockchain and cryptocurrencies, so it is nice to see a solid treatment that provides helpful information to further everyone's understanding of this important technology.

—David Strom, david@strom.com

The author responds:

David,

Thank you for the kind words. A lot of credit goes to the reviewers, who did an awesome job of providing detailed feedback on the first draft.

—Bill Stallings, ws@shore.net

Ole,

I just re-read Geoff Huston's article on *Network Address Translation* (NAT) in the November 2017 issue (Volume 20, No. 3), and needed to applaud—possibly for the second time. Geoff is in a league of his own—the clear thinking, the pragmatism, the ability to communicate clearly and understandably—and not the least, the perspective on history. How we got here and why. Enlightening. NAT changed the Internet and continues to do so in ways I had never thought of until I read the article. Understanding that the 2018 Internet is vastly different in almost every possible way from the 1995 and 2005 Internet is extremely important. Otherwise we continue to plan for an historic model rather than the future. IPJ rocks!

—Helge Skrivervik, helge@mymayday.com

Ole,

As a data-networking engineer and architect of 20 years on the front line of fortune 100 network projects, I would like to offer a counter perspective from the recent article “In Defence of NATs” (Volume 20, No. 3).

The networking world seems to be losing sight that NAT is a “crutch” of sorts, a way of dealing with the primary problem in a lack of IPv4 address space.

By trying to justify NAT as a way to scale up IPv4 future potential scalability by “stealing” port/socket bits for something other than their originally intended purpose is nonsensical. One correction I would like to make on the article is the claim that NAT provides a firewall function. NAT and firewall functionality are mutually exclusive mechanisms, even if they are most often found on the same network device. NAT provides an obscured view of a host from elsewhere via address/port translation but does not by itself provide natural protection from the host on the other side of the translation point. Firewall protection encompasses the scrutinizing and controlling of the traffic that is allowed to traverse the NAT point. Firewalls provide this function with or without NAT, and NAT can function with or without a firewall mechanism.

Slow adoption of IPv6 has nothing to do with any perceived brilliant nature of NAT, and NAT does present real-world problems for several types of very important software; Microsoft’s *Active Directory Replication* and IBM’s *Virtual Tape Library (VTL)/Virtual Tape Server (VTS)* are two examples off the top of my head. When traversing a NAT point, many applications that share their host IP address in the data field with other hosts require active “swapping” of this payload-imbedded IP address or another compensating mechanism. The communication would “break” the application’s intended communication model without the addition of a compensating mechanism.

The question in the article “should I deploy IPv6 now?” is the wrong question. The Internet was first “born” in a practical sense at the point it became commercially available to the world to use. In the first 25 years it grew at an amazing rate, doubling its size several times over in a very natural and organic way. It has been almost a quarter century since IPv6 was released, and its resistance has been *significant* for good reasons beyond the scope of this letter. IPv6 is “The Emperor’s New Clothes,” otherwise IPv4 would have been replaced by now if IPv6 had been a natural and organic progression of IPv4, and there would be no need to “defend” NAT or speak up against its forced ubiquitous overuse. Someone once told me that IPv6 was here to stay. To my way of thinking it has not yet arrived after almost a quarter century.

While IPv4 port/socket numbers can be seen as “borrowed address bits” for NAT, I believe it is a distorted view of port/socket intended use. Fields and protocols are defined and delineated for a reason. If you wish to repurpose bits—for example, *Type of Service (ToS)* to *Differentiated Services Code Point (DSCP)* use—then repurpose them officially. Until then, fields and protocols should be respected as originally intended and not subject to implied de-facto depreciation by NAT’s liberal theft. Field definitions and structure have purpose.

The sirens' song that I believe we collectively are starting to fall for is that NAT is a "one-size-fits-all" solution for all forms of network scalability in a ubiquitous way. This is not the real-world case. Resource hosts need (in a practical sense) globally unique Layer 3 identifiers. Consider corporate mergers/acquisitions as well as divestitures leading to merging of a divested entity. Trying to merge two significant company networks together that both use NAT RFC 1918 on the "inside" for resource hosts is overly complex in a way that it would not have to be if a viable replacement for IPv4 had been rolled out in a manner that world corporations could embrace commercially. That protocol does not exist. While I agree with the notion that the Internet cannot be "stateless," this does not uniquely justify NAT as "middleware." End hosts ideally should be ultimate keepers of their stateful connections; justifying NAT just for the sake of IPv4 continued life-support is nonsensical.

NAT has a proper use; middleware that uses NAT as a complementary protocol along with others, such as "load balancing" [for example, the *Local Traffic Manager* (LTM) product by F5] is justifiable, but in this case application scalability and fault tolerance encompass the direct purpose of this middleware, not compensating for world IPv4 address depletion. Other forms of network middleware are perfectly justifiable, even if NAT did not need to exist; firewalls and *Intrusion Prevention Systems* (IPS) are examples. We should appreciate NAT for its role as a "tactical" compensating mechanism for IPv4 address space depletion, not as a "strategic future-proofing" scalability mechanism for IPv4. People with broken legs appreciate a crutch, but would not appreciate needing to use a crutch for the rest of their lives if their body decided not to heal itself because the body viewed the crutch as "good enough." NAT seen as a long-term way of extending IPv4 scalability and, therefore, lifespan is just putting lipstick on the IPv4 pig.

NAT is a mechanism to be used (like any other protocol) where it makes sense to use it and no further. If IPv6 or some other more sensible replacement for IPv4 were completely rolled out with IPv4 relegated to the history books, then the practical use of NAT in such a future environment would be a single-digit fraction of its current existence. That existence would be primarily in terms of resiliency mechanisms such as load balancing, as previously mentioned, and certain (client) mobility cases. One of the most memorable pieces of wisdom I have ever heard about IT applies here very well: "There is nothing more permanent than a temporary solution." Let us not fall victim to this easy psychological trap only because we seem to have collectively painted ourselves into a corner of sorts.

—Leroy Harvey, leroy.harvey@hotmail.com

The author responds:

The purpose of any opinion piece is to provoke the reader into thinking about the issue, and perhaps looking at it from a set of different perspectives. For more than two decades the *Internet Engineering Task Force* (IETF) viewed NATs as a somewhat ugly hack, and from time to time attempted to discourage its use in various ways. However, the undeniable observation is that NATs keep today's Internet running. Perhaps there is more to NATs than a rather ugly short-term hack that should disappear. What if they are here to stay? The opinion piece was intended to look at NATs from a perspective that shared little with the orthodox view of NATs, and provoke readers to think about this unanticipated direction that the Internet has taken and wonder where it may lead. I'm pleased to see that this provocation has motivated one reader to provide a thoughtful response.

—Geoff Huston, APNIC
gih@apnic.net

Letters may be edited for clarity. We'd love to hear from you. Send us your feedback via e-mail to ipj@protocoljournal.org

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

Coming Soon: Our 20th Anniversary Issue

It is difficult to believe, but another decade has passed and in June we will celebrate 20 years of *The Internet Protocol Journal*. Make sure your subscription is up-to-date so you don't miss this issue!

Ten years ago:



Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Fabrizio Accatino	Dave Crocker	Serge Van Ginderachter	Jonatan Jonasson
Scott Aitken	Kevin Croes	Greg Goddard	Daniel Jones
Antonio Cuñat Alario	John Curran	Octavio Alfageme	Gary Jones
Matteo D'Ambrosio	André Danthine	Gorostiaga	Amar Joshi
Jens Andersson	Morgan Davis	Barry Greene	Merike Kaeo
Danish Ansari	Freek Dijkstra	Martijn Groenleer	Andrew Kaiser
Tim Armstrong	Geert Van Dijk	Geert Jan de Groot	David Kekar
Richard Artes	Richard Dodsworth	Gulf Coast Shots	Shan Ali Khan
David Atkins	Ernesto Doelling	Sheryll de Guzman	Nabeel Khatri
Jaime Badua	Eugene Doroniuk	James Hamilton	Anthony Klopp
John Bigrow	Karlheinz Dölger	Stephen Hanna	Henry Kluge
Axel Boeger	Andrew Dul	Martin Hannigan	Andrew Koch
Gerry Boudreaux	Holger Durer	John Hardin	Carsten Koempe
Kevin Breit	Peter Robert Egli	David Harper	Alexader Kogan
Ilia Bromberg	George Ehlers	Edward Hauser	Antonin Kral
Christophe Brun	Peter Eisses	David Hauweele	Mathias Körber
Gareth Bryan	Torbjörn Eklöv	Headcrafts SRLS	John Kristoff
Stefan Buckmann	ERNW GmbH	Johan Helsingius	Terje Krogdahl
Scott Burleigh	ESdatCo	Robert Hinden	Bobby Krupczak
Jon Harald Bøvre	Steve Esquivel	Alain Van Hoof	Murray Kucherauw
Olivier Cahagne	Mikhail Evstiounin	Edward Hotard	Warren Kumari
Tracy Camp	Paul Ferguson	Bill Huber	Darrell Lack
Fabio Caneparo	Kent Fichtner	Hagen Hultzs	Yan Landriault
Roberto Canonico	Gary Ford	Mika Ilvesmaki	Markus Langenmair
John Cavanaugh	Christopher Forsyth	Karsten Iwen	Fred Langham
Lj Cemer	Craig Fox	Ashford Jaggernaut	Richard Lamb
Dave Chapman	Fausto Franceschini	David Jaffe	Tracy LaQuey Parker
Stefanos Charchalak	Tomislav Futivic	John Jarvis	Simon Leinen
Greg Chisholm	Edward Gallagher	Dennis Jennings	Robert Lewis
Brad Clark	Andrew Gallo	Edward Jennings	Sergio Loreti
Narelle Clark	Chris Gamboni	Aart Jochem	Guillermo a Loyola
Steve Corbató	Xosé Bravo Garcia	Richard Johnson	Hannes Lubich
Brian Courtney	Kevin Gee	Jim Johnston	Dan Lynch

Miroslav Madić	Alexis Panagopoulos	Carsten Scherb	Luca Ventura
Alexis Madriz	Gaurav Panwar	Roger Schwartz	Tom Vest
Carl Malamud	Manuel Uruena Pascual	SeenThere	Dario Vitali
Michael Malik	Ricardo Patara	Scott Seifel	Randy Watts
Yogesh Mangar	Dipesh Patel	Yury Shefer	Andrew Webster
Bill Manning	Alex Parkinson	Yaron Sheffer	Tim Weil
Harold March	Craig Partridge	Tj Shumway	Jd Wegner
Vincent Marchand	Dan Paynter	Jeffrey Sicuranza	Rick Wesson
David Martin	Leif-Eric Pedersen	Thorsten Sideboard	Peter Whimp
Timothy Martin	Juan Pena	Andrew Simmons	Jurrien Wijlhuizen
Gabriel Marroquin	Chris Perkins	Henry Sinnreich	Pindar Wong
Carles Mateu	David Phelan	Geoff Sisson	Bernd Zeimetz
Juan Jose Marin Martinez	Derrell Piper	Helge Skrivervik	廖明沂.
Ioan Maxim	Rob Pirnie	Darren Sleeth	
Miles McCredie	Jorge Ivan Pincay Ponce	Bob Smith	
Brian McCullough	Blahoslav Popela	Mark Smith	
Joe McEachern	Tim Pozar	Job Snijders	
Jay McMaster	David Raistrick	Asit Som	
Carsten Melberg	Priyan R Rajeevan	Ignacio Soto Campos	
Kevin Menezes	Paul Rathbone	Peter Spekrijse	
Bart Jan Menkveld	Bill Reid	Thayumanavan Sridhar	
William Mills	Rodrigo Ribeiro	Matthew Stenberg	
Desiree Miloshevic	Justin Richards	Adrian Stevens	
Thomas Mino	Mark Risinger	Clinton Stevens	
Mohammad Moghaddas	Ron Rockrohr	John Streck	
Charles Monson	Carlos Rodrigues	Viktor Sudakov	
Andrea Montefusco	Lex Van Roon	Edward-W. Suor	
Fernando Montenegro	William Ross	Vincent Surillo	
Soenke Mumm	Boudhayan Roychowdhury	Roman Tarasov	
Tariq Mustafa	Carlos Rubio	David Theese	
Stuart Nadin	Timo Rüter	Sandro Tumini	
Mazdak Rajabi Nasab	RustedMusic	Phil Tweedie	
Krishna Natarajan	Babak Saberi	Steve Ulrich	
Darryl Newman	George Sadowsky	Unitek Engineering AG	
Marijana Novakovic	Scott Sandefur	John Urbanek	
Ovidiu Obersterescu	Sachin Sapkal	Martin Urwaleck	
Mike O'Connor	Arturas Satkovskis	Betsy Vanderpool	
Mike O'Dell	Phil Scarr	Surendran Vangadasalam	
Carlos Astor Araujo Palmeira	Jeroen Van Ingen Schenau	Alejandro Vennera	



Follow us on Twitter and Facebook

@protocoljournal



<https://www.facebook.com/newipj>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Diamond Sponsors



Ruby Sponsor



Sapphire Sponsors

Your logo here!

Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

ADDRESS SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

David Conrad, Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

August 2018

Volume 21, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

FROM THE EDITOR

In This Issue

From the Editor	1
Another 10 Years.....	2
Fileless Malware	17
Fragments	26
Thank You.....	28
Call for Papers.....	30
Supporters and Sponsors	31

In June 1998 we published the first issue of *The Internet Protocol Journal* (IPJ). Since then, we have produced 75 issues and a total of 2,848 pages. Today, IPJ has about 22,000 subscribers all around the world. Although two-thirds of our readers prefer the paper edition, a growing number of subscribers are downloading the PDF version instead. As we remarked in 2008, this shift in reading habits is related to the emergence of low-cost, high-resolution displays and printers, as well as improvements in Internet access technologies, particularly with respect to mobile devices.

In this 20th anniversary year, we decided to ask our most frequent contributor, Geoff Huston, to reflect on Internet developments since 2008. His article, “Another 10 Years,” outlines the many ways in which the Internet has changed in this period, as well as a few areas where developments are still lacking.

The Internet continues to be used for numerous unsavory activities including fraud, identity theft, malicious software intrusion, denial-of-service incursions, and much more. These cyber attacks are getting increasingly sophisticated, as David Strom explains in his article entitled “Fileless Malware.”

As I did 10 years ago, let me take this opportunity to thank all those people who make IPJ possible. Our authors deserve a round of applause for carefully explaining both established and emerging technologies. They are assisted by an equally insightful set of reviewers and advisors who provide feedback and suggestions on every aspect of our publications process. The process itself relies heavily on two individuals: Bonnie Hupton, our copy editor, and Diane Andrada, our designer. Of equal importance are our numerous individual donors and corporate sponsors, without whom we would be unable to publish and distribute the journal. Last, but not least, our readers give us encouragement, suggestions, and feedback that enables us to provide the most relevant material.

If you are wondering why this issue has a cover date of “August” rather than “June,” I can only apologize and blame it on a busy summer and a broken arm. I appreciate your continued patience and support.

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

Another 10 Years

by Geoff Huston, APNIC

The evolutionary path of any technology can often take strange and unanticipated turns and twists. At some points simplicity and minimalism can be replaced by complexity and ornamentation, while at other times a dramatic cut-through exposes the core concepts of the technology and removes layers of superfluous additions. The evolution of the Internet appears to be no exception, and it contains these same forms of unanticipated turns and twists. In thinking about the technology of the Internet over the last 10 years, it appears that it's been a very mixed story about what has changed and what has stayed the same.

A lot of the Internet today looks much the same as the Internet of a decade ago^[0,1]. Much of the infrastructure of the Internet has stubbornly resisted various efforts to engender change. We are still in the middle of the process to transition the Internet to IPv6, as was the case a decade ago. We are still trying to improve the resilience of the Internet to various attack vectors, as also was true a decade ago. We are still grappling with various efforts to provide defined *Quality of Service* (QoS)^[2] in the network, also true a decade ago. It seems that the rapid pace of technical change in the 1990s and early 2000s has simply run out of momentum, and that the dominant activity on the Internet over the past decade was consolidation rather than continued technical evolution. Perhaps this increased resistance to change is because as the size of the network increases, its inertial mass also increases. We used to quote *Metcalf's Law*^[3] to each other, reciting the mantra that the value of a network increases in proportion to the square of the number of users.

A related observation appears to be that the inherent resistance of a network to change, or its inertial mass, is also directly related to the square of the number of users. Perhaps as a general observation, all large, loosely coupled, distributed systems are strongly resistant to efforts to orchestrate a coordinated change. At best, these systems respond to various forms of market pressures, but because the overall system of the Internet is so large and so diverse, these market pressures manifest themselves in different ways in different parts of this network. Individual actors operate under no centrally orchestrated set of instructions or constraints. Where change occurs, it is because some sufficiently large body of individual actors sees opportunity in undertaking the change or perceives unacceptable risk in not changing. The result for the Internet appears to be that some changes are very challenging, while others look like natural and inevitable progressive steps.

But the other side of the story is one that is about as diametrically opposed as it's possible to paint. Over the last decade, we've seen another profound revolution in the Internet as it embraced a combination of wireless-based infrastructure and a rich set of services at a speed that has been unprecedented. We've seen a revolution in content and content provision that has changed the Internet, and as collateral damage the Internet appears to be decimating the traditional newspaper and broadcast television sectors. Social media has all but replaced the social role of the telephone and the practice of letter writing. We've seen the rise of the resurgence of a novel twist to the old central mainframe service in the guise of the *cloud*^[4,5] and the repurposing of Internet devices to support views of a common cloud-hosted content that in many ways mimic the function of display terminals of a bygone past. All of these developments are fundamental changes to the Internet and all of them have occurred in the last decade!

That's a significant breadth of material to cover, so I'll keep the story to the larger themes, and to structure this story, rather than offer a set of unordered observations about the various changes and developments over the past decade, I'll use a standard model of a protocol stack as the guiding template. I'll start with the underlying transmission media and then look at IP, the transport layer, and applications and services, and then close by looking at the business of the Internet to highlight developments of the last decade.

Below the IP Layer

What's changed in network media?

Optical systems have undergone sustained change in the past decade. A little over a decade ago production optical systems used simple on-off keying to encode the signal into the optical channel. The speed increases in this generation of optical systems relied on improvements in the silicon control systems and the laser driver chips. The introduction of *Wavelength-Division Multiplexing* (WDM) in the late 1990s allowed the carriers to greatly increase the carrying capacity of their optical cable infrastructure. The last decade has seen the evolution of optical systems into areas of polarisation and phase modulation to effectively lift the number of bits of signal per baud. These days 100-Gbps optical channels are commonly supportable, and we are looking at further refinements in signal detection to lift that beyond 200 Gbps. We anticipate 400-Gbps systems in the near future, using various combinations of a faster basic baud rate and higher levels of phase amplitude modulation, and dare to think that 1 Tbps is now a distinct near-term optical service.

Radio systems have seen a similar evolution in overall capacity. Basic improvements in signal processing, analogous to the changes in optical systems, have allowed the use of phase modulation to lift the data rate of the radio bearer.

The use of *Multiple-Input* and *Multiple-Output* (MIMO) technology, coupled with the use of higher carrier frequencies, has allowed the mobile data service to support carriage services of up to 100 Mbps in today's *Fourth-Generation* (4G) networks. The push to even higher frequencies promises speeds of up to 1 Gbps for mobile systems in the near future with the deployment of 5G technology.

While optical speeds are increasing, Ethernet packet framing still persists in transmission systems long after the original rationale for the packet format died along with that bright-yellow coaxial cable! Oddly enough, the Ethernet-defined minimum and maximum packet sizes of 64 and 1500 octets still persist. The inevitable result of faster transmission speeds with constant packet sizes results in an upper bound of the number of packets per second increasing more than 100-fold over the past decade, in line with the increase of deployed transmission speeds from 2.5 to 400 Gbps. As a consequence, silicon-based switches are demanding higher packet-processing rates. But one really important scaling factor has not changed for the past decade, namely the clock speed of processors and the cycle time of memory, which have not moved at all. The response so far has been in increasing reliance of parallelism in high-speed digital switching applications, and these days multi-core processors and highly parallel memory systems are used to achieve performance that would be impossible in a single threaded processing model.

In 2018, it appears that we are close to achieving 1-Tbps optical systems and up to 20 Gbps in radio systems. Just how far and how quickly these transmission models can be pushed into supporting ever-higher channel speeds is an open question.

The IP Layer

The most notable aspect of the network that appears to stubbornly resist all forms of pressure over the last decade, including some harsh realities of acute scarcity, is the observation that we are still running what is essentially an IPv4 Internet.

Over the past decade, we have exhausted our pools of remaining IPv4 addresses, and in most parts of the world the IPv4 Internet is running on some form of empty. We had never suspected that the Internet would confront the exhaustion of one of its most fundamental pillars—the basic function of uniquely addressing connected devices—and apparently shrug it off and continue on blithely. But, unexpectedly, that's exactly what has happened.

Today we estimate that some 3.4 billion people regularly use the Internet, and some 20 billion devices are connected to it. We have achieved this feat by using some 3 billion unique IPv4 addresses. Nobody thought that we could achieve this astonishing feat, yet it has happened with almost no fanfare.

Back in the 1990s we had thought that the prospect of address exhaustion would propel the Internet to use IPv6, which was the successor IP protocol that comes with a four-fold increase in the bit width of IP addresses. By increasing the IP address pool to some esoterically large number of unique addresses (340 *undecillion* addresses, or 3.4×10^{38}), we would never have to confront network address exhaustion again. But this transition was not going to be easy. There is no backward compatibility in this protocol transition, so everything has to change. Every device, every router, and even every application needs to change to support IPv6. Rather than perform comprehensive protocol surgery on the Internet and change every part of the infrastructure to support IPv6, we changed the basic architecture of the Internet instead. Oddly enough, it looks like this option was the cheaper one!

Through the almost ubiquitous deployment of *Network Address Translators* (NATs)^[6, 7] at the edges of the network, we've transformed the network from a *peer-to-peer* network into a *client/server* network. In today's client/server Internet clients can talk to servers, and servers can talk back to these connected clients, but that's it. Clients cannot talk directly to other clients, and servers need to wait for the client to initiate a conversation in order to talk to a client. Clients "borrow" an endpoint address when they are talking to a server and release this address for use by other clients when they are idle. After all, endpoint addresses are only useful to clients in order to talk to servers. The result is that we've managed to cram some 20 billion devices into an Internet that has deployed only 3 billion public address slots. We've achieved this result though embracing what could be described as *time-sharing* of IP addresses.

All well and good, but what about IPv6? Do we still need it? If so, then are we going to complete this protracted transition? Ten years later the answer to these questions remains unclear. On the positive side, there is a lot more IPv6 usage around now than there was 10 years ago. *Internet Service Providers* (ISPs) are deploying much more IPv6 today than they did in 2008. When IPv6 is deployed within a service provider's network, we see an immediate uptake from these IPv6-equipped devices. In 2018, it appears that one-fifth of Internet users (that itself is now estimated to number around one-half of the human population on the planet) are capable of using the Internet over IPv6, and most of this capability has developed in the past 10 years. However, on the negative side the question must be asked: What's happening with IPv6 for the other four-fifths of the Internet? Some ISPs have made the case that they would prefer to spend their finite operating budgets on other areas that improve their customers' experience, such as increasing network capacity, removing data caps, or acquiring more on-net content. Such ISPs continue to see deployment of IPv6 as a deferrable measure.

It seems that today we are still seeing a mixed picture for IPv6. Some service providers simply see no way around their particular predicament of IPv4 address scarcity, and these providers see IPv6 as a necessary decision to further expand their network. Other providers are willing to defer the question to some undefined point in the future.

Routing

While we are looking at what's largely unchanged over the past decade we need to mention the routing system. Despite dire predictions of the imminent scaling death of the *Border Gateway Protocol* (BGP)^[8] 10 years ago, BGP has steadfastly continued to route the entire Internet. Yes, BGP is as insecure as ever, and yes, a continual stream of fat-finger foul-ups and less common but more concerning malicious route hijacks continue to plague our routing system, but the routing technologies used in 2008 are the same as those we use in today's Internet.

The size of the IPv4 routing table has tripled in the past 10 years, growing from 250,000 entries in 2008 to slightly more than 750,000 entries today. The IPv6 routing story is more dramatic, growing from 1,100 entries to 52,000 entries. Yet BGP just quietly continues to work efficiently and effectively. Who would've thought that a protocol that was originally designed to cope with a few thousand routes announced by a few hundred networks could still function effectively across a routing space approaching a million routing entries and a hundred thousand networks!

In the same vein, we have not made any major change to the operation of our interior routing protocols. Larger networks still use either *Open Shortest Path First* (OSPF)^[9] or *Intermediate System-to-Intermediate System* (IS-IS) depending on their circumstances, while smaller networks may opt for some distance vector protocol like *Routing Information Protocol Version 2* (RIPv2)^[10] or even *Enhanced Interior Gateway Routing Protocol* (EIGRP)^[11]. The work in the *Internet Engineering Task Force* (IETF) on more recent routing protocols such as the *Locator Identifier Separation Protocol* (LISP)^[12] and the *Babel Routing Protocol*^[13] seem to lack any real traction with the Internet at large. While they both have interesting properties in routing management, neither has a sufficient level of perceived benefit to overcome the considerable inertia of conventional network design and operation. Again, this example looks like another instance where inertial mass is exerting its influence to resist change in the network.

Network Operations

Speaking of network operation, we are seeing some stirrings of change, but it appears to be a rather conservative area, and adoption of new network management tools and practices takes time.

The Internet converged on using the *Simple Network Management Protocol* (SNMP) a quarter of a century ago, and despite its security weaknesses, its inefficiency, its incredibly irritating use of *Abstract Syntax Notation One* (ASN.1), and its use in sustaining some forms of *Distributed Denial-of-Service* (DDoS) attacks, it still enjoys widespread use. But SNMP is only a network monitoring protocol, not a network configuration protocol, as anyone who has attempted to use SNMP write operations can attest.

The more recent *Network Configuration Protocol* (NETCONF) and the *Yet Another Next Generation* (YANG) data modelling language are attempting to pull this area of configuration management into something a little more usable than *Command-Line Interface* (CLI) scripts driving interfaces on switches. At the same time, we are seeing orchestration tools such as *Ansible*, *Chef*, *Network Automation and Programmability Abstraction Layer with Multivendor* (NAPALM) and SALT enter the network operations space, permitting the orchestration of management tasks over thousands of individual components. These network operations management tools are welcome steps forward to improve the state of automated network management, but it's still far short of a desirable endpoint.

In the same time period as we appear to have advanced the state of automated control systems to achieve the driverless autonomous car, the task of fully automated network management appears to have fallen way short of the desired endpoint. Surely it must be feasible to feed an adaptive autonomous control system with the network infrastructure and available resources, and allow the control system to monitor the network and modify the operating parameters of network components to continuously meet the service-level objectives of the network? Where's the driverless car for driving networks? Maybe the next 10 years might get us there.

The Mobile Internet

Before we move up a layer in the Internet Protocol model and look at the evolution of the end-to-end transport layer, we probably need to talk about the evolution of the devices that connect to the Internet.

For many years, the Internet was the domain of the desktop personal computer, with laptop devices serving the needs to those with a desire for a more portable device. At the time the mobile phone was still just a phone, and its early forays into the data world were unimpressive.

Apple's iPhone, released in 2007, was a revolutionary device. Boasting a vibrant-colour touch-sensitive screen, just four keys, a fully functional operating system with Wi-Fi and cellular radio interfaces, and a capable processor and memory, its entry into the consumer market space was perhaps the major event of the decade. Apple's early lead was rapidly emulated by Windows and Nokia with their own offerings.

Google's position was more as an active disruptor, using an open licensing framework for the Android platform and its associated application ecosystem to empower a collection of handset assemblers. Samsung, LG, HTC, Huawei, Sony, and Google, to name a few, all use Android. These days almost 80% of the mobile platforms use Android, and some 17% use Apple's iOS.

For the human Internet, the mobile market is now the Internet-defining market in terms of revenue. There is little in terms of margin or opportunity in the wired network these days, and even the declining margins of these mobile data environments represent a vague glimmer of hope for the one dominant access provider industry.

Essentially, the public Internet is now a platform of apps on mobile devices.

End-to-End Transport Layer

It's time to move up a level in the protocol stack and look at end-to-end transport protocols and changes that have occurred in the past decade.

End-to-end transport was the revolutionary aspect of the Internet, and the *Transmission Control Protocol* (TCP)^[14] was at the heart of this change. Many other transport protocols require the lower levels of the network protocol stack to present a reliable stream interface to the transport protocol. It was up to the network to create this reliability, performing data integrity checks and data flow control, and repairing data loss within the network as it occurred. TCP dispensed with all of that, and simply assumed an unreliable datagram transport service from the network and pushed the responsibility for data integrity and flow control to the transport protocol.

In the world of TCP, not much appears to have changed in the past decade. We've seen some further small refinements in the details of the TCP controlled rate increase and rapid rate decrease, but nothing that shifts the basic behaviours of this protocol. TCP tends to use packet loss as the signal of congestion and oscillates its flow rate between some lower rate and this loss-triggering rate.

Or at least that was the case until quite recently. The situation is poised to change, and change in a very fundamental way, with the debut of Google's offerings of *Bottleneck Bounded Rate* (BBR) and *Quick UDP Internet Connections* (QUIC).

The BBR control algorithm is a variant of the TCP flow-control protocol that operates in a very different mode from other TCP protocols. BBR attempts to maintain a flow rate that sits exactly at the delay bandwidth product of the end-to-end path between sender and receiver. In so doing, BBR tries to avoid the accumulation of data buffering in the network (when the sending rate exceeds the path capacity), and also tries to avoid leaving idle time in the network (where the sending rate is less than the path capacity).

The side effect is that BBR tries to avoid the collapse of network buffering when congestion-based loss occurs. BBR achieves significant efficiencies from both wired and wireless network transmission systems.

The second recent offering from Google also represents a significant shift in the way we use transport protocols. The QUIC protocol looks like a *User Datagram Protocol* (UDP) protocol from the perspective of the network. But in this case looks are deceiving. The inner payload of these UDP packets contain a more conventional TCP flow-control structure and a TCP stream payload. However, QUIC encrypts its UDP payload so the entire inner TCP control is completely hidden from the network. The ossification of the Internet transport is due in no small part to the intrusive role of network middleware that is used to discard packets that it does not recognise. Approaches such as QUIC allow applications to break out of this regime and restore end-to-end flow management as an end-to-end function without any form of network middleware inspection or manipulation. I'd call this development as perhaps the most significant evolutionary step in transport protocols over the entire decade.

The Application Layer

Let's keep on moving up the protocol stack and look at the Internet from the perspective of the applications and services that operate across the network.

Privacy and Encryption

As we noted in looking at developments in end-to-end transport protocols, encryption of the QUIC payload is not just to keep network middleware from meddling with the TCP control state, although it does successfully achieve that objective. The encryption applies to the entire payload, and it points to another major development in the past decade. We are now wary of the extent to which various forms of network-based mechanisms are used to eavesdrop on users and services. The documents released by Edward Snowden in 2013 portrayed a very active US Government surveillance program that used widespread traffic-interception sources to construct profiles of user behaviour and inference profiles of individual users. In many ways this effort to assemble such profiles is not much different from what advertising-funded services such as Google and Facebook have been (more or less) openly doing for years, but perhaps the essential difference is that of knowledge and implied consent. In the advertisers' case this information is intended to increase the profile accuracy and hence increase the value of the user to the potential advertiser. The motivations of government agencies are more open to various forms of interpretation, and not all such interpretations are benign.

One technical response to the implications of this leaked material has been an overt push to embrace end-to-end encryption in all parts of the network. The corollary has been an effort to allow robust encryption to be generally accessible to all, and not just a luxury feature available only to those who can afford to pay a premium.

The *Let's Encrypt*^[15] initiative has been incredibly successful in publishing X.509 domain name certificates that are free, and the result is that all network service operators, irrespective of their size or relative wealth, can afford to use encrypted sessions, in the form of *Transport Layer Security* (TLS)^[16], for their web servers.

The push to hide user traffic from the network and network-based eavesdroppers extends far beyond QUIC and TLS session protocols. The *Domain Name System* (DNS) is also a rich source of information about what users are doing; it also is used in many places to enforce content restrictions. There have been recent moves to try to clean up the overly chatty nature of the DNS protocol, using query name minimisation to prevent unnecessary data leaks, and developing both DNS over TLS and DNS over *Secure HTTP* (HTTPS) to secure the network path between a stub resolver and its recursive server. This effort is very much a work in progress at present, and it will take some time to see if the results of this work will be widely adopted in the DNS environment.^[20, 21]

We are now operating our applications in an environment of heightened paranoia. Applications do not necessarily trust the platform on which they are running, and we are seeing efforts from the applications to hide their activity from the underlying platform. Applications do not trust the network, and are increasingly using end-to-end encryption to hide their activity from network eavesdroppers. The use of identity credentials within the encrypted session establishment also acts to limit the vulnerability of application clients to be misdirected to masquerading servers.

The Rise and Rise of Content

Moving further up the protocol stack to the environment of content and applications, we have also seen some revolutionary changes over the past decade.

For a small period of time the content and carriage activities of the Internet existed in largely separate business domains, tied by mutual interdependence. The task of carriage was to carry users to content, which implied that carriage was essential to content. But at the same time a client/server Internet bereft of servers is useless, so content is essential to carriage. In a world of re-emerging corporate behemoths, such mutual interdependence is unsettling, both to the actors directly involved and to the larger public interest.

The content industry is largely the more lucrative of these two industries and enjoys far less in the way of regulatory constraint. There is no concept of any universal service obligation, or even any effective form of price control in the services content providers offer. Many content service providers use internal cross-funding that allows them to offer free services to the public, as in free e-mail, free content hosting, free storage, and similar services, and fund these services through a second, more occluded, transaction that essentially sells the user's consumer profile to the highest-bidding advertiser.

All this activity happens outside of any significant regulatory constraint, a situation that has given the content-services industry both considerable wealth and considerable commercial latitude.

It should be no surprise that this industry is now using its capability and capital to eliminate its former dependence on the carriage sector. We are now seeing the rapid rise of the *Content Delivery Network* (CDN) model, where instead of an Internet carrying the user to a diverse set of content stores, the content stores are opening local content outlets right next to the user. As all forms of digital services move into CDN hostels, and as the CDN opens outlets that are positioned immediately adjacent to pools of economically valuable consumers, then where does that leave the traditional carriage role in the Internet? The outlook for the public carriage providers is not looking all that rosy given this increasing marginalisation of carriage in the larger content economy.

Within these CDNs we've also seen the rise of a new service model enter the Internet in the form of cloud services. Our computers are no longer self-contained systems with processing and compute resources; they look more and more like a window that sees the data stored on a common server. Cloud services are very similar, where the local device is effectively a local cache of a larger backing store. In a world where users may have multiple devices, this model makes persuasive sense, because the view to the common backing store is constant irrespective of the device used to access the data. These cloud services also make data sharing and collaborative work far easier to support. Rather than creating a set of copies of the original document and then attempting to stitch back all the individual edits into a single common whole, the cloud model shares a document by simply altering the access permissions of the document. There is only ever one copy of the document, and all edits and comments on the document are available to all.

The Evolution of Cyber Attacks

At the same time as we have seen announcements of ever-increasing network capacity within the Internet, we've seen a parallel set of announcements that note new records in the aggregate capacity of *Denial-of-Service* (DoS) attacks. The current peak volume is an attack of some 1.7 Tbps of malicious traffic.

Attacks are now commonplace. Many of them are brutally simple, relying on a tragically large pool of potential zombie devices that are readily subverted and co-opted to assist in attacks. The attacks are often simple, such as UDP reflection attacks where a single UDP query generates a large response. The source address of the query is forged to be the address of the intended attack victim, and not much more needs to be done. A small query stream can result in a massive attack. UDP protocols such as SNMP, the *Network Time Protocol* (NTP)^[17], the DNS, and *memcached* have been used in the past and doubtless will be used again.

Why can't we fix this problem? We've been trying for decades, and we just can't seem to get ahead of the attacks. Advice to network operators to prevent the leakage of packets with forged source addresses was published nearly two decades ago, in 2000.^[18] Yet massive UDP-based attacks with forged source addresses still persist today. Aged computer systems with known vulnerabilities continue to be connected to the Internet and are readily transformed into attack bots.

The picture of attacks is also becoming more ominous. Although we previously attributed these hostile attacks to "hackers," we quickly realised that a significant component of them had criminal motivations. The progression from criminal actors to state-based actors is also entirely predictable, and we are seeing an escalation of this cyber warfare arena with the investment in various forms of vulnerability exploitation that are considered desirable national capabilities.

It appears that a major problem here is that collectively we are unwilling to make any substantial investment in effective defence or deterrence. The systems that we use on the Internet are overly trusting to the point of irrational credulity. For example, the public key certification system used to secure web-based transactions is repeatedly demonstrated to be entirely untrustworthy, yet that's all we trust. Personal data is continually breached and leaked, yet all we seem to want to do is increase the number and complexity of regulations rather than actually use better tools that would effectively protect users.

The larger picture of hostile attack is not getting any better. Indeed, it's getting very much worse. If any enterprise has a business need to maintain a service that is always available for use, then any form of in-house provisioning is just not enough to withstand attack. These days only a handful of platforms can offer resilient services, and even then it's unclear whether they could withstand the most extreme of attacks.

A constant background level of scanning and probing goes on in the network, and any form of visible vulnerability is ruthlessly exploited. One could describe today's Internet as a toxic wasteland, punctuated with the occasional heavily defended citadel. Those who can afford to locate their services within these citadels enjoy some level of respite from this constant profile of hostile attack, while all others are forced to try to conceal themselves from the worst of this toxic environment, while at the same time aware that they will be completely overwhelmed by any large-scale attack.

It is a sobering thought that about one-half of the world's population are now part of this digital environment. A more sobering thought is that many of today's control systems, such as power generation and distribution, water distribution, and road-traffic-control systems are exposed to the Internet.

Perhaps even more of a worry is the increasing use of the Internet in automated systems that include various life-support functions. The consequences of massive failure of these systems in the face of a sustained and damaging attack cannot be easily imagined.

The Internet of Billions of Tragically Stupid Things

What makes this scenario even more depressing is the portent of the so-called *Internet of Things* (IoT). In those circles where Internet prognostications abound and policy makers flock to hear grand visions of the future, we often hear about the boundless future represented by this Internet of Things.^[19] This phrase encompasses some decades of the computing industry's transition from computers as esoteric pieces of engineering affordable only by nations to mainframes, desktops, laptops, handheld devices, and now wrist computers. Where next? In the vision of the IoT we are going to expand the Internet beyond people and press on using billions of these chattering devices in every aspect of our world.

What do we know about the “things” that are already connected to the Internet?

Some of them are not very good. In fact, some of them are just plain stupid. And this stupidity is toxic, in that their sometime-inadequate models of operation and security affect others in potentially malicious ways. If such devices were constantly inspected and managed, we might see evidence of aberrant behaviour and correct it. But these devices are unmanaged and all but invisible. Examples include the controller for a web camera, the so-called “smart” thing in a smart television, or the controls for anything from a washing machine to a goods locomotive. Nobody is looking after these devices.

When we think of an IoT we think of a world of weather stations, webcams, “smart” cars, personal fitness monitors, and similar things. But what we tend to forget is that all of these devices are built on layers of other people's software that is assembled into a product at the cheapest possible price point. It may be disconcerting to realise that the web camera you just installed has a security model that can be summarised with the phrase: “no security at all,” and it's actually offering a view of your house to the entire Internet. It may be slightly more disconcerting to realise that your electronic wallet is on a device that is using a massive compilation of open source software of largely unknown origin, with a security model that is not completely understood, but appears to be susceptible to be coerced into being a “yes, take all you want” device.

It would be nice to think that we've stopped making mistakes in code, and from now on our software in our things will be perfect. But that's hopelessly idealistic. It's just not going to happen. Software will not be perfect. It will continue to have vulnerabilities.

It would be nice to think that this Internet of Things is shaping up as a market where quality matters, and consumers will select a more expensive product even though its functional behaviour is identical to a cheaper product that has not been robustly tested for basic security flaws. But that too is hopelessly naive.

The Internet of Things will continue to be a marketplace where the compromises between price and quality will continue to push us on to the side of cheap rather than secure. What's going to stop us from further polluting our environment with a huge and diverse collection of programmed unmanaged devices with inbuilt vulnerabilities that will be all too readily exploited? What can we do to make this world of these stupid cheap toxic things less stupid and less toxic? So far we have not found workable answers to this question.

The Next 10 Years

The silicon industry is not going to shut down anytime soon. It will continue to produce chips with more gates, finer tracks, and more stacked layers for some years to come. Our computers will become more capable in terms of the range and complexity of the tasks that they will be able to undertake.

At the same time, we can expect more from our network. Higher capacity certainly, but also greater levels of customisation of the network to our individual needs.

However, I find it extremely challenging to be optimistic about security and trust in the Internet. We have made little progress in this area over the last 10 years, and there is little reason to think that the picture will change in the next 10 years. If we can't fix it, then, sad as it sounds, perhaps we simply need to come to terms with an Internet jammed full of tragically stupid things

However, beyond these broad-brush scenarios, it's hard to predict where the Internet will head. Technology does not follow a predetermined path. It's driven by the vagaries of an enthusiastic consumer marketplace that is readily distracted by colourful bright shiny new objects, and easily bored by what we quickly regard as commonplace.

What can we expect from the Internet in the next 10 years that can outdo a pocket-sized computer that can converse with me in a natural language? That can offer more than immersive 3D video with outstanding quality? That can bring the entire corpus of humanity's written work into a searchable database that can answer any of our questions in mere fractions of a second?

Personally, I have no clue what to expect from the Internet. But no matter what manages to capture our collective attention, I am pretty confident that it will be colourful, bright, shiny, and entirely unexpected!

References and Further Reading

- [0] Vint Cerf, “A Decade of Internet Evolution,” *The Internet Protocol Journal*, Volume 11, No. 2, June 2008.
- [1] Geoff Huston, “A Decade in the Life of the Internet,” *The Internet Protocol Journal*, Volume 11, No. 2, June 2008.
- [2] Geoff Huston, “QoS — Fact or Fiction?” *The Internet Protocol Journal*, Volume 3, No. 1, March 2000.
- [3] Metcalf’s Law:
https://en.wikipedia.org/wiki/Metcalf%27s_law
- [4] T. Sridhar, “Cloud Computing—A Primer: Part One,” *The Internet Protocol Journal*, Volume 12, No. 3, September 2009.
- [5] T. Sridhar, “Cloud Computing—A Primer: Part Two,” *The Internet Protocol Journal*, Volume 12, No. 4, December 2009.
- [6] Geoff Huston, “Anatomy: Inside Network Address Translators,” *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [7] Geoff Huston, “In Defence of NATs,” *The Internet Protocol Journal*, Volume 20, No. 3, November 2017.
- [8] Geoff Huston, “The BGP Routing Table,” *The Internet Protocol Journal*, Volume 4, No. 1, March 2001.
- [9] Dennis Ferguson, Acee Lindem, and John Moy, “OSPF for IPv6,” RFC 5340, July 2008.
- [10] Gary Scott Malkin, “RIP Version 2,” RFC 2453, November 1998.
- [11] Donald Slice, Peter Paluch, Donnie Savage, Russ White, Steven Moore, and James Ng, “Cisco’s Enhanced Interior Gateway Routing Protocol (EIGRP),” RFC 7868, May 2016.
- [12] David Meyer, “The Locator Identifier Separation Protocol (LISP),” *The Internet Protocol Journal*, Volume 11, No. 1, March 2008.
- [13] Juliusz Chroboczek, “The Babel Routing Protocol,” RFC 6126, April 2011.
- [14] Geoff Huston, “The Future for TCP,” *The Internet Protocol Journal*, Volume 3, No. 3, September 2000.
- [15] Daniel McCarney, “Automatic Certificate Management,” *The Internet Protocol Journal*, Volume 20, No. 2, June 2017.

- [16] William Stallings, “SSL: Foundation for Web Security,” *The Internet Protocol Journal*, Volume 1, No. 1, June 1998.
- [17] Geoff Huston, “Network Time Protocol,” *The Internet Protocol Journal*, Volume 15, No. 4, December 2012.
- [18] Paul Ferguson and Daniel Senie, “Network Ingress Filtering: Defeating Denial of Service Attacks which Employ IP Source Address Spoofing,” RFC 2827, May 2000.
- [19] Bob Hinden, “The Internet of Insecure Things,” *The Internet Protocol Journal*, Volume 20, No. 1, March 2017.
- [20] Geoff Huston and Joao Luis Silva Dama, “DNS Privacy,” *The Internet Protocol Journal*, Volume 20, No. 1, March 2017.
- [21] Zi Hu, Liang Zhu, John Heidemann, Allison Mankin, Duane Wessels, and Paul Hoffman, “Specification for DNS over Transport Layer Security (TLS),” RFC 7858, May 2016.

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

Fileless Malware

by David Strom

Malware authors have gotten more clever and sneaky over time to make their code more difficult to detect and prevent. One of the more worrying recent developments goes under the name “fileless.” There is reason to worry because these kinds of attacks can do more damage and the malware can persist on your computers and networks for weeks or months until they are finally neutralized. Let’s talk about what this malware is and how to understand it better so we can try to stop it from entering our networks to begin with.

Usually, the goal of most malware is to leave something behind on one of your endpoints: one or more files that contain an executable program that can damage your computer, corral your PC as part of a botnet, or make copies of sensitive data and move them to an external repository. Over the years, various detection products have gotten better at finding these *residues*, as they are called, and blocking them.

But the malware game is one of “cat-and-mouse,” and as defenders get better at stopping the malware, the malware authors get better at evading these blockades. Back in the early days of the Internet, most blocking routines looked for certain signatures, either as the name of one of the running programs on your computer or specific patterns of behavior across your network. These options worked until the malware authors got better at hiding their signature moves.

This point is where the fileless versions come into play. They aim to leave as little residue as possible, so the detection products can’t easily find them. Or better yet, to do something misleading, or under the guise of something that an uninfected operating system might do.

Actually, the fileless designation is somewhat of a misnomer since there is still something left behind. It may not be a complete executable file or *Dynamic Linked Library* (DLL), but enough of some code is used to actually be able to run some series of processes that can do the “dirty work” of the malware. Starting in 2016, researchers began to see more of these fileless efforts from attackers, and they have continued to become more popular because the malware can be a powerful infection that is neither easily found nor prevented.^[1]

Fileless Attack Types

Fileless malware uses three different attack types: *Return-Oriented Programming*, *Scripting-Based Attacks*, and *Polymorphic Attacks*. Each is somewhat different.

The first type of attack is called *Return-Oriented Programming*, which is the most popular and could be considered the “classic” version. The malware can execute standard DLLs and other sequences of code that can compromise an otherwise uninfected system. This code could also be part of your desktop web browser, or common *Operating System* (OS) tools such as desktop applications. Since the code is already present in these operating system functions, there is no particular “file” actually being run that is unique to the malware itself. Instead, the malware author piggybacks on these routines to get the job done.

To make this kind of attack work, you have to be familiar with the program code that you intend to hijack for evil purposes, and be reasonably assured that the target endpoint is running the particular version of code for that operating system. Small variations in OS versions, such as from Windows 7 to 7.1 or MacOS 10.12.5 to 10.12.6, could foil the attack because the code base has changed. Or if Microsoft or Apple (in particular, given their popularity and installed base) has issued a patch to fix the potential exploit.

Scripting-Based Attacks are the second fileless category. Another avoidance technique is to execute malware using built-in Windows scripting engines such as Microsoft Office, Windows *PowerShell*, or Microsoft’s HTML Application Host. These attacks typically take advantage of process hooking and don’t leave any file-based residues on the endpoint. If your detection systems can’t see the script execution or understand the command-line arguments, you can’t readily figure out that it is malware.

For example, a typical malware script allocates memory, resolves Windows application program interfaces, and downloads some executable directly to the memory of the target PC. After it gets in memory, it starts up a malicious service and begins to use the target to explore the local network to find other targets, often by starting other malicious PowerShell scripts that use privilege escalation and remote execution. That is a lot of stuff going on to avoid detection and to do the dirty work of the malware. But most of these activities can operate “under the radar,” and perhaps be discovered only months later with detailed forensic analysis that can capture the sequence of events that played out for the particular attack.

Scripting attacks are gaining favor, mainly because there is so much built-in software on a typical modern PC that can do most of what a piece of malware needs to do: to access a shared network drive, copy portions of files, set up some sort of monitoring tool, and so forth. Why reinvent the criminal’s wheel when it is sitting on the average desktop or laptop?

A third method is called *Polymorphic Attacks*. These attacks adapt to a variety of conditions, operating systems, and circumstances and try to evade security scans and protection products to infect your endpoints. They are called polymorphic because they shift their signatures, attack methods, and targets so that you can't easily identify and catch them. Attackers typically use polymorphism as just one of many code obfuscation methods to hide from defenders, such as determining if they are running inside a *Virtual Machine* (a favorite ploy researchers use), or encrypting their code to mask their executables.

Over the past several years, security vendors have begun to take this notion of polymorphism that attackers use and turn it into a defensive maneuver. The idea is to make a target Web server or other piece of network infrastructure appear to change frequently so it can't be easily identified or infected. Sometimes this method is called a *moving-target defense*, which could be a synonym or refer to some aspect of the defense that changes the nature of your applications or code locations. These vendors include Morphisec, Shape Security, and Polyverse, all startup companies. One startup, CyActive, was successful enough to be purchased by PayPal.

Polymorphic defenses can limit the amount of time a potential attacker can invade a network, since their target system appears to move around the network or change properties.

Researchers are seeing combinations of all three types of attacks to make them even more sophisticated and difficult to track down. Sometimes, malware authors program multiple attack types to ensure that something will evade your defenses and penetrate your network. As I said earlier, it is a game of cat and mouse.

Fileless Samples

Let's look at a few recent examples of fileless malware to illustrate the differences and how fileless malware has evolved over time.

Back in 2014, the retailer Target experienced a now-infamous breach. It turns out malware was placed on its network through a very simple strategy: someone's network access credentials were discovered and copied, in this case belonging to an employee of Target's heating vendor. What is noteworthy about this attack is its simplicity, and the fact that Target's network was a flat topology with no *virtual LANs* (vLANs) or other segments. This example is a reminder that bad network practice can make any kind of malware—fileless or otherwise—dangerous. Brian Krebs studied what went wrong and reported on this attack in his blog.^[2]

As most of us know by now, back in 2016, the *Democratic National Committee* (DNC) was hacked. The attack used a fileless malware product that took advantage of both PowerShell and *Windows Management Instrumentation* (WMI) in order to “get a foot into the door” of the political party's systems.

WMI is commonly used for day-to-day management tasks such as deploying automation scripts, running a process at a given time, or getting information about the installed applications or hardware.

The DNC malware also used PowerShell as a staging tool to execute other scripts to compromise a system. It used WMI to install backdoors that allow persistence by enabling the adversary to launch malicious code automatically, after a specified period of system uptime or according to a specific schedule. Again, the malware used all of these tactics to avoid detection.^[3]

August Stealer was discovered at the end of 2016 and was attributed to the TA530 criminal group. Targeted at customer service and call center staffs, it used infected Word macros and PowerShell scripts that were delivered via phished e-mails. The e-mails were designed to look like queries from users over support issues and used various subject lines such as the following:

- Erroneous charges from [recipient's domain]
- [recipient's domain] - Help: Items vanish from the cart before checkout
- [recipient's domain] Support: Products disappear from the cart during checkout
- Need help with order on [recipient's domain]
- Duplicate charges on [recipient's domain]

August contains stealing functionality targeting credentials, cryptocurrency wallets, and sensitive documents from the infected computers.^[4]

Discovered earlier in 2017, *Duqu2* is a good example of the first fileless malware products. The malware was found in more than 140 enterprise networks of banks, government offices, and telecom companies across 40 different countries.

It takes the form of a malicious PowerShell script and a series of the following Windows Registry values that at the time were unique to the malware and used to identify the infected systems:

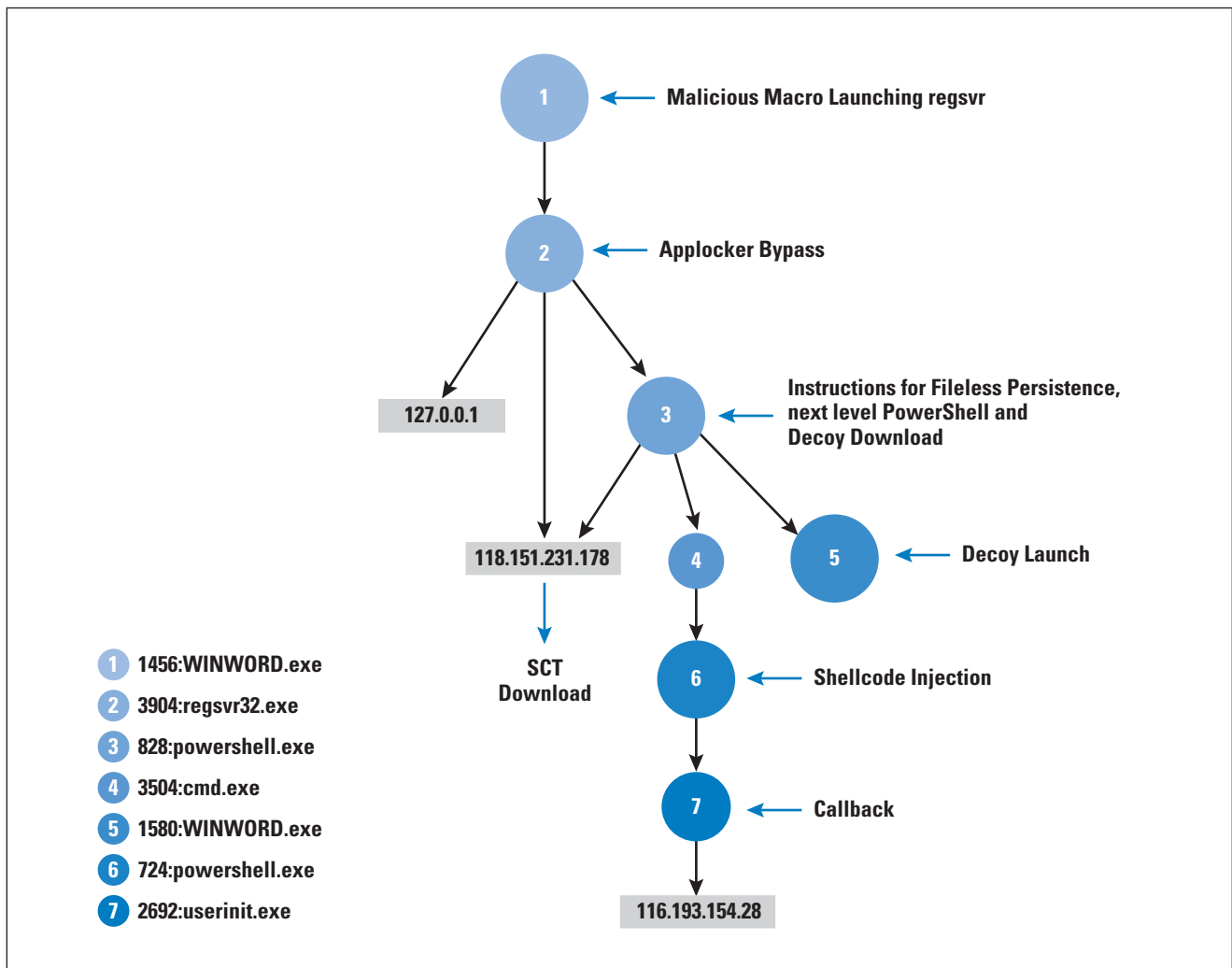
- `HKLM\SYSTEM\ControlSet001\services\` – path is modified after using the SC utility
- `HKLM\SYSTEM\ControlSet001\services\PortProxy\v4tov4\tcp` – path is modified after using the NETSH utility

After finding its way onto the target hard drives, it then starts up via a malicious Windows installer or MSI file, which then deletes itself and renames various files to hide its operations. After the malware is installed on a PC, it just runs in the memory of the PC.

“That’s why memory forensics is critical to the analysis of malware and its functions. In fact, detection of this attack would be possible in memory, network, and Windows Registry only,” says one group of researchers from Kaspersky Labs that studied its operations.^[5] Obviously, running in memory means the Duqu2 malware won’t last after the PC is rebooted—one drawback of many fileless products.

Poison Ivy, also discovered earlier this year, is an example of fileless malware that was used on a specific target, in this case Mongolian government officials. It takes the form of a malicious Microsoft Word macro. If the target has enabled macros—which is a typical setting for most users—it runs and creates a remote-access connection to log keystrokes and capture screens and videos from the PC. All these actions are done from memory-resident programs taking advantage of certain PowerShell command sequences. Figure 1 shows its various modules (courtesy of FireEye).

Figure 1: *Poison Ivy* Modules

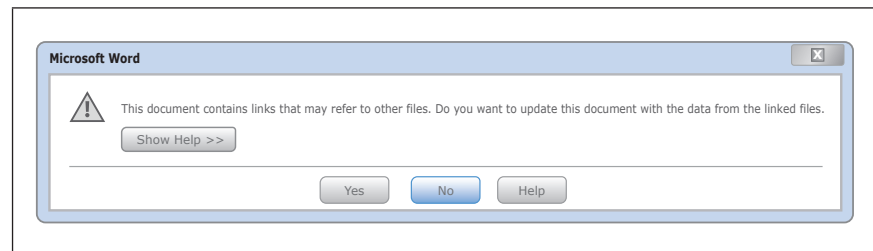


Poison Ivy also tried to evade detection by Microsoft's *AppLocker* protection system by inserting a reference to itself in AppLocker's whitelisted applications using a series of Windows programs and scripts. It also created a series of decoy documents to make its operations seem benign to the infected user. As you can see, this software is very complex, with several different stages and methods to find its way into a user's PC. It has also been used in other circumstances besides the Mongolian case.^[6]

Other targeted fileless campaigns, such as the *OilRig* malware^[7], have been attributed to Iranian state-sponsored actors. This campaign targeted 250 e-mail accounts of various Israelis, including ironically cybersecurity researchers at Ben Gurion University. While Microsoft released a patch back in April 2017 that prevents this malware from spreading, many enterprises haven't yet applied it. Ironically (again), the malware authors used the details from a published proof-of-concept to design their tool accordingly.

This particular malware used an infected Word document that was sent as an attachment and used to steal information from targeted PCs. It used a specialized fileless version of the *Helminth Trojan* malware. Earlier versions of OilRig used infected macros, but this attack used an embedded Web link using an **.HTA** executable file. This type of file is automatically run by the Windows program **MSHTA.EXE** (for Microsoft HTML applications). Normally, when this program runs an **.HTA** file, it displays the following warning message:

Figure 2: This warning about file permissions from Windows is only briefly seen by users when they click on an infected file.



However, this malware anticipates this situation, and automatically sends an “Enter” command so that the warning window is quickly dispatched and the malware does its business. Other targeted malware campaigns include one targeting American restaurant computers using the *Fin7* malware^[8]. In the past, this malware targeted banks and government financial filing documents. And like other fileless attacks, it hides inside a Microsoft Word document that is attached to phishing e-mails. One new twist with the Fin7 restaurant attacks is that it executes various attacks completely in memory, without using any PowerShell commands.

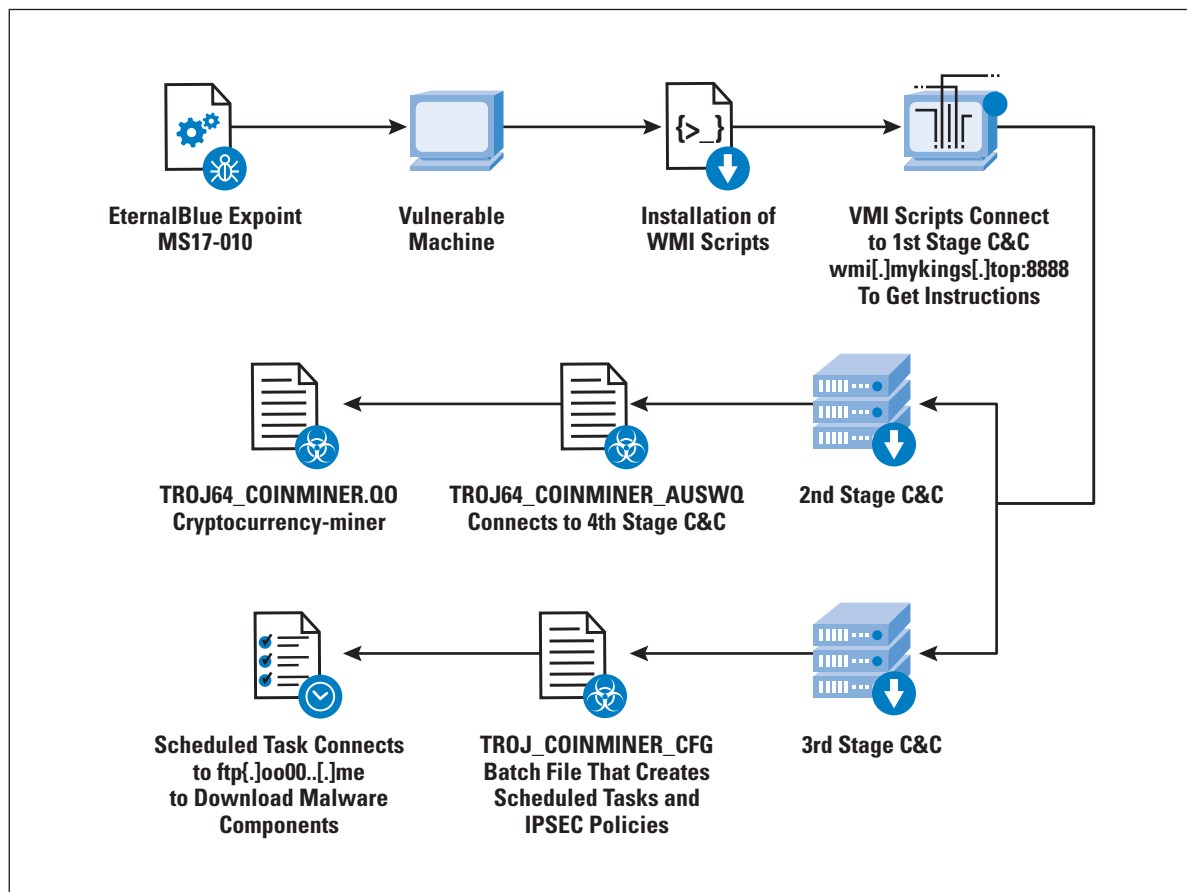
Another obfuscation technique goes by the name of *DoubleAgent*. It takes advantage of an undocumented feature in Microsoft Application Verifier. This verifier is code that has been around since at least Windows XP and is a Windows feature that lets developers do runtime verifications of their applications for finding and fixing security problems.

Unfortunately, it has an undocumented feature that security researchers from Cybellum discovered. The feature gives attackers a way to replace the legitimate verifier with a rogue one so they can gain complete control of the application. Cybellum said, “DoubleAgent gives the attacker the ability to control the AV and perform all the operations above without being detected, while keeping the illusion that the AV is working normally.”^[9]

Security vendors have recently issued patches to correct this flaw, but again this example demonstrates that malware writers are getting better at finding these sorts of hidden mechanisms to avoid discovery and being blocked.

In July 2017, a new fileless malware was discovered called *CoinMiner*^[10], which was found mainly in Japan and Indonesia. The purpose of this malware was to create a hidden bitcoin mining application, to generate cryptocurrency for the attacker. It uses WMI to persist beyond reboots and execute a series of scripts. CoinMiner invades a PC through the *EternalBlue* exploit^[11], which is the same method that was used by the *WannaCry* worm. Figure 3 is a diagram of its logic flow (courtesy of Trend Micro):

Figure 3: CoinMiner Logic Flow



Common Prevention Steps

Given the scope of these exploits, here are a few steps to take to prevent infections across your network:

- Apply patches quickly and across all systems. Microsoft issues regular patches for Windows, and the other operating system vendors do the same for their systems. Don't delay an update, because you can see some criminals take advantage of unpatched systems with their malware. The EternalBlue exploit is a good example: the patch to prevent this attack was available for more than a month before this exploit was launched.
- Segment your network carefully and make sure you understand access rights, especially of third parties.
- Restrict administrator rights to the minimum number of systems. Many of the WMI-based exploits count on profligate use of admin rights that aren't needed.
- Disable Windows programs that aren't needed, such as WMI, PowerShell, and support for ancient protocols such as *Server Message Block* (SMB) v1.
- Whitelist applications to further restrict what can run on most endpoints.

Conclusion

As you can see, the bad guys have gotten better at plying their trade, and through the use of fileless techniques, they are making their malware harder to detect and protect. Hopefully, by learning about some of these past examples, you can tune your own defenses accordingly and do a better job of keeping them infection-free.

References

- [1] Ericka Chickowski, "Fileless Malware Takes 2016 By Storm," DarkReading, December 2016.
<http://www.darkreading.com/vulnerabilities---threats/fileless-malware-takes-2016-by-storm/d/d-id/1327796>
- [2] "Target Hackers Broke in Via HVAC Company," Krebs on Security blog, February 2014.
<https://krebsonsecurity.com/2014/02/target-hackers-broke-in-via-hvac-company/>
- [3] "DNC Hack Exhibits One of 3 Attack Trends To Watch for in 2017," Crowdstrike blog, January 2017.
<https://www.crowdstrike.com/blog/dnc-hack-exhibits-one-of-3-attack-trends-to-watch-for-in-2017/>

- [4] “August in November: New Information Stealer Hits the Scene,” Proofpoint, December 2016.
<https://www.proofpoint.com/us/threat-insight/post/august-in-december-new-information-stealer-hits-the-scene>
- [5] “Fileless attacks against enterprise networks,” Securelist, February 2017.
<https://securelist.com/fileless-attacks-against-enterprise-networks/77403/>
- [6] “Poison Ivy: Assessing Damage and Extracting Intelligence,” FireEye Special Report, 2014.
<https://www.fireeye.com/content/dam/fireeye-www/global/en/current-threats/pdfs/rpt-poison-ivy.pdf>
- [7] Michael Gorelik, “Iranian Fileless Attack Infiltrates Israeli Organizations,” Morphisec Cyber Security Blog, April 2017.
<http://blog.morphisec.com/iranian-fileless-cyberattack-on-israel-word-vulnerability>
- [8] Michael Gorelik, “FIN7 Takes Another Bite at the Restaurant Industry,” Morphisec Cyber Security Blog, June 2017.
<http://blog.morphisec.com/fin7-attacks-restaurant-industry>
- [9] Michael Engstler, “DoubleAgent: Zero-Day Code Injection and Persistence Technique,” Cybellum blog, March 2017.
<https://cybellum.com/doubleagentzero-day-code-injection-and-persistence-technique/>
- [10] Buddy Tancio, “Cryptocurrency Miner Uses WMI and EternalBlue To Spread Filelessly,” TrendLabs Security Intelligence Blog, August 2017.
<http://blog.trendmicro.com/trendlabs-security-intelligence/cryptocurrency-miner-uses-wmi-eternalblue-spread-filelessly/>
- [11] “EternalBlue,” Wikipedia article,
<https://en.wikipedia.org/wiki/EternalBlue>

DAVID STROM has written for *The Internet Protocol Journal* before on e-mail topics; he runs the *Inside.com* Security e-mail newsletter. He was the founding editor-in-chief of *Network Computing* (USA) magazine and is the co-author of the 1998 book, *The Internet Message: Closing the Book with Electronic Mail*, with Marshall T. Rose. E-mail: david@strom.com

Postel Award Presented to Steven G. Huter

The Internet Society, a global non-profit dedicated to ensuring the open development, evolution, and use of the Internet, recently presented the prestigious *Jonathan B. Postel Service Award* to Steven G. Huter, Director for the *Network Startup Resource Center* (NSRC) and a Research Associate at the University of Oregon. For decades he has worked with people around the world to strengthen the infrastructure, partnerships, and expertise upon which the Internet has been developed in more than 120 countries, particularly in support of research and education.



© Stonehouse
Photographic/Internet Society

“Steve Huter is the quintessential candidate for the Postel Award. For a quarter of a century, Steve has enabled hundreds of institutions to build and operate new components of the Internet. His dedication to this task mirrors Postel’s own and continues to this day. Literally millions have benefited from Steve’s work,” explains Vint Cerf, founding president of the Internet Society.

Mr. Huter was selected by an international award committee comprised of former Postel Award winners. The committee placed particular emphasis on candidates who have supported and enabled others in addition to their own contributions. The award is being presented to Mr. Huter in recognition of “his leadership and personal contributions at the Network Startup Resource Center that enabled countless others to develop the Internet in more than 120 countries.” The NSRC was formally begun in 1992 by Randy Bush and John Klensin with funding from a U.S. National Science Foundation grant to provide technical assistance to people setting up networks in developing areas to support scientific collaboration.

“Steve epitomizes the values and spirit of the Postel Award. For more than twenty-five years he has energetically brought the fruits of the Internet to developing countries using his unique combination of a multicultural background, technical knowledge, unfailing energy and commitment,” adds Steve Crocker, CEO and co-founder of Shinkuro, Inc.

Mr. Huter joined the NSRC in 1993, where he has led the development and implementation of programs that provide technical training, equipment, and expertise across Africa, Asia-Pacific, Latin America-Caribbean, and the Middle East.

“It is a tremendous honor to be acknowledged for helping to advance Jon’s vision and philosophy of developing the Internet into a global resource,” said Mr. Huter on receiving the award.

“The most important thing I learned from Jon Postel and the founders of the NSRC is to cultivate a culture of network operators who help each other via technical exchange and resource sharing; this is an effective way to empower more network engineers and enable continuous progress for a community of peers in all regions of the world. Thank you to the NSRC team and all who have contributed over the years towards achieving this objective and enriching the Internet.”

The Postel Award was established by the Internet Society to honor individuals or organizations that, like Jon Postel, have made outstanding contributions to the data communications community. The award is focused on sustained and substantial technical contributions, service to the community, and leadership. Kathy Brown, President and CEO of the Internet Society presented the award including a US\$20,000 honorarium and a crystal engraved globe, during the 102nd meeting of the *Internet Engineering Task Force* (IETF) held in Montreal, Canada, July 14–20, 2018.

Founded by Internet pioneers, the *Internet Society* (ISOC) is a non-profit organization dedicated to ensuring the open development, evolution and use of the Internet. Working through a global community of chapters and members, the Internet Society collaborates with a broad range of groups to promote the technologies that keep the Internet safe and secure, and advocates for policies that enable universal access. The Internet Society is also the organizational home of the IETF.

The NSRC, which is based at the University of Oregon, was established in 1992 to provide technical assistance to organizations setting up computer networks in new areas to connect scientists engaged in collaborative research and education. For the past few decades, the NSRC has helped develop Internet infrastructure and network operations communities in Africa, Asia-Pacific, Latin America-Caribbean, and the Middle East. The NSRC is partially funded by the *International Research Network Connections* (IRNC) program of the U.S. National Science Foundation and Google, with additional contributions from dozens of public and private organizations.

Check your Subscription Details!

If you have a print subscription to this journal, you will find an expiration date printed on the back cover. For the last couple of years, we have “auto-renewed” your subscription. Now we ask that you log in and perform this simple task yourself. The subscription portal is here: <https://www.ipjsubscription.org/> This process will ensure that we have your current contact information, as well as delivery preference (print edition or download). For any questions, contact us by e-mail at: ipj@protocoljournal.org

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Fabrizio Accatino	Narelle Clark	Chris Gamboni	Edward Jennings
Scott Aitken	Steve Corbató	Xosé Bravo Garcia	Aart Jochem
Jacobus Akkerhuis	Brian Courtney	Kevin Gee	Richard Johnson
Antonio Cuñat Alario	Dave Crocker	John Gilbert	Jim Johnston
Matteo D'Ambrosio	Kevin Croes	Serge Van Ginderachter	Jonatan Jonasson
Jens Andersson	John Curran	Greg Goddard	Daniel Jones
Danish Ansari	André Danthine	Octavio Alfageme	Gary Jones
Tim Armstrong	Morgan Davis	Gorostiaga	Amar Joshi
Richard Artes	Jeff Day	Barry Greene	Merike Kaeo
David Atkins	Freek Dijkstra	Martijn Groenleer	Andrew Kaiser
Jaime Badua	Geert Van Dijk	Geert Jan de Groot	Christos Karayiannis
Hidde Beumer	David Dillow	Gulf Coast Shots	David Kekar
John Bigrow	Richard Dodsworth	Sheryll de Guzman	Jithin Kesavan
Axel Boeger	Ernesto Doelling	James Hamilton	Jubal Kessler
Gerry Boudreaux	Eugene Doroniuk	Stephen Hanna	Shan Ali Khan
L de Braal	Karlheinz Dölger	John Handin	Nabeel Khatri
Kevin Breit	Joshua Dreier	Martin Hannigan	Anthony Klopp
Thomas Bridge	Lutz Drink	John Hardin	Henry Kluge
Ilia Bromberg	Andrew Dul	David Harper	Michael Kluk
Christophe Brun	Holger Durer	Edward Hauser	Andrew Koch
Gareth Bryan	Peter Robert Egli	David Hauweele	Carsten Koempe
Caner Budakoglu	George Ehlers	Marilyn Hay	Alexader Kogan
Stefan Buckmann	Peter Eisses	Headcrafts SRLS	Antonin Kral
Scott Burleigh	Torbjörn Eklöv	Hidde van der Heide	Mathias Körber
Jon Harald Bøvre	ERNW GmbH	Johan Helsingius	John Kristoff
Olivier Cahagne	ESdatCo	Robert Hinden	Terje Krogdahl
Antoine Camerlo	Steve Esquivel	Alain Van Hoof	Bobby Krupczak
Tracy Camp	Mikhail Evstiounin	Edward Hotard	Murray Kucherauw
Fabio Caneparo	Paul Ferguson	Bill Huber	Warren Kumari
Roberto Canonico	Kent Fichtner	Hagen Hultzs	Darrell Lack
David Cardwell	Gary Ford	Kevin Iddles	Yan Landriault
John Cavanaugh	Jean-Pierre Forcioli	Mika Ilvesmaki	Markus Langenmair
Lj Cemer	Christopher Forsyth	Karsten Iwen	Fred Langham
Dave Chapman	Craig Fox	David Jaffe	Richard Lamb
Stefanos Charchalak	Fausto Franceschini	Ashford Jaggernauth	Tracy LaQuey Parker
Greg Chisholm	Tomislav Futivic	Jozef Janitor	Simon Leinen
Marcin Cieslak	Edward Gallagher	John Jarvis	Robert Lewis
Brad Clark	Andrew Gallo	Dennis Jennings	Martin Lillep

Sergio Loreti	Mazdak Rajabi Nasab	Boudhayan Roychowdhury	Adrian Stevens
Guillermo a Loyola	Krishna Natarajan	Carlos Rubio	Clinton Stevens
Hannes Lubich	Darryl Newman	Timo Rüter	John Streck
Dan Lynch	Paul Nikolic	RustedMusic	Viktor Sudakov
Miroslav Madić	Marijana Novakovic	Babak Saberi	Edward-W. Suor
Alexis Madriz	David Oates	George Sadowsky	Vincent Surillo
Carl Malamud	Ovidiu Obersterescu	Scott Sandefur	T2Group
Michael Malik	Mike O'Connor	Sachin Sapkal	Roman Tarasov
Yogesh Mangar	Mike O'Dell	Arturas Satkovskis	David Theese
Bill Manning	Carlos Astor Araujo	Phil Scarr	Douglas Thompson
Harold March	Palmeira	Elizabeth Scheid	Rey Tucker
Vincent Marchand	Alexis Panagopoulos	Jeroen Van Ingen Schenau	Sandro Tumini
David Martin	Gaurav Panwar	Carsten Scherb	Phil Tweedie
Jim Martin	Manuel Uruena Pascual	Dan Schrenk	Steve Ulrich
Timothy Martin	Ricardo Patara	Richard Schultz	Unitek Engineering
Gabriel Marroquin	Dipesh Patel	Roger Schwartz	AG
Carles Mateu	Alex Parkinson	SeenThere	John Urbanek
Juan Jose Marin Martinez	Craig Partridge	Scott Seifel	Martin Urwaleck
Ioan Maxim	Dan Paynter	Yury Shefer	Betsy Vanderpool
David Mazel	Leif-Eric Pedersen	Yaron Sheffer	Surendran
Miles McCredie	Juan Pena	Doron Shikmoni	Vangadasalam
Brian McCullough	Chris Perkins	Tj Shumway	Buddy Venne
Joe McEachern	David Phelan	Jeffrey Sicuranza	Alejandro Vennera
Jay McMaster	Derrell Piper	Thorsten Sideboard	Luca Ventura
Mark Mc Nicholas	Rob Pirnie	Andrew Simmons	Tom Vest
Carsten Melberg	Jorge Ivan Pincay Ponce	Pradeep Singh	Dario Vitali
Kevin Menezes	Victoria Poncini	Henry Sinnreich	Randy Watts
Bart Jan Menkveld	Blahoslav Popela	Geoff Sisson	Andrew Webster
William Mills	Tim Pozar	Helge Skrivervik	Tim Weil
Desiree Miloshevic	David Raistrick	Darren Sleeth	Jd Wegner
Thomas Mino	Priyan R Rajeevan	Bob Smith	Rick Wesson
Wijnand Modderman	Paul Rathbone	Mark Smith	Peter Whimp
Mohammad Moghaddas	Bill Reid	Job Snijders	Jurrien Wijlhuizen
Charles Monson	Rodrigo Ribeiro	Ronald Solano	Pindar Wong
Andrea Montefusco	Justin Richards	Asit Som	Romeo Zwart
Fernando Montenegro	Mark Risinger	Ignacio Soto Campos	Bernd Zeimetz
Joel Moore	Ron Rockrohr	Peter Spekrijse	廖明沂.
Soenke Mumm	Carlos Rodrigues	Thayumanavan Sridhar	
Tariq Mustafa	Lex Van Roon	Ralf Stempfer	
Stuart Nadin	William Ross	Matthew Stenberg	



Follow us on Twitter and Facebook

@protocoljournal



<https://www.facebook.com/newipj>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Internet
Society



Diamond Sponsors



Ruby Sponsors

Your logo here!

Sapphire Sponsors

Your logo here!

Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

ADDRESS SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

David Conrad, Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

March 2019

Volume 22, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

F R O M T H E E D I T O R

In This Issue

From the Editor	1
A Quick Look at QUIC	2
Missing Checksums	13
Fragments	21
Thank You	24
Letters to the Editor	26
Supporters and Sponsors	27

The *Transmission Control Protocol* (TCP) is a core component of the Internet Protocol Suite. TCP has proven robust and flexible in the face of changing network infrastructures, but may not be the most efficient way to retrieve the many components of today's complex web pages. The *Quick UDP Internet Connection* (QUIC) protocol is an alternative to TCP for web traffic. QUIC was initially developed and deployed by Google and is now being standardized in the *Internet Engineering Task Force* (IETF). In our first article, Geoff Huston examines the motivations for QUIC and describes the protocol and its implementation.

According to *Wikipedia*: “A *checksum* is a small-sized datum derived from a block of digital data for the purpose of detecting errors that may have been introduced during its transmission or storage. It is usually applied to an installation file after it is received from the download server. By themselves, checksums are often used to verify data integrity but are not relied upon to verify data authenticity.” In preparation for the “rolling” of the root *Key Signing Key* of the *Domain Name System* (DNS), tests were developed to create so-called *key-tags*. This key-tag generation process “...became an adventure in itself that included beautiful discrete math, flawed functions, carefully crafted primes, multiple cryptographic libraries, and some brilliant people,” according to Roy Arends, author of our second article, “The Quest for the Missing Checksums.” IPJ doesn't normally delve into complex mathematics, but in this case the interplay of various software libraries and methods provides some valuable lessons for anyone involved in code generation and testing.

We would like to remind you that this journal depends on the generous support of numerous individuals and organizations. If you would like to help support IPJ, please contact us for further details. Comments, suggestions, book reviews, and articles are always welcome. Send your messages to ipj@protocoljournal.org

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

A Quick Look at QUIC

by Geoff Huston, APNIC

Quick UDP Internet Connection (QUIC) is a network protocol initially developed and deployed by Google, and is now being standardized in the *Internet Engineering Task Force* (IETF). In this article we'll take a quick tour of QUIC, looking at the goals that influenced its design, and the implications QUIC might have on the overall architecture of the Internet Protocol Stack.

QUIC is not exactly a recent protocol, as the concept appears to have been developed by Google in 2012, and initial public releases of this protocol were included in Chromium version 29, released in August 2013. QUIC is one of many transport-layer network protocols that attempt to refine the basic operation of the *Transmission Control Protocol* (TCP).

Why are we even thinking about refining TCP?

TCP is now used in billions of devices and is perhaps the most widely adopted network transport protocol that we've witnessed so far. If this protocol weren't fit for our use, then we would have moved on and adopted some other protocol or protocols instead. Part of the reason for the broad adoption of TCP is its incredible flexibility. The protocol can support a diverse variety of uses, from micro-exchanges to gigabyte data movement, transmission speeds that vary from hundreds of bits per second to tens and possibly hundreds of gigabits per second. TCP is the workhorse of the Internet. But even so, there is room for refinement. TCP is used in many different ways, and its design represents a set of trade-offs that attempt to be a reasonable fit for many purposes but not necessarily an ideal fit for any particular one.

One of the aspects of the original design of the Internet Protocol Suite was that of elegant brevity and simplicity. The specification of TCP^[1] is not a single profile of behavior that has been cast into a fixed form that was chiseled into the granite slab of a rigid standard. TCP is malleable in many important ways. Numerous efforts over the years have shown that it is possible to stay within the standard definition of TCP, in that all the packets in a session use the standard TCP header fields in mostly conventional ways, but also to create TCP implementations that behave radically differently from each other. Critically, the TCP standard does not strictly define how the sender can control the amount of data in flight across the network. There is a convention to adopt an approach of slowly increasing the amount of data in flight while there are no visible errors in the data transfer (as shown by the stream of received acknowledgement [ACK] packets) and quickly responding to signals of network congestion (packet drop, as shown by duplicate acknowledgements) by rapidly decreasing the sending rate.

Variants of TCP use different controls to manage this “slow increase” and “rapid drop” behavior^[2] and may also use different signals to control this data flow. These signals include measurements of end-to-end delay, or inter-packet jitter (such as the recently published *Bottleneck Bandwidth and Round-trip Propagation Time* (BBR) protocol^[3]). All of these variants still manage to fit with the broad parameters of what is conventionally called TCP.

It is also useful to understand that most variants of TCP need to be implemented only on the data sender (the “server” in a client/server environment). The common assumption of all TCP implementations is that clients will send a TCP ACK packet on successful receipt of both in-sequence and out-of-sequence data. It is left to the server’s TCP engine to determine how the received ACK stream will be applied to its internal model of network capability and how it will modify its subsequent sending rate accordingly. The implication is that deployment of new variants of TCP flow control is essentially based on deployment within service-delivery platforms and does not necessarily imply changing the TCP implementations in all the billions of clients. This feature also contributes to the flexibility of TCP.

But despite its considerable flexibility, TCP has its problems, particularly with web-based services. These days most web pages are not simple monolithic objects. They typically contain many separate components, including images, scripts, customized frames, and others. Each of these is a separate web “object,” and if you are using a browser that is equipped with the original implementation of the *HyperText Transfer Protocol* (HTTP) each object will be loaded in a new TCP session, even if the objects are served from the same IP address. The overheads of setting up both a new TCP session and a new *Transport Layer Security* (TLS)^[4] session for each distinct web object within a compound web resource can become quite significant, and the temptation to reuse an already established TLS session is close to overwhelming. But this approach of multiplexing a number of data streams within a single TCP session also has issues. Multiplexing multiple logical data flows across a single session can generate unwanted interdependencies between the flow processors and generate *Head of Line Blocking* situations. It appears that while it makes some logical sense to share a single end-to-end security association and a rate-controlled data-flow state across a network across multiple logical data flows, TCP represents a rather poor way of achieving this outcome. The conclusion is that if we want to improve the efficiency of such compound transactions by introducing parallel behaviors into the protocol, we need to look beyond TCP.

Why not just start afresh and define a new transport protocol that addresses these shortcomings of TCP? The answer is simple: *Network Address Translators* (NATs)!

NATs and Transport Protocols

The original design of IP allowed for a clear separation between the network element that allowed the network to accept an IP packet and forward it onto its intended destination (the “Internet” part of the IP protocol suite) and the end-to-end transport protocol that enabled two applications to communication via some form of “session.” The transport protocol field in the IPv4 packet header and the Next header field of the IPv6 packet header uses an 8-bit field to identify the end-to-end protocol. This design assumed that the network had no need to “understand” what end-to-end protocol was being used within a packet. Ideally an IP packet switch will not differentiate in its treatment of packets depending on the inner end-to-end protocol.

Some 140 protocols are listed in the IP protocol field registry^[5]. TCP and the *User Datagram Protocol* (UDP) are just two of these protocols (protocol values 6 and 17, respectively). In theory at any rate, there is room for a least 100 more. However, in the public Internet the story is somewhat different. TCP and UDP are widely accepted protocols, and the *Internet Control Message Protocol* (ICMP) (protocol 2) is generally accepted, but little else. How did this situation happen?

NATs changed the assumption about network devices not looking inside the packet (to be precise, port-translating NATs changed that assumption). NATs are network devices that look inside the IP packet and re-write the port addresses used by TCP and UDP^[6]. What if an IP packet contains an end-to-end transport protocol identifier value that is neither TCP nor UDP? Most NATs will simply drop the packet, on the basis of a security paradigm that “what you don’t recognize is likely to be harmful.” The pragmatic result is that NATs have limited the choice of transport protocols of an application in the public Internet to just two: TCP and UDP.

If the aim is to deploy a new transport protocol—but not confuse active network elements that are expecting to see a conventional TCP or UDP header—then how can we achieve this goal?

This question was the challenge of the QUIC developers.

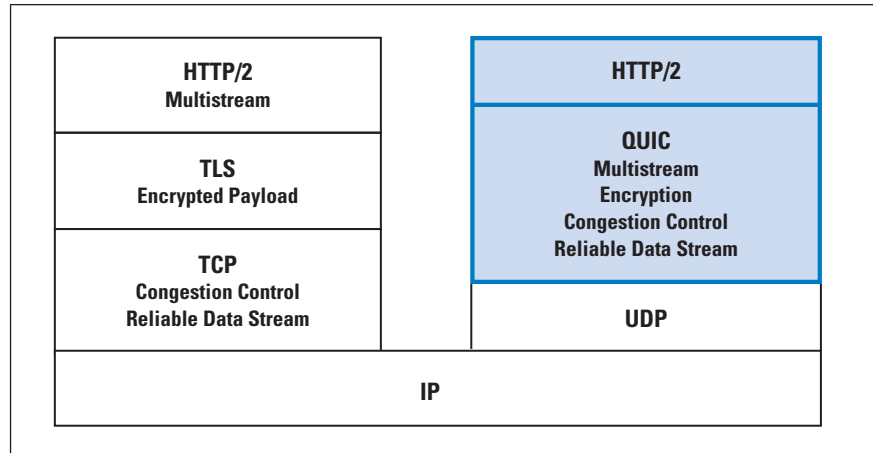
QUIC over UDP

The solution that QUIC chose was a UDP-based approach. UDP is a minimal framing protocol that allows an application to access the basic datagram services that IP offers. Apart from the source and destination port numbers, the UDP header adds a length header and a checksum that covers the UDP header and UDP payload. It is essentially an abstraction of the underlying datagram IP model with just enough additional information to allow an IP protocol stack to direct an incoming packet to an application that has bound itself to a nominated UDP port address. If TCP is an overlay across the underlying IP datagram service, then it’s a small step to think about layering TCP as a payload within a UDP packet.

Using our standard Internet model, QUIC is—strictly speaking—a datagram transport application. An application that uses the QUIC protocol sends and receives packets using UDP port 443.

Technically, this change is very small to an IP packet, adding just 8 bytes to the IP packet by placing a UDP header between the IP and TCP packet headers (Figure 1). The implications of this change are far more significant than these 8 bytes would suggest. However, before we consider these implications, let’s look at some QUIC services.

Figure 1: The QUIC Protocol Architecture



QUIC and the Connection ID

If the choice of UDP as the visible end-to-end protocol for QUIC was a choice dictated by the inflexibility of the base of deployed NAT devices in the public Internet and their collective inability to accommodate new protocols, the way that NATs handle UDP packets has further implications for QUIC.

NATs maintain a *translation table*. In the most general model, a NAT takes the 5-tuple of incoming packets, using the destination and source IP addresses, the destination and source port addresses, and the protocol field, and performs a lookup into the table to find the associated translated fields. The address headers of the packet are rewritten to these new values, *checksums* are recomputed, and the packet is passed onward. Certain NAT implementations may use variants of this model. For example, some NATs use only the source IP address and port address on *outbound* packets as the lookup key, and the corresponding destination IP address and port address in *incoming* packets.

Typically, the NAT generates a new translation table entry when a triggering packet is passed from the *inside* to the *outside* and subsequently removes the table entry when the NAT assumes that the translation is no longer needed. For TCP sessions it is possible to maintain this translation table quite accurately.

New translation-table entries are created in response to *outbound* TCP SYN connection establishment packets and removed either when the NAT sees the TCP FIN exchange or in response to a TCP RST packet or when the session is idle for an extended period.

UDP packets do not have these clear packet exchanges to start and stop sessions, so NATs need to make some assumptions. Most NATs create a new translation table entry when they see an outbound UDP packet that has not matched any existing translation table. The entry is then maintained for some period of time (as determined by the NAT) and is then removed if there are no further packets that match the session signature. Even when there are further matching UDP packets, the NAT may use an overall UDP session timer and remove the NAT entry after some predetermined time interval.

For QUIC and NATs, this situation is a potential problem. The QUIC session is established between a QUIC server on UDP port 443 and the NAT-generated source address and port. However, at some point in the session lifetime the NAT may drop the translation-table entry, and the next outbound client packet will generate a new translation-table entry that may use a different source address and port. How can the QUIC server recognize that this next-received packet, with its new source address and source port number, is actually part of an existing QUIC session?

QUIC uses the concept of *Connection Identifiers* (Connection IDs). Each endpoint generates connection IDs that will allow received packets with that connection ID to be routed to the process that is using that connection ID. During QUIC version negotiation these connection IDs are exchanged, and thereafter each sent QUIC packet includes the current connection ID of the remote party.

This form of semantic distinction between the identity of a connection to an endpoint and the current IP address and port number that QUIC uses is similar to the *Host Identity Protocol* (HIP)^[7]. This protocol also uses a constant endpoint identifier that allows a session to survive changes in the endpoint IP addresses and ports.

QUIC Streams

TCP provides the abstraction of a reliable order byte stream to applications. QUIC provides a similar abstraction to the application, termed within QUIC as *streams*. The essential difference here is that TCP implements a single behavior, while a single QUIC session can support multiple streams profiles.

Bidirectional streams place the client and server transactions into a matched context, as is required for the conventional request/response transactions of HTTP/1. A client would be expected to open a bidirectional stream with a server and then issue a request in a stream which would generate a matching response from the server. It is possible for a server to initiate a bidirectional *push stream* to a client, which contains a response without an initial request.

Control information is supported using unidirectional *control streams*, where one side can pass a message to the other as soon as they are able. An underlying *unidirectional stream* interface, used to support control streams, is also exposed to the application.

Not only can QUIC support many different stream profiles, it can also support different stream profiles within a single end-to-end QUIC session. This concept is not a novel one, of course, and the HTTP/2 protocol is a good example of an application-level protocol adding multiplexing and stream framing in order to carry multiple data flows across a single transport data stream. However, a single TCP transport stream as used by HTTP/2 may encounter *Head of Line Blocking* where all overlay data streams fate-share across a single TCP session. If one of the streams stalls, all overlay data streams could be affected and could stall as well.

QUIC allows for a slightly different form of multiplexing where each overlay data stream can use its own end-to-end flow state, and a pause in one overlay stream does not imply that any other simultaneous stream is affected.

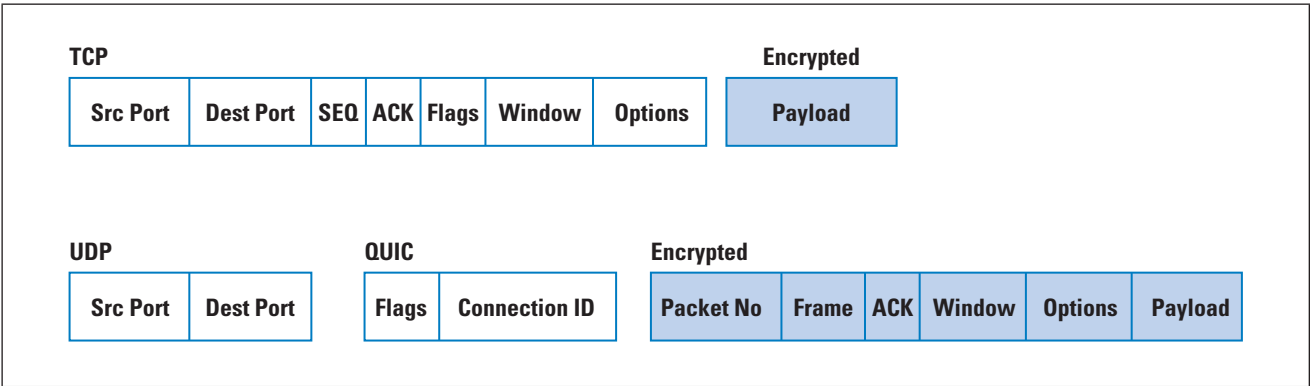
Part of the reason to multiplex multiple data flows between the same two endpoints in HTTP/2 was to reduce the overhead of setting up a TLS security association for each TCP session. This overhead can be quite significant when the individual streams are each sending a small object, and it's possible to encounter a situation where the TCP and TLS handshake component of a compound web object fetch dominates both the total download time and the data volume.

QUIC pushes the security association to the end-to-end state that is implemented as a UDP data flow, so that streams can be started in a very lightweight manner because they essentially reuse the established secure session state.

QUIC Encryption

As is probably clear from the references to TLS already, QUIC uses *end-to-end encryption*. This encryption is performed on the UDP payload, so once the TLS handshake is complete very little of the subsequent QUIC packet exchange is in the clear (Figure 2).

Figure 2: Comparison of TCP and TLS with QUIC



What is exposed in QUIC are the *public flags*. This initial part of a QUIC packet consists of the connection ID, which allows the receiver to associate the packet with an endpoint without decrypting the entire packet. The QUIC version is also part of the public flag set, which is used in the initial QUIC session establishment and can be omitted thereafter.

The remainder of the QUIC packet includes *private flags* and the payload. They are encrypted and are not directly visible to an eavesdropper. This private section includes the packet sequence number. This field is used to detect duplicate and missing packets. It also includes all the flow-control parameters, including window advertisements.

This encryption is one of the critical differences between TCP and QUIC. With TCP the control parts of the protocol are in the clear, so that a network element would be able to inspect the port addresses (and infer the application type), as well as the flow state of the connection. Connection of a sequence of such TCP packets, even if only looking at the packets flowing in one direction within the connection, would allow the network element to infer the round-trip time and the data-transmission rate. And, like a NAT, manipulation of the receive window in the ACK stream would allow a network element to apply a throttle to a connection and reduce the transfer rate in a manner that would be invisible to both endpoints. Placing all of this control information inside the encrypted part of the QUIC packet ensures that no network element has direct visibility to this information, and no network element can manipulate the connection flow.

One could take the view that QUIC enforces a perspective that was assumed in the 1980s: that the end-to-end transport protocol is not shared with the network. All the network “sees” are stateless datagrams, and the endpoints can safely assume that the information contained in the end-to-end transport control fields is carried over the network in a manner that protects it from third-party inspection and alteration.

QUIC and IP Fragmentation

The short answer is “no!” QUIC packets cannot be fragmented^[7, 8]. The way this feature is achieved is by having the QUIC HELLO packet be padded out to the maximal packet size, and not completing the initial HELLO exchange if the maximally sized packet is fragmented.

For IPv4 the maximum QUIC packet size is 1,350 bytes. Adding 8 bytes for the UDP header, 20 bytes for IPv4, and 14 bytes for the Ethernet frame means that a QUIC packet on Ethernet totals 1,392 bytes. There is no particular rationale for this choice of 1,350 other than the results of empirical testing on the public Internet.

For IPv6 the QUIC maximum packet size is reduced by 20 bytes to 1,330. The resultant Ethernet packet is still 1,392 bytes because of the larger IPv6 IP packet header.

What happens if the network path has a smaller *Maximum Transmission Unit* (MTU) than this value? The answer is in the next section.

QUIC and TCP

QUIC is not intended as a replacement for TCP. Indeed, QUIC relies on the continued availability of TCP.

Whenever QUIC encounters a fatal error—such as fragmentation of the QUIC HELLO packet—the intended response from QUIC is to shut down the connection. Since QUIC itself lies in the application space, not the kernel space, the client-side application can be directly informed of this closure of the QUIC connection and it can re-open a connection to the server using a conventional TCP transport protocol.

The implication is that QUIC does not necessarily have to have a robust response for all forms of behavior, and when QUIC encounters a state where it has no clear definition of the desired behavior, it is always an option to signal a QUIC failure to the application. The failure need not be fatal to the application, because such a signal can trigger the application to repeat the transaction using a conventional TCP session.

I can QUIC, do you?

Unlike all other TCP services that use a dedicated TCP port address to distinguish themselves from all other services, QUIC does not advertise itself in such a manner. That reality leaves numerous ways in which a server could potentially advertise itself as being accessible over QUIC.

One such possible path is the use of *Domain Name System* (DNS) *Service Records* (SRV)^[9]. The SRV record can indicate the connection point for a named service using the name of the transport protocol and the protocol-specific service address. This usage may be an option for the future, but no such DNS service record has been defined for QUIC.

Instead, in keeping with the overall QUIC approach of loading up most of the service functionality into the application itself, a server that supports QUIC can signal its capability within HTTP itself. The way it signals is defined in an Internet standard for “Alternative Services”^[10], which is a means to list alternative ways to access the same resources.

For example, the Google homepage, www.google.com, includes the HTTP header:

```
alt-svc: quic=":443"; ma=2592000; v="44,43,39"
```

This entry indicates that the same material is accessible using QUIC over port 443. The “**ma**” field is the time to keep this information on the local client, which in this case is 30 days, and the “**v**” field indicates that the server will negotiate QUIC versions 39, 43, and 44.

QUIC Lessons

QUIC is a rather forceful assertion that the Internet infrastructure is now heavily ossified and more highly constrained than ever. There is no room left for new transport protocols in today’s network. If what you want to do can’t be achieved within TCP, then all that’s left is UDP.

The IP approach to packet-size adaptation through fragmentation was a powerful concept once upon a time. A sender did not need to be aware of the constraints that may apply on a path. Any network-level packet fragmentation and reassembly was invisible to the end-to-end packet transfer. This invisibility is no longer wise. Senders need to ensure that their packets can reach their intended destinations without any additional requirement for fragmentation handling.

Mutual trust is over. Applications no longer trust other applications. They don’t trust the platform that hosts them or the shared libraries that implement essential functions. Applications no longer trust a network to keep their secrets. More and more functions and services are being pulled back into the application and are passed out from an application as much as possible in packets that are cloaked in a privacy shroud.

There is a tension between speed, security, and paranoia. An ideal outcome is one that is faster, private, and secure. Where it is not obvious and the inevitable trade-offs emerge, it seems that we have some minimum security and privacy requirements that simply must be achieved. But once we have achieved these minimum requirements, we are then happy to trade off incremental improvements in privacy and security for better session performance.

The traditional protocol-stack model was a convenient abstraction, not a design rule. Applications do not necessarily need to bind to transport-layer sockets provided by the underlying platform. Applications can implement their own end-to-end transport if necessary.

The infrastructure of the Internet might be heavily ossified, but the application space is seeing a new set of possibilities open up. Applications need not wait for the platform to include support for a particular transport protocol or await the deployment of a support library to support a particular name-resolution function. Applications can solve these issues for themselves directly. The gain in flexibility and agility is considerable.

There is a price to pay for this new-found agility, and that price is broad interoperability. Browsers that support QUIC can open up UDP connections to certain servers and run QUIC, but browsers cannot assume—as they do with TCP—that QUIC is a universal and interoperable lingua franca of the Internet. While QUIC is a fascinating adaptation with some very novel concepts, it is still an optional adaptation. For those clients and servers that do not support QUIC, or for network paths where UDP port 443 is not supported, the common fallback is TCP. The expansion of the Internet is inevitably accompanied by inertial bloat, and as we’ve seen with the extended saga of IPv6 deployment, it is a formidable expectation to think that the entire Internet will embrace a new technical innovation in a timeframe of months, years, or possibly even decades! That does not mean that we can’t think new thoughts, and that we can’t realize these new ideas into new services on the Internet. We certainly can, and QUIC is an eloquent demonstration of exactly how to craft innovation into a rather stolid and resistant underlying space.

Further Reading

QUIC has excited considerable interest over the past couple of years, and there are many posts to be found on the ‘net. Here’s a small sample of this online material that you may find to be of interest:

- A useful consideration of positive and negative aspects of QUIC are in Robin Marx’s post “QUIC and HTTP/3: Too big to fail?”
<https://calendar.perfplanet.com/2018/quic-and-http-3-too-big-to-fail/>
- A slightly older (2014) but useful technical overview of QUIC can be found in Shigeki Ohtsu’s presentation to the HTTP/2 Conference Japan.
https://www.slideshare.net/shigeki_ohtsu/quic-overview
- A commentary on Cloudflare’s investigations with QUIC can be found in a recent blog post: “The Road to QUIC”:
<https://blog.cloudflare.com/the-road-to-quic/>
- A discussion of QUIC work in the IETF by Mark Nottingham, QUIC Working Group Co-Chair: “What’s Happening with QUIC,”
<https://www.ietf.org/blog/whats-happening-quic/>

References

- [1] Jon Postel, “Transmission Control Protocol,” RFC 793, September 1981.
- [2] Geoff Huston, “Faster,” *The ISP Column*, June 2005.
<https://www.potaroo.net/ispcol/2005-06/faster.html>
- [3] Neal Cardwell, Yuchuing Cheng, C. Stephen Gunn, Soheil Hasses Yeganeh, and Van Jacobson, “BBR: congestion-based congestion control,” *Communications of the ACM*, Vol. 60, Issue 2, pp 58–66, February 2017.

- [4] Eric Rescorla, “The Transport Layer Security (TLS) Protocol Version 1.34,” RFC 8446, August 2018.
- [5] IANA Protocol Numbers Registry.
<https://www.iana.org/assignments/protocol-numbers/protocol-numbers.xhtml>
- [6] Geoff Huston, “Anatomy: A Look Inside Network Address Translators,” *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [7] Geoff Huston, “Fragmentation,” *The Internet Protocol Journal*, Volume 19, No. 2, June 2016.
- [8] Geoff Huston, “IPv6 and Packet Fragmentation,” *The Internet Protocol Journal*, Volume 21, No. 1, April 2018.
- [9] Arnt Gulbrandsen, Paul Vixie, and Levon Esibov, “A DNS RR for specifying the location of services (DNS SRV),” RFC 2782, February 2000.
- [10] Mark Nottingham, Patrick McManus, and Julian Reschke, “HTTP Alternative Services,” RFC 7838, April 2016.

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

The Quest for the Missing Checksums

by Roy Arends, ICANN

The *Domain Name System* (DNS) is a hierarchical namespace that provides a method to look up Internet identifiers such as IP addresses using easy-to-remember domain names. This hierarchy starts at the root,^[0] where the actual namespace is delegated to several registries. The data at the root is signed with cryptographic keys, using *Domain Name System Security Extensions* (DNSSEC)^[1, 2, 3]. These cryptographic keys are replaced over time.

In an effort to change the top cryptographic key for the DNS, the so-called root *Key Signing Key*^[4], several testbeds were created to emulate the process in a lab environment. In those testbeds, the actual root DNS keys are not used since the testbed operators do not have control of the private keys; rather keys of the same size using the same cryptographic algorithms and functions are generated. Apart from the fact that the key material is different, this emulated root zone cannot be distinguished from the real root zone.

This effort to generate certain cryptographic keys became an adventure in itself that included beautiful discrete math, flawed functions, carefully crafted primes, multiple cryptographic libraries, and some brilliant people.

The result of this effort shows that using an ancient checksum function to identify cryptographic keys is not optimal.

The problem

DNSSEC protects the DNS. To be precise, it protects validating resolvers' caches. DNSSEC uses cryptographic keys to validate signatures, and these signatures contain a *key-tag* that helps to identify which key to use. This key-tag is merely a hint; it doesn't have to be collision-free, and the function to generate it is similar to an IP header checksum (the difference between the two functions is that the key-tag function does not include a final end-around carry).

Technically, a key-tag is a 16-bit unsigned value. For our testbed, to clearly identify which keys were introduced in what year, the idea was to generate some vanity key-tags with the year in them; that is, "2010" for a key that was introduced in 2010, and "2015" for a key introduced in 2015. One way to generate those key-tags is to simply generate all possible key-tags in order to pick the desired ones. This process can be done by repeatedly generating a single key. Since the key-tag is based on the contents of the key, and since the contents of the key contain a lot of random bits, it was assumed that the resulting key-tag would be as random as the key.

After the process to generate keys ran long enough, the expectation was to have 65,536 keys—one for each tag. Surprisingly, it was possible to generate only 16,387 keys with unique tags, even after generating millions of keys. Specifically, the key-tags “2010” and “2015” were not included. It turns out that key-tag “2015” was excluded for a different reason than why key-tag “2010” was excluded!

Is it the software?

In order to track down this non-intuitive result, suspicion first fell on the software used to generate the keys. The *BIND* software package from *Internet Systems Consortium, Inc.* (ISC) has a command-line tool named *dnssec-keygen*. The convention it uses is to embed the key-tag in the filename. When a new key is generated, *dnssec-keygen* checks to determine if a key with a certain tag already exists to avoid overwriting it.

The *Flags* field in a DNSSEC key influences the value of the key-tag. For instance, if a key is revoked in the future, the “REVOKE” flag is set and that changes the value of the key-tag. To make sure that a new key-tag doesn’t collide with any existing key, *dnssec-keygen* checks if a new key-tag (and its revoked equivalent) matches an existing key-tag (and its revoked equivalent as well). Initially, it was thought that this key-tag collision check was the culprit.

Since those vanity key-tags were still desired, and since revoked equivalents of keys with the 2010 and 2015 key-tag would not collide with any existing key-tags, it was decided to try to work around this specific check.

One way to avoid this check is to simply use another tool. The LDNS library from NLNetLabs comes with a set of examples. One of these examples is a utility named *ldns-keygen*, which produces DNSSEC keys and does not have the key-tag collision check to protect against accidentally overwriting an existing key. However, after generating millions of keys again, it too generated about 16,384 keys.

The two software tools used have no authors in common, but they do share a cryptographic library: *OpenSSL*. Both pieces of software independently had the limitation of producing only a subset of all possible key-tags. Both used a well-known, widely used cryptographic library. At this discovery the worrying started. If it is the library, and the tags are not distributed evenly, is the quality of the entropy in question? Does the library have any bugs?

To make sure this anomaly was not user error, different versions of *OpenSSL* were tested. Additionally, different entropy sources were used, and lastly, different key sizes were tried. Still, the same number of key-tags was generated.

Is it the library?

The folks on DNS-OARC's operations list came to the rescue. Peter van Dijk from PowerDNS used the PowerDNS management tool: *pdnsutil add-zone-key*, and was able to generate 32,769 unique key-tags. More key-tags than before, but still only about 50% of all possibilities. The tools in PowerDNS, BIND, and LDNS do not share any code or any authors. All three tools were written "from scratch." Additionally, PowerDNS does not use OpenSSL at all; rather it uses *mbedTLS*, a different cryptographic library. That means a problem related solely to the cryptographic libraries or the tools can be ruled out. There was still the observation that *pdnsutil* was able to produce twice as many key-tags as the other tools, but we'll get to that later.

Is it the checksum algorithm?

The next step was testing the key-tag function in RFC 4034^[5]. The key-tag function is very similar to the radix-minus-one complement function for the *Internet Header Checksum*—a radix-minus-one complement function. Note that it is not exactly the same, but the minor difference could not fundamentally reduce the possible number of key-tags.

To test this possibility, a loop was created that fed random numbers into the key-tag algorithm. When using 2,048-bit random numbers as the input (instead of cryptographic keys), all possible key-tags could be produced in a short amount of time. This experiment ruled out that the limiting part was the key-tag algorithm itself. However, we'll come back to that later as well.

Is it purely a math problem?

Meanwhile, Florian Maury and Jérôme Plût from ANSSI took a good look at the problem and discovered it was none of the possibilities mentioned previously. It turns out that an interesting combination of the properties of the Internet Header Checksum and RSA moduli rules out certain results.

The input to the Internet Header Checksum function is treated as blocks of 16 bits and the output is a 16-bit checksum. Radix-minus-one complement methods are as old as accounting itself. The nine's complement method (where the radix is base 10) was used in Pascal's calculator. The method of complements is a technique used to subtract one number from another using only addition of positive numbers. We're not using the complements part here, only the part where we add, with carry, a bunch of bits.

A description of the Internet Header Checksum function follows: Add the 16-bit values with end-around carry; that is, if adding two 16-bit values results in a carry, then add that carry bit to the result of the addition.

Following is the end-around-carry part of the checksum function:

```
($sum AND 65535) + ($sum >> 16)
```

What Jérôme Plût observed is that this expression can be reduced to:

```
$sum mod 65535
```

Since modular arithmetic has the addition property, we can also deduce:

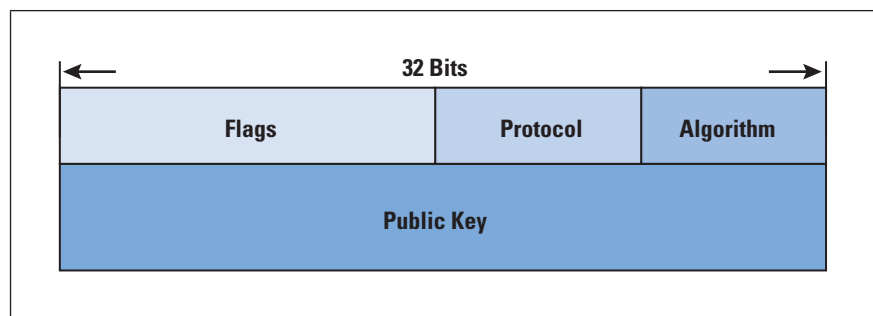
```
($Value1 + $Value2) mod 65535
```

or:

```
($value1 mod 65535) + ($value2 mod 65535)
```

Calculating a key-tag

As said earlier, the Internet Header Checksum is very similar to the key-tag function. The input for this key-tag algorithm is the RDATA part of a DNSKEY record:



For all keys generated in this exercise, all the fields remain the same, except for the modulus in the *Public Key* field.

For a *Key Signing Key*, the value of the *Flags* field is 257, and the *Protocol* field always has the value of 3. The *Algorithm* field has the value 8 (RSASHA256)^[7]. With those parameters, the *Public Key* field consists of an *Exponent* and a *Modulus*. For this exercise, the exponent has value 65537 and is preceded with an *Exponent Length* field (value 3).

The constant part of this input can now simply be added up as a series of 16-bit unsigned values:

```
$value1 = Flags + Protocol*256 + Algorithm + ExpLen*256 + Exponent
$value1 = 257 + 3*256 + 8 + 3*256 + 65537
$value1 = 67338
```

Using the deduction from before:

$$\text{keytag} = (\text{value1} \bmod 65535) + (\text{value2} \bmod 65535)$$

$$\text{keytag} = (67338 \bmod 65535) + (\text{value2} \bmod 65535)$$

$$\text{keytag} = 1803 + (\text{value2} \bmod 65535)$$

The part of the checksum that is not constant is the *RSA-modulus*. The RSA-modulus is a composite number with two very large prime factors. In the previous equation, value2 is the RSA-modulus. The last substitution becomes:

$$\text{keytag} = 1803 + (\text{RSA-modulus} \bmod 65535)$$

Since the value 1803 is constant, it has no influence on the number of possible key-tags, hence the solution to the reduced set of possible key-tags may be found in the RSA-modulus modulo 65535 part of the equation.

Number theory

What Jérôme Plût observed is that the value 65535 is a composite number with four prime factors: 3, 5, 17, and 257. Since the RSA-modulus and 65535 do not share any factors, the RSA-modulus can't be congruent with 0 modulo 65535.

Therefore, the modulus is not congruent with 0 modulo 3, 0 modulo 5, 0 modulo 17, or 0 modulo 257.

All other congruence values are possible, so the set of possible values is simply a combination of the possible values:

$$2 * 4 * 16 * 256 = 32768.$$

We can now check if we indeed can't have 2010 as a value:

Before, we noted that:

$$\text{keytag} = 1803 + (\text{RSA-modulus} \bmod 65535)$$

We can now substitute the key-tag with our desired value:

$$2010 = 1803 + (\text{RSA-modulus} \bmod 65535)$$

$$2010 - 1803 = \text{RSA-modulus} \bmod 65535$$

$$207 = \text{RSA-modulus} \bmod 65535$$

However, 207 is congruent with 0 modulo 3, meaning that in order for 207 to be possible, the RSA-modulus must have 3 as a factor. We know this is not the case, so 2010 (that is, $207 + 1803$) can't be a key-tag.

The remainder of the problem

Remember that the first exercise led to 16,387 key-tags, not 32,768 as predicted before, or 32,769 as found by Peter van Dijk. Additionally, 32,769 is not 32,768 (and 16,387 is not 16,384, half of the 32,768 space).

32,769 is not 32,768

The key-tag function is similar to the Internet Header Checksum, but not the same. The crucial difference is the last end-around carry.

The last part of the key-tag function is defined in RFC 4034, and reads as follows:

```
ac += (ac >> 16) & 0xFFFF;
return ac & 0xFFFF;
```

The first line adds the carry bits to the accumulator. As a result, the accumulator might be a value larger than fits in a 16-bit value. Instead of again adding the carry bits to the value, it ignores those.

Ignoring the carry bits can, in some cases, result in an off by one value, compared to the Internet Header Checksum. With the Internet Header Checksum, only 32,768 values are possible, as we've seen in the previous section. Since the key-tag function might be off by one, a few more key-tag values are possible.

16,387 is not 32,769

Why was Peter able to produce about twice as many key-tags? Assuming that the values could have been 16,384 and 32,768 (as explained before), the only remaining difference is the library used.

OpenSSL generates primes that are congruent with 2 modulo 3. The resulting modulus is thus always congruent with 1 modulo 3, since:

```
(2 modulo 3) * (2 modulo 3) =
4 modulo 3 =
1 modulo 3
```

This formula reduces the possible key-tag space from $2 * 4 * 16 * 256$ to $1 * 4 * 16 * 256$, which is 16384.

This reduction is the reason why it was not possible to generate a key-tag with the value 2015. Using the same reduction as before, we can now substitute key-tag with 2015:

```
2015 = 1803 + (RSA-modulus mod 65535)
2015 - 1803 = RSA-modulus mod 65535
212 = RSA-modulus mod 65535
```

However, 212 is congruent with 2 modulo 3. We now know that RSA moduli from OpenSSL are always congruent with 1 modulo 3, so key-tag 2015 is simply not possible when using OpenSSL.

The library that Peter is using, *mbedTLS*, does generate primes that are congruent with 1 modulo 3.

Conclusion

The limited key-tag space does not present a security issue. The key-tag is merely a hint and it is well known that different cryptographic keys may lead to the same key-tag. However, the decision to use a checksum as an identifier is poor at best. A checksum is designed to check if an error exists in data, and not, in general, designed to be an identifier. Additionally, using a function that is nearly identical to the well-known Internet Header Checksum seems to be an error in the design stage.

Acknowledgements

I cannot begin to thank adequately those who helped me to understand and explain the various compounding issues that resulted in the absence of 75% of all possible key-tags. Florian Maury and Jérôme Plût from ANSSI explained the core issue with the Internet Header Checksum over RSA moduli. Without them, I would still be searching in the dark. Peter van Dijk and Bert Hubert of PowerDNS consumed uncountable electrons and brainwaves to reproduce my findings with different tools and libraries. Google's Ben Laurie held my hand while I was drowning in modular arithmetic and brought me ashore. Finally, it was ICANN's David Conrad who made my broken English and various grammar faux pas readable.

References and Further Reading

- [0] Geoff Huston, "The Root of the Domain Name System," *The Internet Protocol Journal*, Volume 20, No. 2, June 2017.
- [1] Miek Gieben, "DNSSEC: The Protocol, Deployment, and a Bit of Development," *The Internet Protocol Journal*, Volume 7, No. 2, June 2004.
- [2] Donald E. Eastlake 3rd, "Domain Name System Security Extensions," RFC 2535, March 1999.
- [3] Scott Rose, Matt Larson, Dan Massey, Rob Austein, and Roy Arends, "DNS Security Introduction and Requirements," RFC 4033, March 2005.
- [4] George Michaelson, Patrick Wallström, Roy Arends, and Geoff Huston, "Rolling Over DNS Keys," *The Internet Protocol Journal*, Volume 13, No. 1, March 2010.
- [5] Scott Rose, Matt Larson, Dan Massey, Rob Austein, and Roy Arends, "Resource Records for the DNS Security Extensions," RFC 4034, March 2005.
- [6] Scott Rose, Matt Larson, Dan Massey, Rob Austein, and Roy Arends, "Protocol Modifications for the DNS Security Extensions," RFC 4035, March 2005.

- [7] Wes Hardaker, “Use of SHA-256 in DNSSEC Delegation Signer (DS) Resource Records (RRs),” RFC 4509, May 2006.
- [8] Olaf Kolkman and Miek Gieben, “DNSSEC Operational Practices, Version 2,” RFC 6781, December 2012.
- [9] Samuel Weiler and David Blacka, “Clarifications and Implementation Notes for DNS Security (DNSSEC),” RFC 6840, February 2013.
- [10] OpenSSL: <https://www.openssl.org/>
- [11] BIND: <https://www.isc.org/downloads/bind/>
- [12] PowerDNS: <https://www.powerdns.com/>
- [13] LDNS: <https://www.nlnetlabs.nl/projects/ldns/about/>
- [14] Paul Hoffman, Andrew Sullivan, and Kazunori Fujiwara, “DNS Terminology,” RFC 8499, January 2019.
- [15] *L’Agence nationale de la sécurité des systèmes d’information* (ANSSI): <https://www.ssi.gouv.fr/>
- [16] *DNS Operations, Analysis, and Research Center* (DNS-OARC): <https://www.dns-oarc.net/>
- [17] Michael StJohns, “Automated Updates of DNS Security (DNSSEC) Trust Anchors,” RFC 5011, September 2007.
- [18] Duane Wessels, Paul Hoffman, and Warren Kumari, “Signaling Trust Anchor Knowledge in DNS Security Extensions (DNSSEC),” RFC 8145, April 2017.

ROY ARENDS serves as a Principal Research Scientist at *The Internet Corporation for Assigned Names and Numbers* (ICANN). Roy is responsible for successfully delivering research projects; undertaking research design, data collection, and analysis; and producing insightful, stimulating reports that expand knowledge related to the system of unique identifiers on the Internet. E-mail: roy.arends@icann.org

New DNS Terminology RFC

A *Request For Comments* (RFC) updating *Domain Name System* (DNS) terminology was recently published^[10], continuing a decades-long IETF practice of publishing documents to help introduce interested readers to protocol topics by going through the most important terms.

The list of topics with terminology documents includes general terminology^[1], *Network Address Translators* (NATs)^[2], *Diffserv*^[3], Internet connectivity^[4], internationalization^[5], and *Internet of Things* (IoT) networks^[6]. Although these documents are not meant to be step-by-step introductions to the topics, they help someone who already has some understanding go deeper into the topic, and often help clarify terms that are often misused in common writing.

There are many dozens of RFCs defining the DNS, so the terminology is often hard to find. Some common terms such as “host name” are not defined in any RFCs; some are defined only by example; worse, some are defined differently in different RFCs. RFC 8499, “DNS Terminology,” was published as an update to an earlier work to address these issues.

This document is the result of long discussions in the *Domain Name System Operations* (DNSOPS) Working Group^[7], where dozens of DNS operators, software developers, and other experts brought up terms to be covered and argued over the current meaning of terms that are more than 30 years old. A common glossary is necessary to operate the DNS, and to continue to develop the DNS, so that people know what each other mean. The Working Group also hoped that the document would be useful to people who used the DNS tangentially, such as developers of other protocols and non-technical people who interact with the DNS in their work.

RFC 8499 is an update to the first DNS terminology document, RFC 7719^[8]. While the first document was being written, the Working Group agreed that some definitions (such as for “domain name”) needed more work, and it was so difficult to get consensus on other terms that they were left out. The new document is much more complete, and contains some common terms not covered in the earlier document, such as “recursive query,” “lame delegation,” and “split DNS.”

Another significant addition to the document is the first definition of a standards-track document of “the global DNS” and “private DNS.” Many people think they know what “the DNS” is but may not have a specific definition for it; these new terms helps get everyone using the same definitions. Overall, nearly 40 terms that are not defined in other RFCs are defined in this document.

Of course, the DNS will continue to evolve, and new terminology may appear. RFC 8499 is stable, but it might be revised a few years down the road to add these new terms.

- [0] Paul Hoffman, Andrew Sullivan, and Kazunori Fujiwara, “DNS Terminology,” RFC 8499, January 2019.
- [1] Gary Scott Malkin, “Internet Users’ Glossary,” RFC 1983, August 1996.
- [2] Matt Holdrege and Pyda Srisuresh, “IP Network Address Translator (NAT) Terminology and Considerations,” RFC 2663, August 1999.
- [3] Dan Grossman, “New Terminology and Clarifications for Diffserv,” RFC 3260, April 2002.
- [4] John C Klensin, “Terminology for Describing Internet Connectivity,” RFC 4084, May 2005.
- [5] Paul Hoffman and John C Klensin, “Terminology Used in Internationalization in the IETF,” RFC 6365, September 2011.
- [6] Carsten Bormann, Ari Keranen, and Mehmet Ersue, “Terminology for Constrained-Node Networks,” RFC 7228, May 2014.
- [7] DNSOPS Working Group:
<https://datatracker.ietf.org/wg/dnsop/charter/>
- [8] Kazunori Fujiwara, Paul Hoffman, and Andrew Sullivan, “DNS Terminology,” RFC 7719, December 2015.

(Source: <https://www.ietf.org/blog/>)

DNS-OARC

The *DNS Operations, Analysis, and Research Center* (DNS-OARC) brings together key operators, implementers, and researchers on a trusted platform so they can coordinate responses to attacks and other concerns, share information and learn together. DNS-OARC has five key functions:

Information Sharing: DNS-OARC provides a trusted, shared platform to allow the DNS operations community to share information and data. Stringent confidentiality requirements and secure communications mean that proprietary information can be shared on a bilateral basis.

Operational Characterization: As Internet traffic levels continue to grow, the demand on root and other key name servers will outgrow the current infrastructure: this year's DDoS attack traffic levels will become next year's steady state load. DNS-OARC measures the performance and load of key name servers and publish statistics on both traffic load and traffic type (including error types).

Workshops: DNS-OARC organizes semi-annual workshops where members and the public are invited to give presentations on timely topics relevant to DNS both operations and research.

Analysis: Leading researchers and developers provide long-term analysis of DNS performance and post-mortems of attacks so that institutional learning occurs. A well-provisioned system allows members to upload traces and logs, and to perform their own analysis.

Tools and Services: As vulnerabilities and DNS problems come to light, DNS-OARC develops publicly available tools and services to assist with highlighting, diagnosing, and remedying such problems.

DNS-OARC participants fall into one or more of the following categories:

- Operators of root, TLD, or large commercial name servers who consume DNS technology and produce DNS services.
- Implementers who produce DNS technology including software, appliances, and network elements such as load balancing hardware
- Researchers whose work has a strong DNS emphasis and who need access to trace and log data about the global DNS under both “normal” and “abnormal” conditions.
- Security Providers whose companies offer products and services that utilize DNS information to improve the security of their customers.

For more information, visit: <https://www.dns-oarc.net/>

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Fabrizio Accatino	Dave Chapman	Christopher Forsyth	Ashford Jaggernauth
Michael Achola	Stefanos Charchalakis	Andrew Fox	Jozef Janitor
Scott Aitken	Greg Chisholm	Craig Fox	John Jarvis
Jacobus Akkerhuis	David Chosrova	Fausto Franceschini	Dennis Jennings
Antonio Cuñat Alario	Marcin Cieslak	Tomislav Futivic	Edward Jennings
Matteo D'Ambrosio	Brad Clark	Edward Gallagher	Aart Jochem
Jens Andersson	Narelle Clark	Andrew Gallo	Richard Johnson
Danish Ansari	Steve Corbató	Chris Gamboni	Jim Johnston
Tim Armstrong	Brian Courtney	Xosé Bravo Garcia	Jonatan Jonasson
Richard Artes	Dave Crocker	Kevin Gee	Daniel Jones
David Atkins	Kevin Croes	John Gilbert	Gary Jones
Jaime Badua	John Curran	Serge Van Ginderachter	Jerry Jones
Hidde Beumer	André Danthine	Greg Goddard	Amar Joshi
Pier Paolo Biagi	Morgan Davis	Octavio Alfageme	Merike Kaeo
John Bigrow	Jeff Day	Gorostiaga	Andrew Kaiser
Axel Boeger	Freek Dijkstra	Barry Greene	Christos Karayiannis
Keith Bogart	Geert Van Dijk	Martijn Groenleer	David Kekar
Mirko Bonadei	David Dillow	Geert Jan de Groot	Jithin Kesavan
Roberto Bonalumi	Richard Dodsworth	Christopher Guemez	Jubal Kessler
Julie Bottorff	Ernesto Doelling	Gulf Coast Shots	Shan Ali Khan
Photography	Eugene Doroniuk	Sheryll de Guzman	Nabeel Khatri
Gerry Boudreaux	Karlheinz Dölger	James Hamilton	Anthony Klopp
L de Braal	Joshua Dreier	Stephen Hanna	Henry Kluge
Kevin Breit	Lutz Drink	Martin Hannigan	Michael Kluk
Thomas Bridge	Andrew Dul	John Hardin	Andrew Koch
Ilia Bromberg	Holger Durer	David Harper	Ia Kochiashvili
Václav Brožík	Mark Eanes	Edward Hauser	Carsten Koempe
Christophe Brun	Peter Robert Egli	David Hauweele	Alexader Kogan
Gareth Bryan	George Ehlers	Marilyn Hay	Antonin Kral
Caner Budakoglu	Peter Eisses	Headcrafts SRLS	Mathias Körber
Stefan Buckmann	Torbjörn Eklöv	Hidde van der Heide	John Kristoff
Scott Burleigh	ERNW GmbH	Johan Helsingius	Terje Krogdahl
Jon Harald Bøvre	ESdatCo	Robert Hinden	Bobby Krupczak
Olivier Cahagne	Steve Esquivel	Asbjorn Hojmark	Murray Kucherauw
Antoine Camerlo	Jay Etchings	Alain Van Hoof	Warren Kumari
Tracy Camp	Mikhail Evstiounin	Edward Hotard	Darrell Lack
Ignacio Soto Campos	Paul Ferguson	Bill Huber	Yan Landriault
Fabio Caneparo	Kent Fichtner	Hagen Hultzs	Markus Langenmair
Roberto Canonico	The Flirble	Kevin Iddles	Fred Langham
David Cardwell	Organisation	Mika Ilvesmaki	Andrew Lamb
John Cavanaugh	Gary Ford	Karsten Iwen	Richard Lamb
Lj Cemerar	Jean-Pierre Forcioli	David Jaffe	Tracy LaQuey Parker

Simon Leinen	Tariq Mustafa	Carlos Rodrigues	Paul Stancik
Robert Lewis	Stuart Nadin	Lex Van Roon	Ralf Stempfner
Martin Lillepuu	Mazdak Rajabi Nasab	William Ross	Matthew Stenberg
Sergio Loreti	Krishna Natarajan	Boudhayan Roychowdhury	Adrian Stevens
Guillermo a Loyola	Darryl Newman	Carlos Rubio	Clinton Stevens
Hannes Lubich	Paul Nikolich	Timo Rüter	John Streck
Dan Lynch	Travis Northrup	RustedMusic	Viktor Sudakov
Miroslav Madić	Marijana Novakovic	Babak Saberi	Edward-W. Suor
Alexis Madriz	David Oates	George Sadowsky	Vincent Surillo
Carl Malamud	Ovidiu Obersterescu	Scott Sandefur	T2Group
Michael Malik	Tim O'Brien	Sachin Sapkal	Roman Tarasov
Yogesh Mangar	Mike O'Connor	Arturas Satkovskis	David Theese
Bill Manning	Mike O'Dell	PS Saunders	Douglas Thompson
Harold March	Jim Oplotnik	John Sayer	Lorin J Thompson
Vincent Marchand	Carlos Astor Araujo	Phil Scarr	Joseph Toste
David Martin	Palmeira	Elizabeth Scheid	Rey Tucker
Jim Martin	Alexis Panagopoulos	Jeroen Van Ingen Schenau	Sandro Tumini
Timothy Martin	Gaurav Panwar	Carsten Scherb	Angelo Turetta
Gabriel Marroquin	Manuel Uruena Pascual	Ernest Schirmer	Phil Tweedie
Carles Mateu	Ricardo Patara	Dan Schrenk	Steve Ulrich
Juan Jose Marin Martinez	Dipesh Patel	Richard Schultz	Unitek Engineering
Ioan Maxim	Alex Parkinson	Roger Schwartz	AG
David Mazel	Craig Partridge	SeenThere	John Urbanek
Miles McCredie	Dan Paynter	Scott Seifel	Martin Urwaleck
Brian McCullough	Leif Eric Pedersen	Yury Shefer	Betsy Vanderpool
Joe McEachern	Juan Pena	Yaron Sheffer	Surendran
Jay McMaster	Chris Perkins	Doron Shikmoni	Vangadasalam
Mark Mc Nicholas	David Phelan	Tj Shumway	Buddy Venne
Carsten Melberg	Derrell Piper	Jeffrey Sicuranza	Alejandro Vennera
Kevin Menezes	Rob Pirnie	Thorsten Sideboard	Luca Ventura
Bart Jan Menkveld	Marc Vives Piza	Andrew Simmons	Tom Vest
William Mills	Jorge Ivan Pincay Ponce	Pradeep Singh	Dario Vitali
David Millsom	Victoria Poncini	Henry Sinnreich	Laurence Walker
Desiree Miloshevic	Blahoslav Popela	Geoff Sisson	Randy Watts
Joost van der Minnen	Eduard Llull Pou	Helge Skrivervik	Andrew Webster
Thomas Mino	Tim Pozar	Darren Sleeth	Tim Weil
Wijnand Modderman	David Raistrick	Bob Smith	Jd Wegner
Mohammad Moghaddas	Priyan R Rajeevan	Courtney Smith	Rick Wesson
Charles Monson	Paul Rathbone	Mark Smith	Peter Whimp
Andrea Montefusco	Bill Reid	Job Snijders	Jurrien Wijnhuizen
Fernando Montenegro	Rodrigo Ribeiro	Ronald Solano	Pindar Wong
Joel Moore	Glenn Ricart	Asit Som	Romeo Zwart
Maurizio Moroni	Justin Richards	Ignacio Soto Campos	Bernd Zeimetz
Brian Mort	Mark Risinger	Peter Spekrijse	廖明沂.
Soenke Mumm	Ron Rockrohr	Thayumanavan Sridhar	



Follow us on Twitter and Facebook

@protocoljournal



<https://www.facebook.com/newipj>

Letters to the Editor

Ole,

Geoff Huston's most recent article on the last 10 years of the Internet is absolutely brilliant (IPJ Volume 21, No. 2, August 2018). As one of the early implementers of our dear Internet, I am of course amazed at its evolution these past decades, and Geoff has more than "kept up"! His ability to summarize quickly and accurately is without peer. Thank you all.

—Dan Lynch
dan@lynch.com

Geoff,

Thank you very much for your article "Another 10 Years" in *The Internet Protocol Journal*. I enjoyed your perspective and your writing style very much. You have a great skill at explaining a great amount of information.

I subscribed to the early *ConneXions*—*The Interoperability Report* and later IPJ. I've been glad to see your articles over the many years.

Sincerely,

—Richard Berke
Richard_Berke@troweprice.com

The author responds:

I really appreciate your kind words, and I am glad you liked the article.

—Geoff Huston
gih@apnic.net

Letters may be edited for clarity. We'd love to hear from you. Send us your feedback via e-mail to ipj@protocoljournal.org

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

Check your Subscription Details!

If you have a print subscription to this journal, you will find an expiration date printed on the back cover. For the last couple of years, we have "auto-renewed" your subscription, but now we ask you to log in to our subscription system and perform this simple task yourself. The subscription portal is here: <https://www.ipjsubscription.org/> This process will ensure that we have your current contact information as well as delivery preference (print edition or download). For any questions, contact us by e-mail at: ipj@protocoljournal.org

Supporters and Sponsors

Supporters



Internet
Society



Diamond Sponsors



Ruby Sponsors

Your logo here!

Sapphire Sponsors

Your logo here!

Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

ADDRESS SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

David Conrad, Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

July 2019

Volume 22, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

FROM THE EDITOR

In This Issue

From the Editor	1
DNS Privacy and the IETF	2
Improving Routing Security..	14
Fragments	22
Thank You!	24
Call for Papers	26
Supporters and Sponsors	27

Security and privacy have received much attention and treatment in this journal over the years. The original ARPANET protocol suite had few if any security features, but over time a great deal of effort has gone into retrofitting existing protocols with security and privacy features, and adding new technologies such as encryption and authentication mechanisms. In this issue we look at two areas of protocol development related to security and privacy.

The *Domain Name System* (DNS) provides a vital function for everything we do on the Internet, namely translating human-friendly names such as **google.com** to machine-friendly numbers such as **17.172.224.47** or **2001:4860:4860::8888**. A typical DNS entry not only contains the IP address for the server you are trying to reach, but also tells you how to send e-mail to that server. If you tried to contact us between May 31st and June 14th using any of our e-mail addresses such as **ipj@protocoljournal.org**, your message did not get delivered or was delayed. This happened because the DNS registrar for **protocoljournal.org** was changed and the corresponding *Mail Exchange* (MX) records were not updated accordingly. We apologize for this glitch; service has now been restored.

The topic of *DNS Privacy*, originally discussed in this journal in our March 2017 issue, has recently sparked considerable debate following the specification and deployment of *DNS over Hypertext Transfer Protocol Secure* (DoH). In our first article, Geoff Huston explains the motivations for DoH and explores its wider implications.

Routing Security is also an important component for a stable and reliable Internet. The *Mutually Agreed Norms for Routing Security* (MANRS) are a set of “best practice” operational agreements as explained in our second article, by Andrei Robachevsky.

We welcome two new sponsors of IPJ: Akamai and PKNIC. Publication of this journal is made possible by the generous support of numerous individuals and organizations. If you would like to help support IPJ, please contact us for further details. Comments, suggestions, book reviews, and articles are always welcome.

Send your messages to **ipj@protocoljournal.org**

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

DNS Privacy and the IETF

by Geoff Huston, APNIC

From time to time the *Internet Engineering Task Force* (IETF) seriously grapples with its role with respect to technology relating to users' privacy. Should the IETF publish standard specifications of technologies that facilitate third-party eavesdropping on communications, or should it refrain from working on such technologies? Should the IETF take a further step and publish standard specifications of technologies that directly impede various forms of third-party eavesdropping on communications? Is a consistent position from the IETF on personal privacy preferred? Or should the IETF be as agnostic as possible and publish protocol specifications based solely on technical coherency and interoperability without particular regard to issues of personal privacy?

These are not new questions for the IETF. Going back some twenty years, the IETF was working on a standardization of a suite of media gateway protocols when the request was raised to make the protocols compliant with the US *Communications Assistance for Law Enforcement Act* (CALEA)^[16]. This situation excited passions both within the IETF and in the broader circle of observers and commentators. The *Electronic Privacy Information Center* (EPIC) communicated to the IETF its position, which resonated with many IETF participants at the time: "We are writing to urge the IETF not to adopt new protocols or modify existing protocols to facilitate eavesdropping. [...] we believe that such a development would harm network security, result in more illegal activities, diminish users' privacy, stifle innovation, and impose significant costs on developers of communications."^[10]. After much angst and debate, the IETF refused to act on this request, and published its position in RFC 2804: "The IETF has decided not to consider requirements for wiretapping as part of the process for creating and maintaining IETF standards."^[11].

To put this situation into some context, the telephone networks that preceded the Internet typically operated under a framework of interception capability, and this capability was a mandatory requirement for licensed service operators for both their voice and data services. For the IETF to place interception capabilities out of scope for their standards work was not only a strong break from an established public carriage function, but it also threw into some confusion how vendors and operators could define an interoperable standard for interception requests. The *European Telecommunications Standards Institute* (ETSI) has evidently filled this gap with a set of standards for lawful interception^[2]. However, this set of standards still presents real issues to both network operators and law enforcement agencies. One interesting approach in the New Zealand networking community was to support the development of a tool called *OpenLI*, an open source implementation of the ETSI protocols^[13] for use by local network operators.

The IETF's position of refusal to standardize surveillance-enabling architecture modifications twenty years ago did not settle the matter then and hasn't settled it now. Code and standard specifications of network protocols do not necessarily usurp our laws, and code, law, and markets are all elements in a political tussle over what ultimately determines social policies and practices.

The time following the CALEA matter saw an uneasy stand-off between the IETF, as the most visible body associated with the Internet code base, and various public bodies wanting to undertake various forms of surveillance on the Internet. The situation changed in response to the revelations in the documents leaked by Edward Snowden in 2013. Snowden's disclosures of mass surveillance by the US *National Security Agency* (NSA)—evidently working in close cooperation with related agencies in Australia, the United Kingdom, and Canada—prompted the IETF to take a very strong public position in RFC 7258: “Pervasive monitoring is a technical attack that should be mitigated in the design of IETF protocols, where possible.”^[3] This position means that the IETF has crossed into the second of the previous questions. Rather than simply refusing to work on interception technologies, as espoused in RFC 2804, this later RFC advocates that the IETF should publish standard specifications of technologies that directly impede third-party eavesdropping on communications.

It's a noble position that the IETF has taken, but it is perhaps a rather unworldly one in the light of subsequent concerns about the extent of corporate activities in this same area, activities that now have their own name: *surveillance capitalism*. The world of the Internet is now a world where surveillance dominates every aspect of its environment. The online market for goods and services is distorted by the presence of “free” products and services that are funded through a back flow of advertising revenue based on a thorough and comprehensive knowledge of individual users, gained only by using thorough and comprehensive surveillance frameworks that target every user.

The Internet is largely dominated, and indeed driven, by surveillance, and pervasive monitoring is a feature of this network, not a bug. Indeed, perhaps the only debate left today is one over the respective merits and risks of surveillance undertaken by private actors and surveillance by state-sponsored actors. The pronouncement of the IETF denouncing state-sponsored surveillance can only generate a wry smile in retrospect. Sadly, pervasive monitoring is what generates the revenue that propels today's Internet, and the IETF is a coerced fellow traveler, despite the occasional bursts of sometimes hysterical rhetoric that attempts to disavow any such relationship. We have come a very long way from this lofty moral stance on personal privacy into a somewhat tawdry and corrupted digital world, where “do no evil!” has become “don't get caught!”

It has been five years since RFC 7258 was published, and the privacy issue refuses to go away. It seems that the IETF is heading into this turgid and complex field of privacy once more, this time because of the *Domain Name System* (DNS).

DNS Privacy

The DNS has always been a fertile field of opportunity for both surveillance and access control. The basic DNS name-resolution protocol has always worked in a totally unencrypted mode, so that queries and responses are available to any party who can see these transactions on the wire. The wire protocol has no authentication, so network actors can intercept DNS queries addressed to any IP address and provide a response in their name, while the querier may be none the wiser that this substitution has occurred. This idea may sound somewhat esoteric, but every Internet transaction starts with a DNS name-resolution query. The DNS is a timely and accurate indicator of everything we do online, and it's an entirely unprotected and open protocol. What a rich environment for a network eavesdropper! Little wonder that many service operators, and many nation states for that matter, use the DNS for all kinds of purposes relating to both surveillance and access control.

The intersection of RFC 7258 and the DNS has generated the topic of *DNS Privacy*, complete with an IETF Working Group and a worthy collection of drafts of ideas of how to improve the privacy aspects of the DNS.

The first steps in this activity were to look at the interaction between *end clients* and their chosen recursive *resolver*. This element is a critical one of the larger picture, because it is the only part of the DNS resolution service where the IP address of the end client is contained in the query. Once the query is passed within the DNS infrastructure, the query contains no direct identifying link to the client.

Client Subnet

As an aside it is worth mentioning the *Client Subnet* extension to DNS queries and the tension between privacy and performance levers that are accessible with such end-user information leakage (RFC 7871)^[4].

The rise of *Content Distribution Networks* (CDNs) and multiple points of presence has led to a technique, commonly used by Akamai today as well as some others, where the assumed geolocation of the DNS resolver posing the question is a reasonable facsimile to the location of the end client. The concurrent rise of the use of open DNS resolvers, most notably the **8.8.8.8** service from Google, negated this assumption.

In response to the frustrations on the CDN side of misdirected users and woefully inefficient content delivery, the IETF standardized a mechanism to attach the subnet of the end client to the query, RFC 7871^[4].

The attachment of the client's credentials was made by using the *Extension Mechanisms for DNS* (EDNS)^[17], and the idea was to put the IP address of the end client making the query (or an IP prefix) into the query that both survived recursive resolver hand-offs and could be used as a distinguishing label in local cache lookups to perform content steering via the DNS.

Semantically a bridge is being crossed here. Previously the DNS could be thought of as an invariant distributed database. No matter who posed a name query, the response was always the same. Client Subnet is an overt admission that some folks want the DNS to be inconstant, such that the value of the response may depend on the identity of the querier. More importantly, a major privacy bridge is also being crossed. Previously, authoritative name servers were not exposed to the identity of the original client making the query, because they were masked by the intermediary recursive resolvers. With Client Subnet, the authoritative server is aware of the original client. Interception and eavesdropping undertaken at the server end will enjoy a richer view of the end clients that are expressing some level of interest in the names served by this authoritative server.

Perhaps in deference to RFC 7258 it should be noted that the IETF appeared to be reluctant to reference it when specifying this Client Subnet extension, but nevertheless the organization ended up doing it! I quote here Section 2 of RFC 7871, which is a good description of the level of compromise and discomfort that lies just beneath the surface of this DNS privacy debate in the IETF:

“If we were just beginning to design this mechanism, and not documenting existing protocol, it is unlikely that we would have done things exactly this way.

The IETF is actively working on enhancing DNS privacy and the reinjection of metadata has been identified as a problematic design pattern.

As noted above however, this document primarily describes existing behavior of a deployed method to further the understanding of the Internet community.

We recommend that the feature be turned off by default in all nameserver software, and that operators only enable it explicitly in those circumstances where it provides a clear benefit for their clients. We also encourage the deployment of means to allow users to make use of the opt-out provided. Finally, we recommend that others avoid techniques that may introduce additional metadata in future work, as it may damage user trust.

Regrettably, support for the opt-out provisions of this specification are currently limited. Only one stub resolver, *getdns*, is known to be able to originate queries with anonymity requested, and as yet no applications are known to be able to indicate that user preference to the stub resolver.”^[4]

DNS over TLS

The *Transport Layer Security* (TLS) protocol can both encrypt the communication between a client and a server and provide some assurance to the client that the server is operated under the authority of the named entity that the client intended to connect to. In much the same manner as TLS is used to protect HTTP sessions and provide some assurance that the service point is an authorized agent of the named service, this protocol can also be used in the DNS context between end users' client stub resolver and their chosen recursive resolver service.

The IETF *DNS PRIVate Exchange* (DPRIVE) Working Group^[12] has worked on *DNS over TLS* (DoT)^[9, 10] and we are now seeing numerous DNS recursive resolver services that support DoT. Resolver code for Unbound, PowerDNS, and Knot exists, and BIND can be configured with TLS use through a *stunnel* configuration. So if you are prepared to set up your own DNS-resolution environment on your device, you can bypass the open DNS-resolution system provided by your ISP and use a DoT service that will hide your DNS queries and responses from your ISP, and any interested onlookers.

However, it has to be said that using DoT constitutes a highly qualified form of privacy. It's not a solution for everyone. Adding DoT support to your platform may require the installation of a third-party app on your device (which may or may not be possible on your device), and in any case the number of users who are willing to alter the DNS configuration of their device is very limited. Even when the packing of the TLS service is quite seamless, such as in Android Pie's DNS privacy option, it probably still will not be broadly used. In Android's case it is an esoteric feature buried a few levels deep in menus, it is not necessarily supported on all Android platforms, and unless you already know about it you will probably never stumble over it when poking around in your device.

But configuring the client is only half the story. Whom are you going to talk to? Which recursive resolvers support client connections using DoT?

It's an important question. While you are stopping others from looking over your shoulder at your online DNS activity, you are still telling your chosen DNS recursive resolver your complete DNS profile. Today, Google, Cloudflare, and Quad9 all provide open DNS resolver servers.

Sharing your secrets with Google may sound a bit like dancing with the devil. Google's advertising platform generates comprehensive user profiles and its ad support systems are certainly expert and capable practitioners of the art of surveillance capitalism!

In their defense, I must note that Google clearly states that it does not use its public DNS service to reap user profile data and it exercises strict controls over access to DNS data, but that itself raises the question of how such unilateral undertakings are enforced within the company. Google does not open itself up for any form of third-party compliance inspection. Although its DNS practice statement is an excellent statement of noble intent, how can a user be assured that Google is thoroughly and completely committed to every detail in the practice statement?

Let's look at it from the user's perspective. When you configure your system to use a third-party open DNS resolver, you may also be leaving aside your local national regulatory framework. It's a mixed package, because you may be circumventing what you might think of as onerous national content controls, including DNS censorship, but at the same time you may also be circumventing any rights and protections you may have under these same national regulatory structures. When you are outside of any national jurisdiction, then who is left to ensure that service providers adhere to their stated practices in providing the service?

It's not just trust in the service provider at the other end of the TLS connection. Even accessing such a privacy-oriented service may present a problem. In its wisdom, the IETF's DPRIVE Working Group standardized DoT over TCP port 853. This port is not port 443 as used by TLS in supporting HTTP. Any network operator can prevent users from applying this DNS overlay service by simply blocking all traffic to TCP port 853.

DNS over TLS represents a specialized service accessible to just a few. It's a service that is readily blocked. It's a service that may prevent surveillance on the wire, but still ends up sharing your DNS activity with the DoT service provider of your choice. You may well still be compromised in terms of assured privacy protection, but does it make you feel better having a choice as to which service operator you choose to expose yourself to?

DNS over HTTPS

What caused all the current fuss in the IETF was a variant of this DoT approach, termed *DNS over Hypertext Transfer Protocol Secure* (DoH).

In terms of the carriage of DNS on the wire there is almost nothing that differs between DoT and DoH. Both take wire format DNS messages, encrypt them using TLS, and use a TCP session between the client and resolver. In protocol terms of packets on the wire the only difference between the two approaches is that DoH uses the same TCP port number as HTTPS, namely port 443. It may sound like a cosmetic change, but two very fundamental differences transcend this simple protocol tweak.

Firstly, DoH is very difficult to detect. It looks like HTTPS traffic and uses the same port as HTTPS traffic. One could make assumptions in the opening TLS handshake where the name of the server is carried in the clear, but work on encrypted *Server Name Indication* (SNI) in TLS 1.3 is proceeding, and it is reasonable to believe that even this small aperture of visibility will be sealed up in the near future. If you also add TLS padding to the mix, then even traffic profile analysis would not necessarily reveal that it is a DNS session within the TLS stream.

If privacy is the goal, then what's to complain about with this picture? Surely DoH offers the end user a package of encryption, mimicry, and obfuscation that hides the DNS to all but the endpoints of the session.

The answer to this question leads to the second fundamental difference between DoT and DoH. We are no longer talking about an esoteric feature knob that requires a knowledgeable, or even fool-hardy, user to turn it on. The DNS session may look like just another HTTPS session to the network, but it also looks like just another HTTPS session to the host platform. In other words, a browser may just turn on DoH all by itself. It's not the user turning it on, nor the platform turning it on, but the browser itself. No special configuration needs to be in place by the platform of the local network to support the operation of DoH. If a browser chooses to use DoH, then there is little that the platform or the network can do to prevent it. If a browser has installed DoH support, then control over the DNS name-resolution function has passed from the user to the browser provider, and rather than being an esoteric function enabled by a handful of users, it becomes a "mainstream" service used by potentially billions of end users. For example, it appears that Google's Chrome browser enjoys a 60% market share of browsers^[5]. If Chrome enabled DoH by default, then what would that mean for the entire DNS? Would it literally disappear from sight?

The second concern is the choice of DoH server. Instead of using a locally configured DNS-resolver service provided by the ISP, DoH switches the situation to use a service configured by the browser. The early implementation of this service in Firefox requires the explicit configuration of a trusted recursive resolver, in a manner similar to the configuration of the DoT server in Android Pie. What if the DoH resolver is configured by the browser by default?

Let's just pause for a second to think about this notion. DoH can place the control of the privacy setting for DNS queries into the hands of the browser, bypassing both the user and the local internet infrastructure, and can do so in a way that intertwines secure web services with secure DNS service. In privacy terms it sounds very enticing.

The downside is that the user's browser is now sharing all of its local activity with the configured DoH server. To put it a different way, what part of "sharing your entire personal profile with the browser-selected DoH server" is consistent with our traditional concepts of personal privacy and informed choice?

Consider a second concern here as well. This ability for a browser and a DoH resolver to combine and thereby effectively dominate the Internet namespace is a legitimate concern. Few companies are in such a position, but there are few companies left in the Internet ecosystem. A very small number of digital behemoths inhabits the core of the Internet, and these entities could potentially take advantage of such an opportunity, were it offered to them. Google is the dominant provider of the platform in Android, the browser in Chrome, and the DNS resolver in the 8.8.8.8 service. Would this scenario be a case of a single corporate entity being in a position of overarching control of the entire namespace of the Internet? Netflix already fielded an app that used its own DNS resolution mechanism independent of the platform upon which the app was running. What if the Facebook app included DoH? What if Apple's iOS used a DoH-resolution mechanism to bypass local DNS resolution and steer all DNS queries from Apple's platforms to a set of Apple-operated name resolvers?

We'll find out some answers to these questions in the near future. On April 9, 2019, Mozilla announced its plan to enable DoH by default in the Firefox browser^[6], committing to an earlier informal description of its plans that were outlined by Mozilla's Eric Rescorla at the end of March 2019^[7].

To place this announcement into a broader perspective, it should be noted that the market share of Mozilla's Firefox browser, while large, is by no means dominant. The StatCounter site reported a market share of 4.69% for Firefox in March 2019^[5], so these moves by Mozilla are not intrinsically all that significant in terms of the profile of the larger Internet and the average Internet user. A major concern with this announcement is that the move by the Firefox browser to make DoH the default means of DNS name resolution is a precursor for similar changes to the Chrome browser. Chrome is definitely the dominant browser in today's Internet ecosystem, with some 62.63% market share according to StatCounter. If Chrome were to use a default setting that pushed all its DNS name-resolution activities to a Chrome-selected DoH server, then the implications for the DNS are very significant.

Will the other browsers follow Mozilla's lead with DoH enabled by default? The experience so far would support a "yes" answer. Browser vendors have been enthusiastic to integrate changes to their platform that decrease page load times, and they are equally keen to integrate changes that protect the browser activity against various forms of surveillance.

DoH does not necessarily make DNS resolution quicker, although it does put the browser in more control over its use of the DNS and allows the browser to control its own local DNS cache. But, of course, DoH plugs a critical DNS information leak in the current browser architecture. Third-party observers can infer browser activity by looking at the browser DNS query stream. DoH prevents any such observation in both the user's platform and the local network. So "yes" is a likely answer to this question.

Can such positions be regulated? How can we be assured that transactions that now have disappeared from sight, and from any meaningful form of oversight, are still conducted with all due integrity? We have already seen many national regimes struggle with very real questions concerning the limitations of imposing constraints on the actions of these entities. Have the concerns of the U.S. Supreme Court's Louis Brandeis in the first half of the twentieth century over the rise of industrial and financial behemoths that in his view were too big to effectively regulate at all come full circle?

What Does DoH Mean?

Here is the core of the collective angst and disquiet in the IETF when considering the implications of DoH and the "centrality" of Internet infrastructure.

We are attempting to actively withhold the DNS from the traditional forms of inspection and interception using access carriers and wire-based mechanisms. In so doing we are looking to counter what was perceived as a state-based surveillance operation that had assumed too much capability.

But in the case of the DNS have we over-achieved? In withholding our DNS secrets from one party, have we instead handed the entire plate to another? Have we now provided the private surveillance framework with a whole new trove of personal data to mine by ruthlessly exploiting the DNS in a manner that is entirely out of sight? When the browsers and even the apps direct their name queries through encrypted channels to resolvers operated by the same browser and app providers, then have we dealt a body blow to any efforts to safeguard personal privacy on the Internet?

At least RFC 7871 on Client Subnet included an admonition to operators to turn it off and a tacit apology for specifying a tool that had serious issues relating to erosion of user privacy in the DNS infrastructure. The DoH specification in RFC 8484^[11] contains no such considerations. It fails to mention the security and privacy issues if a browser invisibly co-opted the name-resolution function and passed all its DNS traffic in a secure encrypted tunnel to a cooperating resolver using DoH that faithfully mimics conventional content transactions. It fails to mention the risks of increasing the "centrality" of the Internet when the DNS name resolution is forcibly sucked into the browser and application space and then concealed behind a veil of strong encryption.

It's incredibly challenging to make the case that DoH enhances personal privacy. It probably doesn't. It's easier to sustain a case that DoH has the potential to change the parties whom you bring into your trust circle by virtue of their being privy to your private profile, and not necessarily in a good way. In and of itself such a substitution of trust should not necessarily be of concern. But now it's your browser that can make the decision as to whom you are trusting with your personal data, not you. And the parties who are looking to be your DoH trust partner are the same parties who have a direct and overbearing interest in selling you to the highest bidding advertiser.

Privacy Undertakings

These open DNS providers appear to have a clear view of user concerns over personal privacy. Their privacy policies implicitly acknowledge that the DNS query stream could be used to provide insights into the personal profile of users and assert that they have no such intent to do so. Such noble intentions to operate a free public service and refrain from any form of monetization of the service are entirely laudable.

However, from an historical perspective these undertakings appear to be unrealistic and unsustainable. We should remember the events of a century ago with Theodore Vail and the *Kingsbury Commitment* in 1913 in the United States. His key commitment was a profession of noble intent to enrich the public space. AT&T was to be an "enlightened monopoly" that served the public in close cooperation with the state while at the same time serving the interests of AT&T shareholders. His view of the telephone service as a privately operated public utility is, to quote Tim Wu in his treatise on Vail and AT&T, "...at once the most sympathetic and scariest element of his vision. Vail saw no harm in, and indeed believed in, giants, so long as they be friendly giants. He believed power should be beneficently concentrated, and that with great power came great responsibility."^[8]

As we observe the aggregation of this critical part of the Internet infrastructure in the centralization of the DNS, it cannot be ignored that these grand statements of respect for the public interest and undertakings that safeguard personal privacy sound scarily similar to the espoused public benefactor vision of AT&T in 1913 as it embarked on a course of establishing a national monopoly. But it is perhaps not today's operators and today's commitments that should concern us, but where this condition may lead. Again, quoting Tim Wu: "[Theodore Vail] presents us therefore with a challenging figure: an unabashed monopolist, but a benign one, who lived up to his own ideals of enlightened despotism. The fault in this arrangement therefore lay not so much with Theodore Vail as with the men who would succeed him."^[8]

Perhaps the same is true of these current undertakings relating to protection of personal privacy and their perception of the greater public interest.

Over time these earnest undertakings in the provision of free services may well be eroded by the inevitable pressures that every private enterprise is prone to, namely those of paying the bills and maximizing shareholder value. After the DNS is placed under an all-encompassing shroud of deep encryption, then both good and dark deeds will be both indistinguishable and undetectable.

It appears that the original disquiet on the part of the IETF was not that state-sponsored intelligence agencies collected intelligence, because, after all, that is their role, but a perception that the public accountability of some of these agencies had, in the IETF's view, failed. It is ironic that the IETF's response appears to literally hand the keys to an encrypted DNS over to a handful of private sector entities that appear to have no enduring public accountability whatsoever.

References and Further Reading

- [0] Electronic Privacy Information Center (EPIC), "An Open Letter to the Internet Engineering Task Force," November 8, 1999.
https://www.epic.org/privacy/internet/letter_to_ietf.html
- [1] IAB and IESG, "IETF Policy on Wiretapping," RFC 2804, May 2000.
- [2] ETSI, "Lawful Interception (LI)," <https://www.etsi.org/technologies/lawful-interception>
- [3] Stephen Farrell and Hannes Tschofenig, "Pervasive Monitoring Is an Attack," RFC 7258, May 2014.
- [4] Wilmer van der Gaast, Carlo Contavalli, and Warren Kumari, "Client Subnet in DNS Queries," RFC 7871, May 2016.
- [5] StatCounter, "Browser Market Share Worldwide," <http://gs.statcounter.com/browser-market-share>
- [6] Marshall Erwin, "DNS-over-HTTPS Policy Requirements for Resolvers," Mozilla Security Blog, April 9, 2019.
<https://blog.mozilla.org/security/2019/04/09/dns-over-https-policy-requirements-for-resolvers/>
- [7] Eric Rescorla, "Mozilla's plans re: DoH," IETF Mailing List Archive, March 27, 2019.
<https://mailarchive.ietf.org/arch/msg/doh/po6GCAJ52BAKuyL-dZiU91v6hLw>
- [8] Tim Wu, *The Master Switch: The Rise and Fall of Information Empires*, Borsoi Books, ISBN-13: 978-0307269935, 2010.
- [9] John Heidemann, Duane Wessels, Allison Mankin, Paul Hoffman, and Liang Zhu, "Specification for DNS over Transport Layer Security (TLS)," RFC 7858, May 2016.

- [10] Sara Dickinson, Tirumaleswar Reddy, and Daniel Gillmor, “Usage Profiles for DNS over TLS and DNS over DTLS,” RFC 8310, March 2018.
- [11] Patrick McManus and Paul Hoffman, “DNS Queries over HTTPS (DoH),” RFC 8484, October 2018.
- [12] IETF DNS PRIVate Exchange (DPRIVE) Working Group:
<https://datatracker.ietf.org/wg/dprive/charter/>
- [13] The OpenLI Project: <https://openli.nz/>
- [14] Dan Wing, Tirumaleswar Reddy, and Prashanth Patil, “DNS over Datagram Transport Layer Security (DTLS),” RFC 8094, February 2017.
- [15] Stephane Bortzmeyer, “DNS Query Name Minimisation to Improve Privacy,” RFC 7816, March 2016.
- [16] CALEA:
<https://www.fcc.gov/public-safety-and-homeland-security/policy-and-licensing-division/general/communications-assistance>
- [17] Paul Vixie, Joao Damas, and Michael Graff, “Extension Mechanisms for DNS (EDNS(0)),” RFC 6891, April 2013.
- [18] Geoff Huston and Joao Luis Silva Dama, “DNS Privacy,” *The Internet Protocol Journal*, Volume 20, No. 1, March 2017.
- [19] Patrick McManus and Paul E. Hoffman, “DNS-over-HTTPS (DoH) Operational and Privacy Issues,” IETF Blog,
<https://www.ietf.org/blog/doh-operational-and-privacy-issues/>
- [20] Catalin Cimpanu, “First-ever malware strain spotted abusing new DoH (DNS over HTTPS) protocol,” *ZDNet*, July 3, 2019.
<https://www.zdnet.com/article/first-ever-malware-strain-spotted-abusing-new-doh-dns-over-https-protocol/>
- [21] Catalin Cimpanu, “UK ISP group names Mozilla ‘Internet Villain’ for supporting ‘DNS-over-HTTPS’,” *ZDNet*, July 4, 2019.
<https://www.zdnet.com/article/uk-isp-group-names-mozilla-internet-villain-for-supporting-dns-over-https/>

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

Improving Routing Security

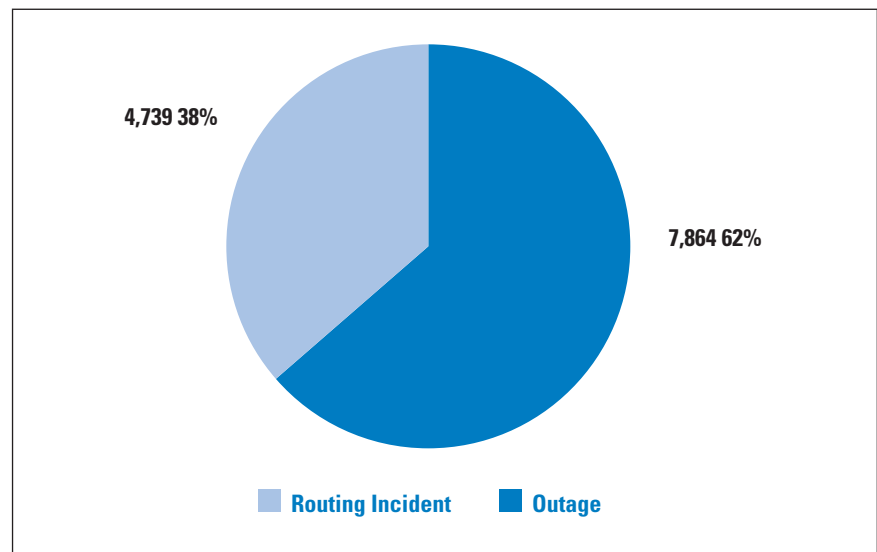
by Andrei Robachevsky, *The Internet Society*

Not a single day goes by without dozens of incidents affecting the routing system of the Internet. Route hijacking, route leaks, IP address spoofing, and other harmful activities can lead to *Denial of Service* (DoS) attacks, traffic inspection and surveillance, lost revenue, reputational damage, and more.

According to our analysis based on BGPStream^[0] data, the following numbers indicate the scale of the problem along with the comparison to data from 2017:

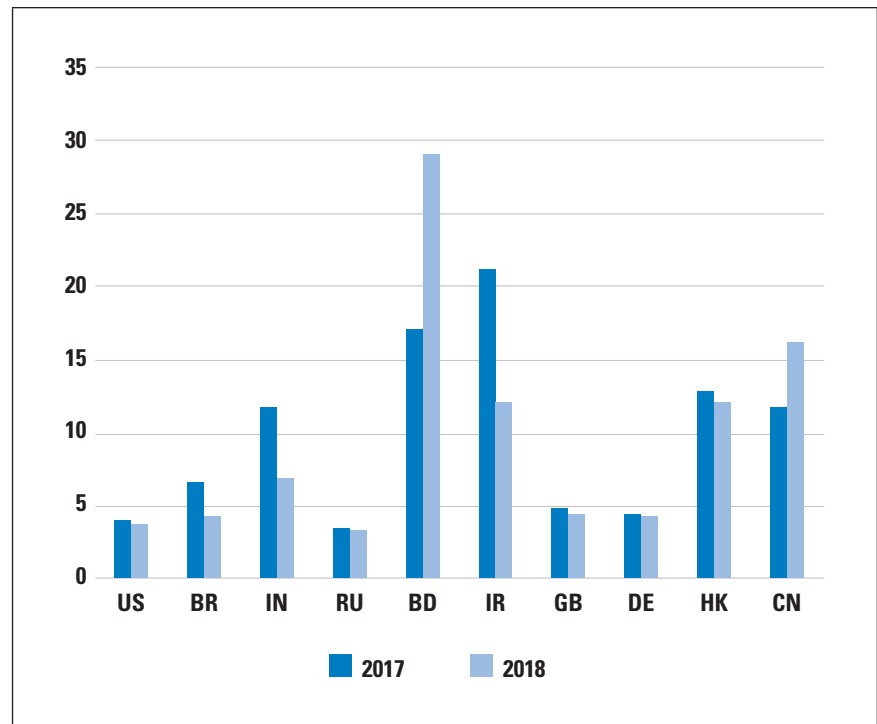
- In 2017, there were 12,600 (a 9.6% decrease) incidents (either outages or attacks such as route leaks and hijacks). Figure 1 shows the number of routing incidents by type in 2018.
- Although the overall number of incidents was reduced, the ratio of outages vs. routing security incidents remained unchanged — 62/38.
- About 4.4% (a decrease of 1%) of all Autonomous Systems on the Internet were affected.
- 2,737 (a decrease of 12%) Autonomous Systems experienced at least one routing incident.
- 1,294 (a 17% decrease!) networks were responsible for 4,739 routing incidents (a 10.6% decrease).

Figure 1: Routing incidents by type in 2018; almost 40% of all incidents were due to routing security issues.



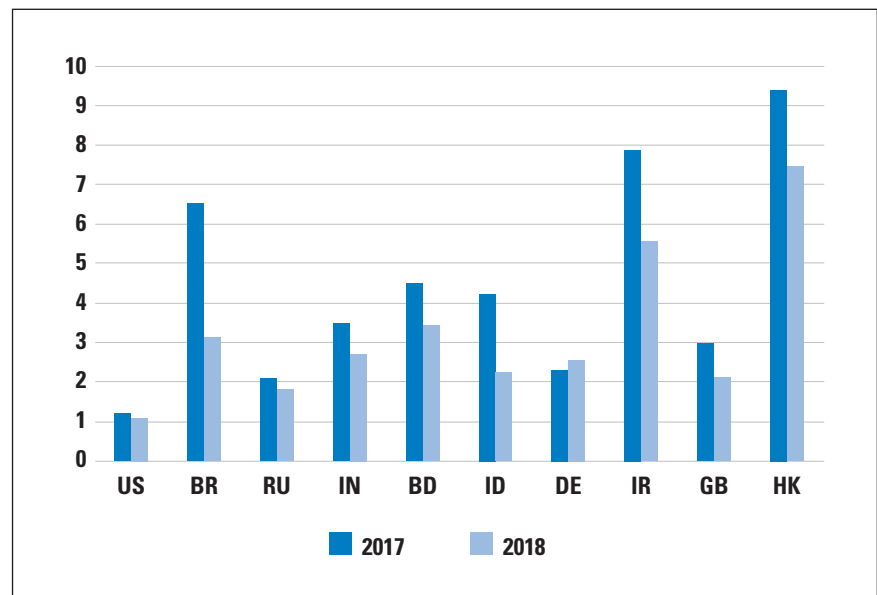
Looking further at the dynamics of the evolution of routing security from the perspective of affected networks (Figure 2), we can see that Bangladesh, mainland China, and Hong Kong appear to be the most vulnerable, with up to 30% of all networks affected by a routing mishap. We can also see positive dynamics in countries like Brazil, India, and Iran, where the percentage of “victimized” networks decreased significantly.

Figure 2: Changes in the percentage of affected networks in a country. Top-10 most affected countries.



Let's look at networks whose configuration mistakes or intentional acts caused these incidents. In absolute numbers, 36% of all “culprits” operate in the US, Brazil, and Russia, but if we normalize this number by the total number of networks in a country, mainland China and Hong Kong are at the top (see Figure 3). On a positive side, the situation has improved in most of the top-10 countries with the highest number of culprits. For instance, in Brazil the number of misbehaving networks has been cut by more than half.

Figure 3: Changes in the percentage of networks in a country responsible for a routing incident, top 10 with most incidents detected.



Why Is Routing Security Hard?

Despite the positive trend, it is too early to celebrate victory; vulnerabilities still exist, and too many networks are not applying required controls to prevent incidents from happening. What is holding the networks from resolving this problem once and forever?

There are several reasons. In the Internet, as a decentralized system, the overall level of routing security depends on the individual actions of all network operators, and incidents in most cases are impossible to address by your own operator. The economics favor insecurity, as the impacts of routing incidents are often felt by others and not by the culprit, and security has not yet emerged as a market differentiator. To put it another way—the controls that are necessary to reduce routing incidents, and that the operator should apply, improve the overall security, and to a much lesser extent they offer protection to their own networks. In other words, the security of your network is in the hands of other network operators, with whom you may not have any relation-ship. Therefore, addressing security issues in the Internet routing system requires a collective action.

Another challenge is related to the fact that security in general is not a state, but a process. Implementing security requires a systemic approach, and that is why corporate security relies on frameworks and established processes. How is it possible to apply such a systemic approach in a decentralized system with more than 60K independent networks?

MANRS^[1], *Mutually Agreed Norms for Routing Security*, attempts to address both challenges.

How MANRS Can Help

MANRS, a global initiative driven by network operators and *Internet Exchange Points* (IXPs) all over the globe and supported by the Internet Society, outlines simple but concrete actions that different types of network operators should take. The actions are limited in scope, and backed up by a growing community they have a good chance to become true norms of security in network operations.

Norms are often seen as possible solutions to a so-called *Collective Action Problem*. The name of this social phenomenon, known for centuries, was coined by Mancur Olson in 1965 in his book *The Logic of Collective Action*^[2]. Not really a problem in small communities, it becomes a real challenge as the number of entities grows, resulting in the failure to cooperate because of conflicting interests, despite a clear common benefit. That phenomenon is exactly what we observe in the area of routing security in the Internet.

Let's look at a set of actions that MANRS offers. Four actions are defined for *Internet Service Providers* (ISPs):

- *Action 1. Filtering:* Ensure the correctness of your own routing announcements and of announcements from your customers to adjacent networks with prefix and *Autonomous System (AS)-path granularity*.
- *Action 2. Anti-spoofing:* Enable source address validation for at least single-homed stub customer networks, your own end users, and infrastructure.
- *Action 3. Coordination:* Maintain globally accessible up-to-date contact information.
- *Action 4. Global Validation:* Publish your routing policy, including the intended announcements, so others can validate routing information on a global scale.

These actions represent a minimum baseline that yields significant improvements to the routing system with relatively little effort from individual players. The actions also provide a global reference that other initiatives or corporate improvement projects can use as a starting building block. This process can help focus various efforts in the area of routing security for steady and continuous improvement on a global scale.

Another aspect of MANRS is related to the interdependency and the fact that only a collective solution is possible. Not only does MANRS serve as a recommendation of what to do, but it also builds a community of security-minded operators committed to the common cause. The community is crucial in reinforcing the baseline and transforming it into norms of operational behavior.

Network operators join MANRS not out of pure altruism. Many understand that a stable and secure communication fabric is an essential component for their growing business. Many of the operators that joined MANRS were already implementing good routing security and even exceeding the requirements of the actions. However, MANRS, as a global reference point, provides them an opportunity to signal their security posture to customers and regulators. Moreover, the growing MANRS community is a clear demonstration that the industry is taking action to address these complex security issues.

IXPs Onboard

ISPs are not the only players that affect routing security. For instance, IXPs form local communities of ISPs with a common operational objective. They are in an excellent position to support the reciprocity of network protection that routing security requires and create a “safe neighborhood” at the exchange. To take advantage of the impact IXPs can have in the area of routing security, the MANRS community set the goal to get IXPs on board.

But IXPs are not exactly ISPs. And since MANRS membership requires demonstration of commitment with a tangible contribution, the community has created a related but separate set of MANRS actions for participating IXPs:

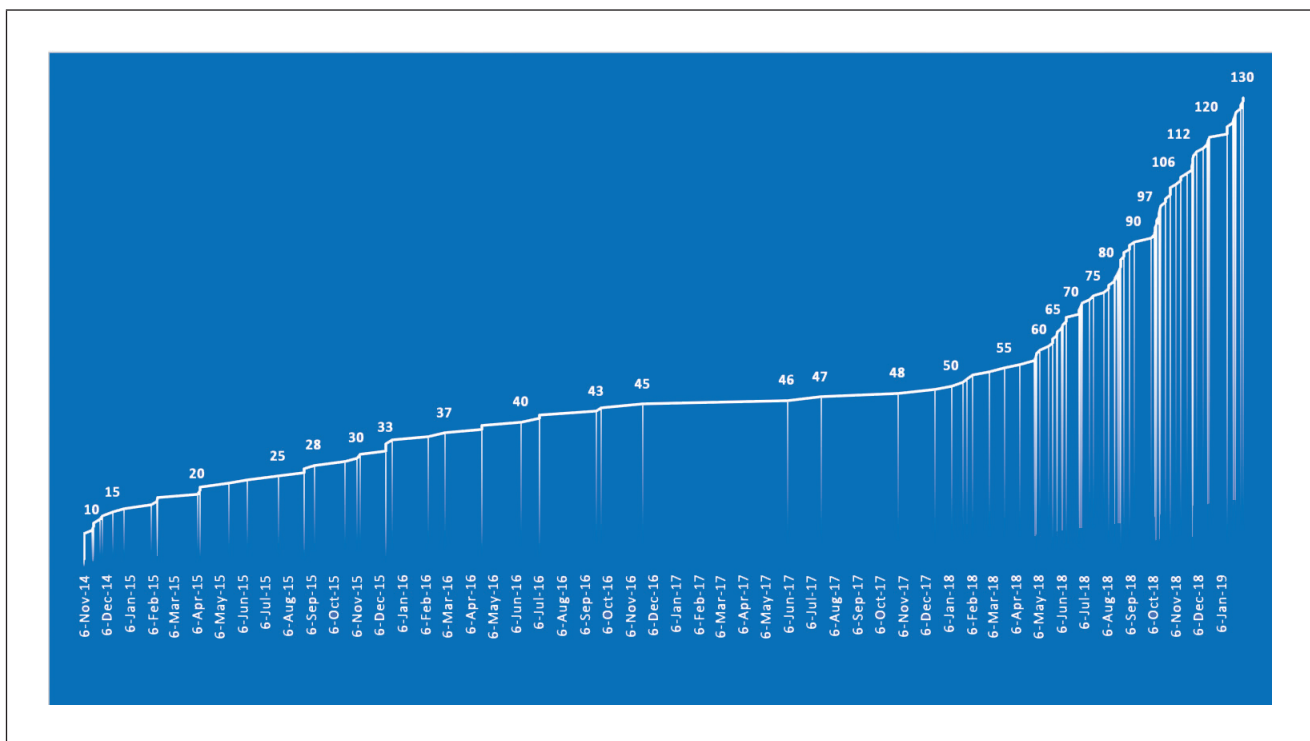
- *Action 1. Prevent propagation of incorrect routing information:* This mandatory action requires IXPs to implement filtering of route announcements at the Route Server based on routing information data (*Internet Routing Registry* [IRR] and/or *Resource Public Key Infrastructure* [RPKI]). A Route Server is a proxy network used to facilitate multilateral peering between the operators. Instead of setting up multiple *Border Gateway Protocol* (BGP) peering sessions with various operators at the exchange, an operator can peer with only the Route Server to accomplish this goal.
- *Action 2. Promote MANRS to the IXP membership:* IXPs joining MANRS are expected to provide encouragement or assistance for their members to implement MANRS actions.
- *Action 3. Protect the peering platform:* This action requires that the IXP have a published policy of traffic not allowed on the peering fabric and perform filtering of such traffic.
- *Action 4. Facilitate global operational communication and coordination among its members by providing necessary mailing lists and member directories.*
- *Action 5. Provide monitoring and debugging tools, such as the Looking Glass (LG)^[3]:* BGP LG servers are computers on the Internet running one of a variety of publicly available Looking Glass software implementations. A LG server is accessed remotely for the purpose of viewing routing information. Essentially, the server acts as a limited, read-only portal to routers of whatever organization is running the LG server. Typically, publicly accessible LG servers are run by ISPs or *Network Operations Centers* (NOCs).

Membership Growth

MANRS has seen a steady growth of its membership since its launch in November 2014. At the time of this writing, 130 operators cover more than 250 ASNs. Since the launch of the IXP Programme^[4] in April 2018, the number of participating IXPs has reached 30 (Figure 4).

MANRS needs partners to scale up adoption. IXPs are a great example of such collaboration. MANRS also partners with organizations such as APNIC, *Latin America and Caribbean Network Information Centre* (LACNIC), *Latin American and Caribbean Internet Exchange Association* (LAC-IX), the *Brazilian Network Information Center* (NIC.BR), Internet2, GÉANT, and RedCLARA to reach out to regional communities to grow the MANRS membership.

Figure 4: MANRS membership growth, ISPs.



Education and training play a very important role helping to lower the threshold for adoption. Based on the “Implementation Guide” developed by the community, MANRS Online Training^[5] contains six modules to help engineers understand the implementation details of the actions. These online modules can be completed either individually or as part of a moderated class, earning a certificate of completion from the Internet Society.

Next step in the capacity-building program is the release of the online hands-on lab; its development is being finalized.

In the area of capacity building, MANRS partners with training organizations such as the *Network Startup Resource Center* (NSRC) and the *Asia Pacific Network Information Centre* (APNIC)—reaching out to hundreds of network engineers.

As the awareness grows, more and more organizations are evaluating their readiness for MANRS actions, making necessary adjustments and joining the effort. The capacity-building efforts help networks that lack necessary expertise to implement the actions quickly.

It is our collective responsibility as participants of the Internet global routing system to ensure the reliability and security of the Internet. Help us make the Internet a safer place. Only together can we protect the core.

References and Further Reading

- [0] BGPStream: <https://bgpstream.com/>
- [1] MANRS: <https://www.manrs.org/>
- [2] Mancur Olson, *The Logic of Collective Action: Public Goods and the Theory of Groups*, Harvard University Press, ISBN 0-674-53751-3, 1971.
- [3] BGP: the Border Gateway Protocol Advanced Internet Routing Resources: <http://www.bgp4.as/looking-glasses>
- [4] MANRS IXP Programme: <https://www.manrs.org/ixps/>
- [5] MANRS Tutorials: <https://www.manrs.org/tutorials/>
- [6] Salam Yamout, “Improving Routing Security: Microsoft Joins MANRS,” Internet Society Blog, May 22, 2019.
<https://www.internetsociety.org/blog/2019/05/improving-routing-security-microsoft-joins-manrs/>
- [7] Dan Goodin, “Google goes down after major BGP mishap routes traffic through China,” Ars Technica, November 12, 2018.
<https://arstechnica.com/information-technology/2018/11/major-bgp-mishap-takes-down-google-as-traffic-improperly-travels-to-china/>
- [8] “Routing Security for Policymakers,” An Internet Society White Paper, October 2018.
<https://www.internetsociety.org/resources/doc/2018/routing-security-for-policymakers/>

- [9] Matthew Lepinski, “BGPsec Protocol Specification,” RFC 8205, September 2017.
- [10] Wesley George and Sandra Murphy, “BGPsec Considerations for Autonomous System (AS) Migration,” RFC 8206, September 2017.
- [11] Randy Bush and Geoff Huston, “Securing BGP with BGPsec,” *The Internet Protocol Journal*, Volume 14, No. 2, June 2011.
- [12] Stephen Kent, “Securing the Border Gateway Protocol,” *The Internet Protocol Journal*, Volume 6, No. 3, September 2003.

ANDREI ROBACHEVSKY is the Senior Technology Programme Manager for the Internet Society. Andrei’s primary area of interest is security and resilience of the Internet infrastructure. This work is based on active engagement with the operator, research, and policy communities. Prior to joining ISOC, Andrei was Chief Technical Officer of the RIPE NCC, leading the development of the company’s IT strategy, external and internal IT services, and work of the engineering departments. He was responsible for the deployment of DNSSEC for the reverse DNS tree and deployment of anycast instances of the K-root DNS server. Andrei brings to the Internet Society more than 10 years of experience in the Internet technical community. For more than a decade he has been actively following Regional Internet Registry (RIR) and Internet Engineering Task Force (IETF) activities. He was Chair of the Number Resource Organization’s (NRO) Engineering Coordination Group (ECG), which is responsible for various technical inter-RIR activities and projects. In 2010–2012 Andrei was a member of the Internet Architecture Board (IAB). Andrei is based in Amsterdam, The Netherlands. E-mail: robachevsky@isoc.org

ISOC Signs Letter Opposing GCHQ Proposal for Weakening Encryption

In late 2018, The British *Government Communications Headquarters* (GCHQ) published an essay^[1] on *Lawfare* outlining its principles for “exceptional” or “lawful” access to encrypted information, alongside a proposed use case—the “ghost proposal.” (Generally, when people speak of lawful or exceptional access they refer to some means of allowing law enforcement the ability to lawfully access the content of encrypted communications and encrypted data in an unencrypted form. For example, by asking companies to have the technical ability to access encrypted content.)

The GCHQ proposal would add a silent (or ghost) user to end-to-end encrypted messaging services, such as *WhatsApp*, and allow the government to listen in to ongoing encrypted conversations secretly for law enforcement or national security purposes. The *Internet Society* is pleased to add its name to an open letter^[2] outlining the dangers that this proposal, and techniques like it, pose to the Internet and to users everywhere.

All exceptional or lawful access proposals put users, the economy, the services we depend on and the Internet itself at greater risk to security threats. GCHQ’s “ghost proposal” is no exception. As stated in the open letter, the ghost proposal would:

“...introduce potential unintentional vulnerabilities, and increase risks that communications systems could be abused or misused ... [and] mean that users cannot trust that their communications are secure.”

Protected communications are a matter of security. Whether they are used to keep critical infrastructure running, safeguard our financial information, or keep personal information from those who would use it to do us harm, protected communications keep us all safe. All of these rely on encryption and other digital security tools.

The ISOC is proud to add its voice to a diverse group of stakeholders from civil society, industry and academia calling on GCHQ to abandon the ghost proposal and avoid any alternate approaches that would similarly threaten digital security and human rights. We must strengthen, not weaken encryption. By whatever name, any point of entry to a secure service is a weakness.

[1] Ian Levy and Crispin Robinson, “Principles for a More Informed Exceptional Access Debate,” *Lawfare*, November 2018.
<https://www.lawfareblog.com/principles-more-informed-exceptional-access-debate>

[2] Open Letter to GCHQ:
https://newamericadotorg.s3.amazonaws.com/documents/Coalition_Letter_to_GCHQ_on_Ghost_Proposal_-_May_22_2019.pdf

ICANN Publishes Updated Domain Name Marketplace Indicators

The *Internet Corporation for Names and Numbers* (ICANN) recently announced publication of the first wave of the *Domain Name Marketplace Indicators* report, which presents statistics related to *generic top-level domains* (gTLDs) and *country code top-level domains* (ccTLDs).^[1]

This report is an evolution of the previous *gTLD Marketplace Health Index* report (Beta), which was first published in July 2016 with twice annual reports through June 2018. This report includes expanded coverage to include ccTLD data. ICANN plans to further expand its coverage of shortlisted indicators and continue to publish these statistics twice a year to track progress against its goal of supporting the evolution of the domain name marketplace to be robust, stable and trusted.

A community Advisory Panel worked with ICANN to refine these indicators in preparation for publishing this version. Concurrent to the release of these Version 1.0 marketplace indicators, ICANN org will continue to work with the community and the Advisory Panel to evaluate additional enhancements that might be incorporated into this initiative in the future.

ICANN's mission is to help ensure a stable, secure and unified global Internet. To reach another person on the Internet, you need to type an address—a name or a number—into your computer or other device. That address must be unique so computers know where to find each other. ICANN helps coordinate and support these unique identifiers across the world. ICANN was formed in 1998 as a not-for-profit public-benefit corporation with a community of participants from all over the world.

[1] <https://www.icann.org/resources/pages/metrics-gdd-2015-01-30-en>

Check your Subscription Details!

If you have a print subscription to this journal, you will find an expiration date printed on the back cover. For the last couple of years, we have “auto-renewed” your subscription, but now we ask you to log in to our subscription system and perform this simple task yourself. The subscription portal is here: <https://www.ipjsubscription.org/> This process will ensure that we have your current contact information as well as delivery preference (print edition or download). For any questions, contact us by e-mail at: ipj@protocoljournal.org

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Fabrizio Accatino	David Cardwell	Andrew Fox	John Jarvis
Michael Achola	John Cavanaugh	Craig Fox	Dennis Jennings
Martin Adkins	Lj Cemerar	Fausto Franceschini	Edward Jennings
Christopher Affleck	Dave Chapman	Tomislav Futivic	Aart Jochem
Scott Aitken	Stefanos Charchalakakis	Edward Gallagher	Brian Johnson
Jacobus Akkerhuis	Greg Chisholm	Andrew Gallo	Curtis Johnson
Antonio Cuñat Alario	David Chosrova	Chris Gamboni	Richard Johnson
Nicola Altan	Marcin Cieslak	Xosé Bravo Garcia	Jim Johnston
Matteo D'Ambrosio	Brad Clark	Kevin Gee	Jonatan Jonasson
Jens Andersson	Narelle Clark	John Gilbert	Daniel Jones
Danish Ansari	Steve Corbató	Serge Van Ginderachter	Gary Jones
Tim Armstrong	Brian Courtney	Greg Goddard	Jerry Jones
Richard Artes	Dave Crocker	Octavio Alfageme	Amar Joshi
David Atkins	Kevin Croes	Gorostiaga	Merike Kao
Jaime Badua	John Curran	Barry Greene	Andrew Kaiser
Eric Baker	André Danthine	Richard Gregor	Christos Karayiannis
Santosh Balagopalan	Morgan Davis	Martijn Groenleer	David Kekar
David Belson	Jeff Day	Geert Jan de Groot	Jithin Kesavan
Hidde Beumer	Freek Dijkstra	Christopher Guemez	Jubal Kessler
Pier Paolo Biagi	Geert Van Dijk	Gulf Coast Shots	Shan Ali Khan
John Bigrow	David Dillow	Sheryll de Guzman	Nabeel Khatri
Orvar Ari Bjarnason	Richard Dodsworth	James Hamilton	Dae Young Kim
Axel Boeger	Ernesto Doelling	Stephen Hanna	Anthony Klopp
Keith Bogart	Eugene Doroniuk	Martin Hannigan	Henry Kluge
Mirko Bonadei	Karlheinz Dölger	John Hardin	Michael Kluk
Roberto Bonalumi	Joshua Dreier	David Harper	Andrew Koch
Julie Bottorff Photography	Lutz Drink	Edward Hauser	Ia Kochiashvili
Gerry Boudreaux	Andrew Dul	David Hauweele	Carsten Koempe
L de Braal	Holger Durer	Marilyn Hay	Alexander Kogan
Kevin Breit	Mark Eanes	Headcrafts SRLS	Antonin Kral
Thomas Bridge	Peter Robert Egli	Hidde van der Heide	Mathias Körber
Ilia Bromberg	George Ehlers	Johan Helsingius	John Kristoff
Václav Brožík	Peter Eisses	Robert Hinden	Terje Krogdahl
Christophe Brun	Torbjörn Eklöv	Asbjorn Hojmark	Bobby Krupczak
Gareth Bryan	Y Ertur	Damien Holloway	Murray Kucherauw
Stefan Buckmann	ERNW GmbH	Alain Van Hoof	Dirk Kurfuerst
Caner Budakoglu	ESdatCo	Edward Hotard	Warren Kumari
Darrell Budic	Steve Esquivel	Bill Huber	Darrell Lack
Scott Burleigh	Jay Etchings	Hagen Hultzs	Yan Landriault
Jon Harald Bøvre	Mikhail Evstiounin	Kevin Iddles	Sig Lange
Olivier Cahagne	Paul Ferguson	Mika Ilvesmaki	Markus Langenmair
Antoine Camerlo	Kent Fichtner	Karsten Iwen	Fred Langham
Tracy Camp	The Flirble Organisation	David Jaffe	Andrew Lamb
Ignacio Soto Campos	Gary Ford	Ashford Jaggernaut	Richard Lamb
Fabio Caneparo	Jean-Pierre Forcioli	Martijn Jansen	Tracy LaQuey Parker
Roberto Canonico	Christopher Forsyth	Jozef Janitor	Rick van Leeuwen

Simon Leinen	Tariq Mustafa	Ron Rockrohr	Paul Stancik
Robert Lewis	Stuart Nadin	Carlos Rodrigues	Ralf Stempfer
Martin Lillepuu	Michel Nakhla	Magnus Romedahl	Matthew Stenberg
Roger Lindholm	Mazdak Rajabi Nasab	Lex Van Roon	Adrian Stevens
Sergio Loreti	Krishna Natarajan	William Ross	Clinton Stevens
Eric Louie	Darryl Newman	Boudhayan Roychowdhury	John Streck
Guillermo a Loyola	Thomas Nikolajsen	Carlos Rubio	Viktor Sudakov
Hannes Lubich	Paul Nikolich	Timo Rüter	Edward-W. Suor
Dan Lynch	Travis Northrup	RustedMusic	Vincent Surillo
Miroslav Madić	Marijana Novakovic	Babak Saberi	T2Group
Alexis Madriz	David Oates	George Sadowsky	Roman Tarasov
Carl Malamud	Ovidiu Obersterescu	Scott Sandefur	David Theese
Michael Malik	Tim O'Brien	Sachin Sapkal	Douglas Thompson
Yogesh Mangar	Mike O'Connor	Arturas Satkovskis	Lorin J Thompson
Bill Manning	Mike O'Dell	PS Saunders	Joseph Toste
Harold March	Jim Oplotnik	John Sayer	Rey Tucker
Vincent Marchand	Carlos Astor Araujo Palmeira	Phil Scarr	Sandro Tumini
David Martin	Alexis Panagopoulos	Elizabeth Scheid	Angelo Turetta
Jim Martin	Gaurav Panwar	Jeroen Van Ingen Schenau	Phil Tweedie
Ruben Tripiana Martin	Manuel Uruena Pascual	Carsten Scherb	Steve Ulrich
Timothy Martin	Ricardo Patara	Ernest Schirmer	Unitek Engineering AG
Gabriel Marroquin	Dipesh Patel	Dan Schrenk	John Urbanek
Carles Mateu	Alex Parkinson	Richard Schultz	Martin Urwaleck
Juan Jose Marin Martinez	Craig Partridge	Roger Schwartz	Betsy Vanderpool
Ioan Maxim	Dan Paynter	SeenThere	Surendran
David Mazel	Leif Eric Pedersen	Scott Seifel	Vangadasalam
Miles McCredie	Rui Sao Pedro	Yury Shefer	Buddy Venne
Brian McCullough	Juan Pena	Yaron Sheffer	Alejandro Vennera
Joe McEachern	Chris Perkins	Doron Shikmoni	Luca Ventura
Jay McMaster	David Phelan	Tj Shumway	Tom Vest
Mark Mc Nicholas	Derrell Piper	Jeffrey Sicuranza	Dario Vitali
Carsten Melberg	Rob Pirnie	Thorsten Sideboard	Lakhinder Walia
Kevin Menezes	Marc Vives Piza	Andrew Simmons	Laurence Walker
Bart Jan Menkveld	Jorge Ivan Pincay Ponce	Pradeep Singh	Randy Watts
William Mills	Victoria Poncini	Henry Sinnreich	Andrew Webster
David Millsom	Blahoslav Popela	Geoff Sisson	Tim Weil
Desiree Miloshevic	Eduard Llull Pou	Helge Skrivervik	Jd Wegner
Joost van der Minnen	Tim Pozar	Darren Sleeth	Westmoreland
Thomas Mino	David Raistrick	Richard Smit	Engineering Inc.
Wijnand Modderman	Priyan R Rajeevan	Bob Smith	Rick Wesson
Mohammad Moghaddas	Balaji Rajendran	Courtney Smith	Peter Whimp
Charles Monson	Paul Rathbone	Mark Smith	Russ White
Andrea Montefusco	William Rawlings	Job Snijders	Jurrien Wijnhuizen
Fernando Montenegro	Bill Reid	Ronald Solano	Derick Winkworth
Joel Moore	Rodrigo Ribeiro	Asit Som	Pindar Wong
Maurizio Moroni	Glenn Ricart	Ignacio Soto Campos	Romeo Zwart
Brian Mort	Justin Richards	Peter Spekrijse	Bernd Zeimetz
Soenke Mumm	Mark Risinger	Thayumanavan Sridhar	廖明沂.



Follow us on Twitter and Facebook

@protocoljournal



<https://www.facebook.com/newipj>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Diamond Sponsors



Ruby Sponsors

Your logo here!

Sapphire Sponsors

Your logo here!

Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

David Conrad, Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

December 2019

Volume 22, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

FROM THE EDITOR

In This Issue

From the Editor	1
MSS Values of TCP	2
50 Years of the Internet	12
Book Review.....	15
Letter to the Editor	17
Fragments	18
Thank You!	24
Call for Papers.....	26
Supporters and Sponsors	27

“A major design feature of the *Internet Protocol* (IP) is its ability to run over a variety of underlying network technologies. If you look through the *Request For Comments* (RFC) document series, you will find numerous specifications of the form “IP over xxx,” where “xxx” is anything from Ethernet to X.25, Frame Relay, Bluetooth, WiFi, and even “Avian Carriers” (pigeons), the latter being one of the more famous April Fools RFCs. Because each of these technologies has different capabilities in terms of how much data can be carried in a “packet” or datagram, IP employs the concept of *fragmentation* and *reassembly* in cases where the originating datagram is larger than what the underlying network medium can support.”

That paragraph is a quote from our June 2016 issue (Volume 19, No. 2) in which Geoff Huston described various aspects of IPv4 and IPv6 fragmentation. In this issue he explains how the *Transmission Control Protocol* (TCP) and its concept of a *Maximum Segment Size* (MSS) might interact with IP fragmentation even if this interaction is technically a “layer violation.” His article presents measurement data on TCP MSS handshakes recorded by APNIC in August 2019.

The Internet has its origins in the *Advanced Research Projects Agency Network* (ARPANET), which began operation just over 50 years ago, in October 1969, with only two nodes. We asked Vint Cerf, one of the “Fathers of the Internet,” to reflect on this milestone.

As always, we welcome your feedback and suggestions on anything you read in this journal. Letters to the Editor may be edited for clarity and length and can be sent to ipj@protocoljournal.org. Please make sure your subscription details are accurate. In this issue you will also find a summary of our Privacy Policy.

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

MSS Values of TCP

by Geoff Huston, APNIC

The *Transmission Control Protocol* (TCP) *Maximum Segment Size* (MSS) has been under some examination in recent months because an operating system vulnerability related to the Linux implementation of TCP occurred; it is described in CVE-2019-11477, 11478, and 11479^[1]. One of the effective work-arounds to avoid this problem is to block all TCP connection attempts that use a MSS value of 500 or lower.

What is the impact of such a TCP filter policy? What is being used as MSS values? How will a drop filter of TCP sessions with an MSS value of 500 or lower affect the Internet user base? In the *Asia-Pacific Network Information Centre* (APNIC) measurement platform we have assembled a large collection of recorded TCP handshakes, each of which contains a record of the TCP MSS exchange. Let's look at the MSS settings.

The TCP MSS Parameter

The MSS parameter is a part of the *Options* field in the TCP initial handshake that specifies the largest amount of data that a TCP speaker can receive in a single TCP segment^[2]. The MSS relates to the TCP input buffer size within the implementation as packets are passed from the IP module to the TCP module. Each direction of TCP traffic uses its own MSS value, as this value is receiver-specified. The two ends don't have to agree on a common value because it acts as a constraint on the sender to send TCP segments no larger than this MSS value. But of course smaller TCP segments can always be sent. This MSS value can vary between the forward and reverse directions of a TCP data flow.

The MSS value does not count the TCP header or the IP header. The received IP datagram containing a TCP segment may be self-contained within a single packet, or it may be reconstructed from several fragmented pieces.

Because IP packet fragmentation is an IP-level issue, TCP should not directly concern itself with IP fragmentation in any case. In theory. In practice, a judicious setting of TCP MSS sizes that attempts to avoid sending TCP packets that incur IP-level packet fragmentation should be avoided!

Conventionally, the platform rather than the application sets the MSS value for a connection and the setting is applied to all TCP connections. But many operating system platforms provide a hook in the connection *Application Programming Interface* (API) for an application to specify the MSS value for a connection (such as the TCP_MAXSEG socket option).

What Is a “Good” MSS Value?

Getting the MSS value “just right” is important. While in theory the IP and TCP layers are largely independent, the practical reality is quite the opposite. Too high a value can lead to inefficient and even wedged TCP sessions due to issues with mishandling of IP fragmentation. The problem is that the sender may perform TCP segmentation by using the received MSS value as its guide and the resultant TCP packet is then far larger than the outgoing IP interface *Maximum Transmission Unit* (MTU) size, entailing the sender to perform IP-level fragmentation on the TCP packet.

It’s also worth remembering that many of the TCP congestion control protocols use a rate acceleration based on an increase in the sending rate of 1 MSS of data per round-trip-time interval. Larger MSS values imply a faster rate of acceleration in such protocols, while smaller MSS values will lead to inefficiencies and may stall the sender, potentially leading to some congestion issues within the sender. The IPv4 packet contains a 16-bit packet identification number, implying that in order to avoid fragmentation reassembly issues not more than 65,536 IP packets should be in flight at any point in time.

Therefore, the combination of very small MSS values, long-held TCP sessions, and long-delay bandwidth network paths is certainly inefficient, but it should not necessarily represent any form of attack vector if the implementation of TCP is suitably robust. The recent security notices point to some platform vulnerabilities within the sender that are exposed by low MSS values.

What guidance is there in the RFCs on setting the TCP MSS value?

RFC 791^[3] provides IP MTU guidance, stating that:

"All hosts must be prepared to accept datagrams of up to 576 octets (whether they arrive whole or in fragments). It is recommended that hosts only send datagrams larger than 576 octets if they have assurance that the destination is prepared to accept the larger datagrams."

Given that IP has no explicit MTU signalling capability, this explicit recommendation of obtaining assurance of the receiver’s preparedness to accept larger IP datagrams presumably refers to the TCP MSS value.

RFC 879^[4] provided quite explicit guidance about the TCP MSS value:

"THE TCP MAXIMUM SEGMENT SIZE IS THE IP MAXIMUM DATAGRAM SIZE MINUS FORTY. The default IP Maximum Datagram Size is 576. The default TCP Maximum Segment Size is 536."

These documents were written prior to the specification of IPv6 of course, and in RFC 2460^[5] the following guidance was given for IPv6:

"When using TCP over IPv6, the MSS must be computed as the maximum packet size minus 60 octets."

It also states that:

"IPv6 requires that every link in the internet have an MTU of 1280 octets or greater."

Taken together, RFC 2460 asserts that for TCP over IPv6 the MSS value would be expected to be 1,220 or greater.

These days the now-ancient Ethernet packet-framing specification still dominates the networking environment (although the old thick yellow coaxial cables and even the *Carrier Sense Multiple Access/Collision Detection* [CSMA/CD] 10-Mbps common bus protocol were both consigned to the networking section of silicon heaven years ago!). Thus the most common IP packet MTU is 1,500 octets.

Further clarification was provided in RFC 6691^[6], "TCP Options and Maximum Segment Size," (published in July 2012):

"When calculating the value to put in the TCP MSS option, the MTU value SHOULD be decreased by only the size of the fixed IP and TCP headers and SHOULD NOT be decreased to account for any possible IP or TCP options; conversely, the sender MUST reduce the TCP data length to account for any IP or TCP options that it is including in the packets that it sends. [...] the goal is to avoid IP-level fragmentation of TCP packets."

That information implies that the most common anticipated TCP MSS values would correspond to a 1,500-octet MTU in both IPv4 and IPv6, further implying that we should see a MSS value of 1,460 in IPv4 and 1,440 in IPv6.

How well does practice line up with the theory?

Measuring TCP MSS Values

We looked at the MSS sizes in the HTTP(S) sessions offered by clients who connected to servers with our measurement as part of our large measurement program into IPv6 deployment. We collected all TCP handshakes that occurred in August 2019 and recorded the MSS values from the SYN packets received from the client systems.

We saw some 3B TCP sessions over this period, and after we removed the duplicate entries for multiple TCP sessions from the same end-point within a similar timeframe with a common MSS value, we were left with some 551M unique TCP sessions.

Surprisingly enough, 202 endpoints offered an MSS value of 0, and 284 endpoints offered a value of 1. To put this data into perspective, this count of 486 endpoints represents 0.0001% (slightly less than 1 per million) of all observed TCP sessions. Small MSS values exist, but they are very much a rarity in the larger population of the Internet. A total of 20,488 sessions were opened with MSS values of 500 or lower (0.004%).

At the other end of the range of observed MSS values, three sessions used a value of 65,516 (the IPv4 maximum MTU minus 40). If we categorise any MSS value greater than 1,460 as some form of jumbo MSS, then we observed that 68,278 sessions used jumbo MSS values, or 0.012% of all TCP sessions.

As a side note, the network industry has never reached a clear agreement on exactly what a *jumbo frame* size should be. A value of 9,216 octets has been commonly quoted, as has the Internet2-defined value of 9,000 octets. The lack of agreement within the *Institute of Electrical and Electronics Engineers* (IEEE) on a single definition of a jumbo frame is not entirely unique, as many media-level protocols have used what could only be described in retrospect as idiosyncratic maximum MTU values. IEEE 802.5 *Token Ring* used an MTU of up to 4,464 octets, *Fiber Distributed Data Interface* (FDDI) used 4,532 octets, and IEEE 802.11 used 7,935 octets. Perhaps this diversity in media-based MTU values is not all that surprising, and what is perhaps more surprising is a current rough consensus of a commonly assumed MTU of 1,500 octets in the Internet, irrespective of the underlying media capabilities.

Table 1 shows the most common observed MSS values.

Table 1: Most Common MSS Values

Rank	MSS	Ratio	Rank	MSS	Ratio
1	1,460	17.6%	14	1,390	0.8%
2	1,400	16.4%	15	1,358	0.7%
3	1,370	11.3%	16	1,368	0.6%
4	1,452	8.7%	17	1,388	0.6%
5	1,440	8.3%	18	1,350	0.5%
6	1,360	6.9%	19	1,394	0.4%
7	1,412	5.1%	20	1,312	0.4%
8	1,300	4.3%	21	1,220	0.4%
9	1,380	3.7%	22	1,362	0.3%
10	1,420	3.5%	23	1,240	0.3%
11	1,432	1.8%	24	1,414	0.3%
12	1,410	1.3%	25	1,344	0.3%
13	1,340	1.3%			

The 1,460 value appears to correlate with a 1,500-octet MTU and a 40-octet IPv4 and TCP packet header. If the MSS value is calculated from the interface MTU size less the size of the IP and TCP headers, then we would expect the IPv6 MSS sizes to be 20 bytes less than the IPv4 MSS sizes.

We can separate the TCP MSS values used in IPv4 and IPv6, as shown in Table 2.

Table 2: Most Common MSS Values in IPv4 and IPv6

IPv4			IPv6	
Rank	MSS	Ratio	MSS	Ratio
1	1,460	17.6%	1,370	28.1%
2	1,400	16.4%	1,440	20.5%
3	1,370	11.3%	1,432	7.8%
4	1,452	8.7%	1,300	6.2%
5	1,440	8.3%	1,400	5.2%
6	1,360	6.9%	1,380	4.7%
7	1,412	5.1%	1,340	4.2%
8	1,300	4.3%	1,412	2.9%
9	1,380	3.7%	1,368	2.9%
10	1,420	3.5%	1,358	2.8%
11	1,432	1.8%	1,390	2.5%
12	1,410	1.3%	1,420	2.2%
13	1,340	1.3%	1,220	1.6%
14	1,390	0.8%	1,312	1.6%
15	1,358	0.7%	1,350	1.4%
16	1,368	0.6%	1,362	1.2%
17	1,388	0.6%	1,360	0.8%
18	1,350	0.5%	1,426	0.5%
19	1,394	0.4%	1,428	0.5%
20	1,312	0.4%	1,240	0.4%
21	1,220	0.4%	1,394	0.3%
22	1,362	0.3%	1,404	0.2%
23	1,240	0.3%	1,200	0.2%
24	1,414	0.3%	1,140	0.2%
25	1,344	0.3%	1,330	0.1%

The 1,370 value in IPv6 is somewhat unusual, as it corresponds to a MTU of 1,430 octets. It appears to be a common situation to use a 1,430-octet MTU in hosts, presumably as such a value (and any MTU value less than 1,460 octets) would minimise both the risks of IP fragmentation and path MTU issues that may arise from path element encapsulation that could be encountered when using IPv6.

If one takes the 1,460 MSS value in IPv4 and the 1,440 MSS value in IPv6 as an indicator of an underlying 1,500 MTU size, then relatively more endpoints are using a 1,500 MTU in IPv6 than in IPv4. (17.6% in IPv4 vs 20.5% in IPv6). In IPv6 there is a stronger consensus to use a single, smaller MSS value of 1,370 (28.1%) than there is in IPv4, where there is significant use of both 1,400 (16.4%) and 1,370 (11.3%) as MSS values.

The range of observed MSS values between 1,300 and 1,440 in both IPv4 and IPv6 points to the existence of a common action of constraining the IP MTU size in order to circumvent the possibility of IP fragmentation in both IPv4 and IPv6. I described the problem in 2009 in “A Tale of Two Protocols: IPv4, IPv6, MTUs and Fragmentation,”^[7] and pointed out why a reduced MTU setting would be a reasonable response to this problem.

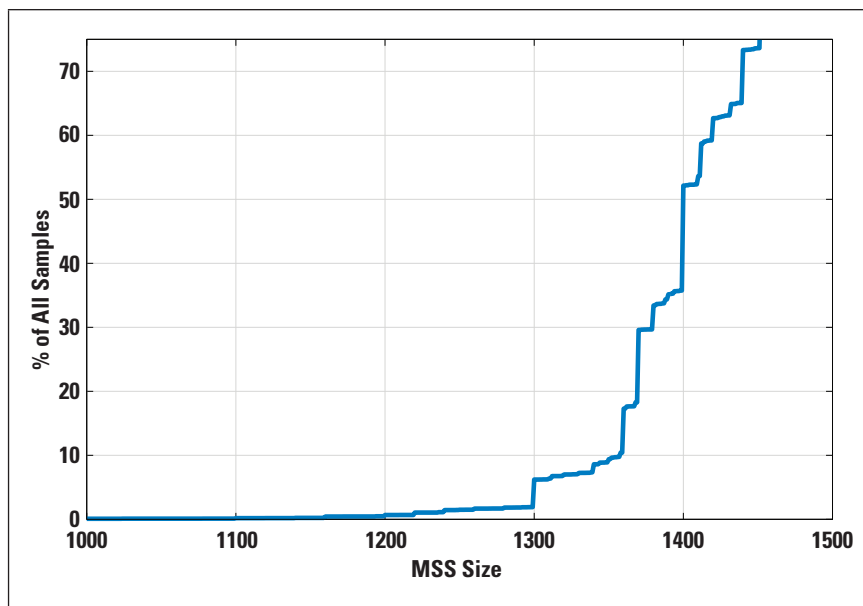
We also saw the MSS value of 536 in 38,359 cases, 38,245 of which were in IPv4 and 14 in IPv6, a value that appears to be derived from an assumed interface MTU of 576 octets and a 20-octet IPv4 packet header and 20-octet TCP header.

Oddly enough, the MSS value of 512 is more prevalent than that of 536, observed in 319,612 cases in IPv4 and not observed at all in IPv6.

While there is no particular media type that uses an MTU of 1,280 (which is the minimum unfragmented packet MTU size in IPv6), we had observed a minor clustering of MSS values at 1,220, with 1,892,966 IPv6 samples using this MSS value. Oddly enough, there were 7,268,197 cases of a 1,300-byte MSS value in IPv6.

Figure 1 on the following page shows the overall distribution of observed IPv4 MSS sizes.

Figure 1: Cumulative Distribution
of TCP MSS Values



Incidence of Low MSS Values

Where might we see hosts with low (less than 500) MSS values?

In IPv4 almost half of all such systems are located in Germany and the Netherlands. Adding the data from systems offering small MSS values from France, South Korea, Bangladesh, Indonesia, Pakistan, and Brazil to the set encompasses some 90% of all hosts with MSS values less than 500. In the case of IPv6, more than half of the low MSS values are from hosts located in Germany, and 90% of all such hosts are located in Germany, Indonesia, the United States, Canada, Malaysia, and Brazil.

In terms of origin network in IPv4, the networks that contain the most hosts with observed low MSS values are operated by a large web-hosting enterprise. In IPv6 the majority of instances originate from a research centre in Germany. It may be that this high incidence of these very low MSS values in these networks could be due to some bug or operational misconfiguration in web-hosting equipment, or an unintentional configuration choice made by a client of this virtual system hosting service.

Incidence of High MSS Values

And where are hosts that use large TCP MSS values (values greater than 1,460)?

In IPv4 the United States contains 21% of all such hosts, and the somewhat diverse collection of India, Russia, and Ireland also each host some 4 to 6% of the total count of such hosts. The picture alters with IPv6, with half of all such hosts located in Japan and a little under one-quarter located in the United States.

At the network level the Amazon *Autonomous System Numbers* (ASNs) were most commonly found to be hosting high MSS-valued TCP stacks in IPv4, while the Japanese *Internet Service Providers* (ISPs) KDDI and NTT's OCN and Comcast in the United States were hosting high-valued MSS hosts in IPv6.

It appears that while some form of hosting or cloud system generates a large MSS value in IPv4, some form of configuration of ISP server product might be the cause of this behaviour in IPv6.

Incidence of Adjusted MSS Values

We can make the supposition that an MSS value of between 1,260 and 1,440 has been the result of a deliberate adjustment of the host MTU value in order to reduce the risk of packet fragmentation and path MTU black holes.

A path MTU black hole occurs when a server emits a packet that is too large for a network path element on the path to the receiver and, in the case of IPv4, either the *Don't Fragment* bit of the packet is set or the packet is an IPv6 packet, and the return path to the server is blocking *Internet Control Message Protocol* (ICMP) messages for some reason.

At this point the connection will stall. The sender is waiting for either an ACK of the data sequence number in the dropped packet or an ICMP packet to indicate that there is an MTU problem. The ACK will never arrive as the packet has been dropped and the ICMP message has been blocked within the network.

The server will timeout and retransmit the large packet, to no effect. It may do so indefinitely unless some local overall session timeout is in effect, or TCP keepalives are in use.

The client has no outstanding data, so it will not retransmit and will just hang, waiting for a packet that will never arrive. TCP keepalives may identify this hung state and kill the hung TCP session.

We see these adjusted MSS values in those locations with a high IPv6 deployment volume—including India, the United States, Japan, and Vietnam.

What Values Should Be Used for TCP MSS?

A decade ago, the best advice around was to use a down-adjusted MSS value such as 1,300, 1,380, or even 1,400. The reason was to avoid path MTU issues, particularly when using IPv6, and the reason why path MTU issues were encountered in IPv6 was the prevalent use of IPv6-in-IPv4 encapsulation tunnels in IP transit paths and the widespread practice of firewall filtering of ICMPv6 *Packet Too Big* messages.

I'm not sure that ICMP filtering has improved or worsened in the last decade, but what has improved markedly is the use of "native" IPv6 transit paths in the Internet.

While it was probably foolhardy to use a 1,500 MTU and a 1,440 MSS with IPv6 a decade ago, it appears now to be not quite so foolhardy. Of course, not every tunnel has been removed and not every potential path MTU issue has been eliminated from the network, not every ICMPv6 filter has been removed, and not every fragment discard rule has been purged from firewalls, and they will probably never be completely purged. In relative terms the situation is better than it was a decade ago, and the expectation of encountering MTU-related problems is far lower than it was when a MSS based on a 1,500-octet MTU setting was used.

But there are still outstanding issues here, and a more reliable service can be staged using a slightly reduced local MTU and MSS setting. If we used a 1,480-octet MTU and corresponding TCP MSS values of 1,420 in IPv4 and 1,400 in IPv6, we could reasonably anticipate that the resultant TCP service would be adequately reliable.

As for the CVE mitigation advice to refuse a connection attempt when the remote-end MSS value is 500 or lower, I'd say that's good advice. It seems that the low MSS values are the result of some form of misconfiguration or error, and rather than attempting to mask over the error and persisting with an essentially broken TCP connection that is prone to generating a packet deluge, the best option is to just say "no" at the outset. If we all do that, then the misconfiguration will be quickly identified and fixed, rather than being silently masked over.

References and Further Reading

- [1] CVE-2019-11477, National Vulnerability Database, Information Technology Laboratory, National Institute of Standards and Technology (NIST), June 2019.
<https://nvd.nist.gov/vuln/detail/CVE-2019-11477>
- [2] J. Postel, "Transmission Control Protocol," RFC 793, September 1981.
- [3] J. Postel, "Internet Protocol," RFC 791, September 1981.
- [4] J. Postel, "The TCP Maximum Segment Size and Related Topics," RFC 879, November 1983.
- [5] Stephen E. Deering, "Internet Protocol, Version 6 (IPv6) Specification," RFC 2460, December 1998.
- [6] David Borman, "TCP Options and Maximum Segment Size (MSS)," RFC 6691, July 2012.

- [7] Geoff Huston, “A Tale of Two Protocols: IPv4, IPv6, MTUs and Fragmentation,” *The ISP Column*, January 2009.
<https://www.potaroo.net/ispcol/2009-01/mtu6.html>
- [8] Christopher A. Kent and Jeffrey C. Mogul, “Fragmentation Considered Harmful,” Proceedings of Frontiers in Computer Communications Technology, ACM SIGCOMM '87, August 1987.
- [9] Geoff Huston, “Fragmentation,” *The Internet Protocol Journal*, Volume 19, No. 2, June 2016.
- [10] Geoff Huston, “IPv6 and Packet Fragmentation,” *The Internet Protocol Journal*, Volume 21, No. 1, April 2018.

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

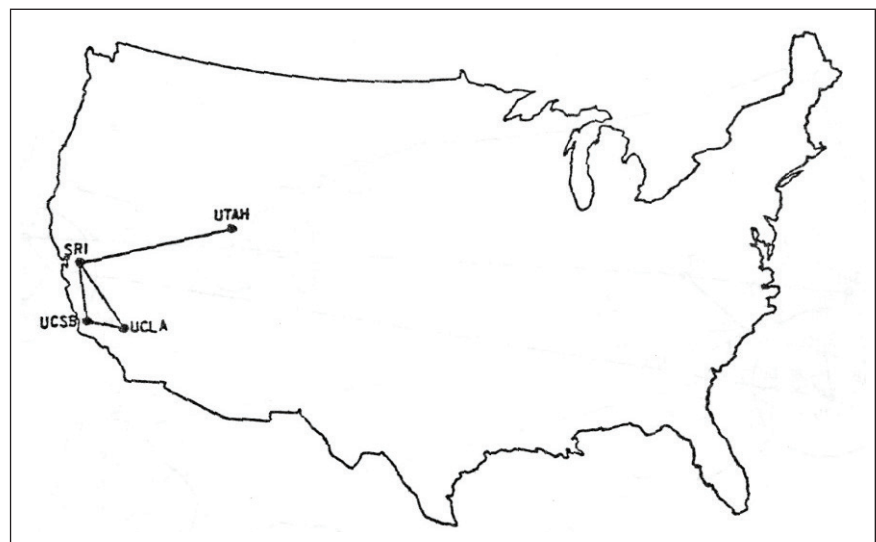
Looking Back on 50 Years of the Internet Era

by Vint Cerf, Internet Pioneer

It isn't possible in a short essay to really cover 50 years of the evolution of the Internet. Books have been written on the subject. We have just celebrated the October 29, 1969, milestone that linked an XDS Sigma-7 computer at the *University of California, Los Angeles* (UCLA) to a time-shared SDS 940 computer at *Stanford Research Institute* (SRI). At that time, there were only two nodes of the planned *Advanced Research Projects Agency Network* (ARPANET)^[0], but two more at *University of California, Santa Barbara* (UCSB) and at *University of Utah, Salt Lake City* were added by the end of the year. What I think is most interesting about the ARPANET part of the Internet saga is the trailblazing that project did in heterogeneous computer networking. It was among the first networks to use *packet switching* as the communications mechanism. Some historians might reasonably argue that the US *Semi-Automated Ground Environment* (SAGE) system developed in the 1950s and in operation until the 1980s and the US *AUTODIN* message switching system also developed in the late 1950s and 1960s both represented then state-of-the-art, wide-area computer communications. These systems were closer in spirit to automated store-and-forward teletype/telegraph messaging services than to the subsequent packet switching of the ARPANET, the French *Cyclades/Cigale* experimental network, and the UK National Physical Laboratory local-area network.

The heterogeneous computers (called *hosts*) of the ARPANET were interconnected to each other through a subnet of identical *Interface Message Processors* (IMPs) that were, in turn, interconnected by dedicated 50 kbps telephone circuits.

Figure 1: The ARPANET in December 1969^[4]



The ARPANET project launched a documentation series called *Request for Comments* (RFCs)^[1] that continues to this day to document the protocols of the Internet.

To coordinate the development of the host-to-host protocols and other applications, a *Network Working Group* (NWG) was established. This concept influenced the creation of an *International Network Working Group* (INWG) that became Working Group 6.1 of Technical Committee 6 of the *International Federation of Information Processing* (IFIP) and also influenced the creation of the *Internet Configuration Control Board* that morphed into the *Internet Architecture Board* (IAB), which gave rise to the *Internet Engineering and Research Taskforces* (IETF and IRTF), all of which are still active today.

A spirit of open sharing and cooperation permeated the participants and organizations that were involved in the ARPANET project and that also influenced the early developers of the Internet and the World Wide Web—which emerged in the early 1990s as the most popular application of the system. Even in today’s competitive environment, we find the engineers of the Internet cooperating to deal with a plethora of malicious attacks against the infrastructure and applications of the Internet and to fashion new protocols to support a growing collection of applications.

As the protocol experiments and research unfolded on the ARPANET, the concept of protocol layering emerged and strongly influenced both the Internet design and the development of the *Open Systems Interconnection* (OSI) Model for computer networking. The layering concepts also influenced the basic Internet architecture in the form of encapsulation and decapsulation of Internet packets in the frames and packet payloads of lower-level protocols in the underlying networks of the Internet. Gateways received Internet packets in the payloads of lower layers, extracted them, decided where to route them, encapsulated them in the appropriate payloads of the next packet network, and sent them on their way.

Electronic messaging, which had been developed in the course of implementing time-shared computer systems, was extended as networked electronic mail in the early years of the ARPANET to work across the network among the cooperating hosts. A *File Transfer Protocol* (FTP)^[2] and a remote-access *telecommunications network* (TELNET)^[3] protocol were among the early applications of the ARPANET and were eventually translated into the Internet.

A look back to the early years of the ARPANET shows that we owe much to those pioneering researchers and engineers who blazed trails into terra incognita for the rest of us to follow and extend. Even today, there is still an enormous frontier of unexplored conceptual space waiting to be discovered. As we collectively struggle to deal with emergent challenges of misinformation, disinformation, denial-of-service attacks, and fragmentation of the Internet, I still remain hopeful that the utility of willing global collaboration will inform the Internet governance policies under consideration around the world and bend them towards positive and fruitful outcomes.

References and Further Reading

- [0] “ARPANET,” Wikipedia entry:
<https://en.wikipedia.org/wiki/ARPANET>
- [1] “History,” RFC Editor webpage:
<https://www.rfc-editor.org/history/>
- [2] J. Postel and J. Reynolds, “File Transfer Protocol,” RFC 959, October 1985.
- [3] A. M. McKenzie, “Telnet Protocol Specifications,” RFC 495, May 1973.
- [4] Image courtesy of J. Noel Chiappa, MIT Advanced Network Architecture Group.
- [5] Daniel Dern, “The ARPANET Is Twenty: What We Have Learned and The Fun We had,” *ConneXions—The Interoperability Report*, Volume 3, No. 10, October 1989. Archive available from The Charles Babbage Institute at the University of Minnesota:
<http://www.cbi.umn.edu/hostedpublications/Connexions/index.html>

VINTON G. CERF is vice president and Chief Internet Evangelist for Google. He contributes to global policy development and continued spread of the Internet. Widely known as one of the “Fathers of the Internet,” Cerf is the co-designer of the TCP/IP protocols and the architecture of the Internet. He has served in executive positions at MCI, the Corporation for National Research Initiatives and the Defense Advanced Research Projects Agency and on the faculty of Stanford University.

Vint Cerf served as chairman of the board of the *Internet Corporation for Assigned Names and Numbers* (ICANN) from 2000–2007 and has been a Visiting Scientist at the Jet Propulsion Laboratory since 1998. Cerf served as founding president of the *Internet Society* (ISOC) from 1992–1995. Cerf is a Foreign Member of the British Royal Society and Swedish Academy of Engineering, and Fellow of IEEE, ACM, and American Association for the Advancement of Science, the American Academy of Arts and Sciences, the International Engineering Consortium, the Computer History Museum, the British Computer Society, the Worshipful Company of Information Technologists, the Worshipful Company of Stationers and a member of the National Academy of Engineering. He currently serves as Past President of the Association for Computing Machinery, Past Chairman of the *American Registry for Internet Numbers* (ARIN) and completed a term as Chairman of the Visiting Committee on Advanced Technology for the US National Institute of Standards and Technology. President Obama appointed him to the National Science Board in 2012.

Cerf is a recipient of numerous awards and commendations in connection with his work on the Internet, including the US *Presidential Medal of Freedom*, US *National Medal of Technology*, the *Queen Elizabeth Prize for Engineering*, the *Prince of Asturias Award*, the *Tunisian National Medal of Science*, the *Japan Prize*, the *Charles Stark Draper* award, the *ACM Turing Award*, Officer of the Legion d’Honneur and 29 honorary degrees. In December 1994, *People* magazine identified Cerf as one of that year’s “25 Most Intriguing People.” In 2012, he was inducted to the *Internet Hall of Fame*. E-mail: vint@google.com

Book Review

Confessions of a Crypto Millionaire

Confessions of a Crypto Millionaire: My Unlikely Escape from Corporate America, by Dan Conway, Zealot Publishing, September 2019, ISBN-13: 978-1733171700.

You probably have read your fill of business books. Author tries to make it big, leverages tons of his money and time, hires the wrong people, fires them, then goes it alone before striking it rich and motoring off into the sunset in some expensive car. Dan Conway's *Confessions of a Crypto Millionaire* is not one of these books. Most business books offer just enough advice to fill a chapter, maybe two. Conway has a lot more to say about his obsession and investments in cryptocurrency, in particular *Ethereum*. Over a period of several years, he used his home mortgage equity loan and borrowed additional funds because he believed blockchain held the future model for decentralized corporations and the way that we will all work together. He ended up cashing out \$14M ahead. It is his obsession that drives the book's narrative, along with the crazy up-and-down valuation of Ether, where you can gain and lose millions in a matter of minutes.

What isn't in this book is also notable: sordid tales of wretched excess of "tech-bros partying on yachts" or trashing expensive Vegas hotel suites. Conway is a father of three, and still married to their mother.

Conway's confessions make a refreshing tale about his fighting his demons, his addictions (alcohol and pills), his insecurities, and his almost always-on self-destructive alter-ego he calls his "Flip Side." This side rears its ugly head during client presentations where he fumbles and fails and during periods of self-doubt when he tries to reassure himself his huge bet on Ether isn't about to land him in the poor house.

"The book forced me to make sense of how my addictive personality played a part in my undoubtedly reckless crypto investments," he told me via an e-mail interview. He is part visionary, buying Ether at a time and at a level few people had the courage, vision, or just dumb luck to do. "It took everything admirable and loathsome about me to make the plunge into Ether. The loathsome part includes my addictive personality. While betting everything was an extreme risk, all risk requires insight, courage, and maybe a little recklessness." He hopes his story will get others to think about how they formulate their own risk taking.

Conway begins his story "working for the man," doing marketing and public relations for large corporations, one of whom he calls Acme. He wasn't a good fit as the organization man to be sure. And since his windfall with Ether, he is unlikely to return to corporate America "unless we suffer a financial catastrophe."

He still believes that the decentralized blockchain can disrupt the traditional corporate power structure and has a lot of merit as an organizing principle. One example he cites is the MakeDAO, where ordinary folks can originate loans and handle other financial transactions without any financial institutional limits. It could pay off; it could fall flat: that is the challenge of cryptocurrency.

One aspect of his book is dealing very honestly with two situations: first, with his addictions. “This undoubtedly played a part in my reckless crypto investments, and writing the book helped force me to make sense of it all.”

Second, the book also describes how his financial windfall changed his family dynamics and the relationships with his circle of friends. Even though Conway lived in Silicon Valley, he was very firmly rooted in the middle class before he made it big with Ether. He writes:

“Crypto was suddenly like an overexposed celebrity, and everyone was rooting for it to fail,” but then realizes, “one of the bittersweet feelings about making a bunch of money is that you can’t bring your (less fortunate) friends with you.” That takes some adjustment, both for him and his family. Still, don’t be too sad: Now he takes long exotic vacations, buys his kids “name-brand clothes” instead of Sears knock-offs, and does car pool duty with a vengeance. “It’s absolutely nice to have the car-ride conversations rather than pinning all parent/child bonding on the “how was your day?” question when everyone is exhausted.” True that.

Conway is committed to Ethereum because of its disruptive ability to change the way companies operate, the way companies get Venture Capital funding (the parts about the ICO shysters alone are worth reading), and the way the early pioneers—which Conway counts as himself—had to try to separate the criminals from the legit businesses. This book is well worth reading, even if your own exposure to bitcoin and other cryptocurrencies is minimal.

—David Strom, david@strom.com

Ed.: This book review originally appeared in David Strom’s *Web Informant*, available at: <https://blog.strom.com/wp/>

See also: William Stallings, “A Blockchain Tutorial,” *The Internet Protocol Journal*, Volume 20, No. 3, November 2017.

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. For more information, contact us at ipj@protocoljournal.org

Letter to the Editor

Hi Geoff,

I have read the article “DNS Privacy and the IETF,” in the latest issue of *Internet Protocol Journal* (Volume 22, No. 2, July 2019). Thank you for the excellent insights.

I am intrigued about the discussion of a world where apps would control the resolver. I am wondering how these apps would work in an IPv6-only world where the *Internet Service Provider* (ISP) does *Carrier Grade Network Address Translation* (CGNAT). If a DoH server responds with only an A record to say `ipv4google.com`, how does the handset make the connection? Or perhaps I misunderstand CGNAT—in reality the Internet-facing device/router has an RFC 1918 address and a global IPv6 address—IPv4 connectivity would use *Network Address and Port Translation* (NAPT) from the ISP public pool, whereas IPv6 connectivity can pass through without issue.

Thanks,

— Naveen Nathan, naveen@lastninja.net

The author responds:

Great question. Some transition mechanisms “crossed the beams” and relied on a DNS resolver that had knowledge of the transition mechanism and deliberately lied in their responses in order to steer the end host’s traffic to a protocol translator/encapsulator. Obviously if the application selected a DoH resolver that was not part of the local environment and was unaware of the need to provide NAT64 responses to these hosts, then the application would be unable to communicate.

It leads to the interesting outcome where the host (non-DoH) would look just fine and certain applications when going to certain remote services would fail. I have some sympathy for the help desk staff trying to identify and solve this problem.

Evidently some plans for DoH use involve application testing the existing configured DNS resolver for DoH capability and turning on DoH only in that case; that is, encapsulate in HTTP only the first DNS “hop” from the stub resolver to the recursive resolver. This solution would certainly avoid the NAT64 issue but would not really prevent the “my ISP is spying on my DNS transactions and possibly using this data in ways I am unaware of” scenario.

More generally, the more we adorn the network infrastructure and the more we add elements that create dependencies on other elements in novel ways, the more fragile the network becomes. The end point is a network that only barely functions and resists any modification—however slight—as the modification causes it to fail. It’s an odd situation to get to when the original concepts behind the Internet were thoughts about creating a level of resiliency of the network as a service that exceeded the resiliency of any component of the network system.

Regards,

— Geoff Huston, gih@apnic.net

Fragments

Postel Service Award Presented to Alain Aina

The Internet Society, a global nonprofit dedicated to ensuring the open development, evolution, and use of the Internet, recently presented the prestigious *Jonathan B. Postel Service Award* to Alain Aina, who serves as the chief technology officer of the *West and Central Africa Research and Education Network* (WACREN).

Aina has been building a Regional Research and Education Network to interconnect *National Research and Education Networks* (NRENs) in the region and connect them to the global Research and Education Network. He wants the world to see the work of Africa's premier researchers and carve out its spot in the academic world—in a way that would be impossible without the resources of this new network and community. He also contributes to *AfricaConnect2*, a project that supports the development of high-capacity networks for research and education across Africa, by building on existing networks in Eastern, Northern, and Southern Africa to connect to West and Central Africa's WACREN.



Photo by Minzayar Oo © IETF LLC 2019.

Aina fell into this work after graduating in the early 90s with a degree in electrical engineering and in the maintenance and analysis of computer systems. He was hired to be a technical seller for a company in the Togolese Republic, which had a branch in Benin, where he is from. The owner of the company had recently returned home from the United States and was anxious about computing and internet-working. He noticed Aina's talent and added him to the technical team, where he ended up building the first *Bulletin Board System* (BBS) in the area.

“People used the modem to dial in, then people on the same server could talk to each other,” he said. “Then we decided to put in the first e-mail gateway, connecting to someone in Accra and later in Montreal twice a day to drop mail and download mail. But the cost was so high, it was not sustainable. The delegation of the country-code TLD in 1996 changed the paradigm for the e-mail service and we were proud to demonstrate the first local web server and intranet.”

By the mid-90s there wasn’t a lot of support for people working on Internet access and connection, but there was ever-growing interest and demand. This meant that Aina and his colleagues often worked around the clock to set up networks and services in communities, then trained the local population on how to use what they had made.

“The Internet became so popular that the demand was suddenly so high, and it was putting pressure on us,” he said.

It was about this time that Aina started collaborating with the *Network Startup Resource Center* (NSRC), where he now serves as a part-time network engineer and trainer. He later launched the first full IP services in the Togolese Republic and then in other countries in West Africa.

“At that time, most of the world did not believe that Africa could have the Internet and play a role. When you’d go to places, you’d have to train people,” Aina said. “Training materials were rare. We were lucky to have some books and some knowledgeable friends far away. The people you trained only knew you, so if something broke they called you to fix it.”

Aina helped build large parts of the Internet ecosystem throughout Africa, setting up networks, contributing to the creation of the regional Internet registry and the network operator group, and building ccTLD registries. He also started a consulting firm and became active in the private sector.

He eventually started attending *Internet Society* (ISOC) network technology workshops and getting involved with the organization in other ways. From 2011 to 2014, he served as a trustee for the organization. Active in the Internet community, he is also involved with ICANN, the *African Network Information Centre* (AFRINIC), the *African Network Operators Group* (AFNOG), and other organizations. He helped found AFNOG, where he’s been an instructor since 2000, and he is one of the founders of AFRINIC, where he’s served in several roles, including acting chief technology officer, acting chief executive officer, and director of research and innovation. Aina is a key technical resource for the DNS community, including *Africa Top Level Domains Organization* (AFTLD). A big part of his life has been Internet related, but he feels there is still so much more to do for Africa.

Mr. Aina was selected by an international award committee comprised of former Postel award winners. The committee placed particular emphasis on candidates who have supported and enabled others in addition to their own contributions. The award was presented to Mr. Aina in recognition of his leadership in pioneering the Internet in Africa and building technical communities that helped connect countless others across the continent and beyond. Aina helped build large parts of the Internet ecosystem throughout Africa, setting up networks, contributing to the creation of the regional Internet registry and the network operator group, and building ccTLD registries.

“This award encourages me to continue the work, to grow and help others spread the Internet continent-wide, and to help break down barriers for the engineers and scientists in Africa,” Aina said. “I feel happy and honored to be recognized for this work.”

The Postel Award was established by the Internet Society to honor individuals or organizations that, like Jon Postel, have made outstanding contributions to the data communications community. The award is focused on sustained and substantial technical contributions, service to the community, and leadership. Andrew Sullivan, President and CEO of the Internet Society presented the award, including a US\$20,000 honorarium and a crystal engraved globe, during the 106th meeting of the *Internet Engineering Task Force* (IETF) held in Singapore, 16–22 November 2019.

ICANN Calls for Full DNSSEC Deployment

The *Internet Corporation for Assigned Names and Numbers* (ICANN) believes that there is an ongoing and significant risk to key parts of the *Domain Name System* (DNS) infrastructure. In the context of increasing reports of malicious activity targeting the DNS infrastructure, ICANN is calling for full deployment of the *Domain Name System Security Extensions* (DNSSEC) across all unsecured domain names. The organization also reaffirms its commitment to engage in collaborative efforts to ensure the security, stability and resiliency of the Internet’s global identifier systems.

As one of many entities engaged in the decentralized management of the Internet, ICANN is specifically responsible for coordinating the top-most level of the DNS to ensure its stable and secure operation and universal resolvability.

On 15 February 2019, in response to reports of attacks against key parts of the DNS infrastructure, ICANN offered a checklist^[1] of recommended security precautions for members of the domain name industry, registries, registrars, resellers, and related others, to proactively take to protect their systems, their customers’ systems and information reachable via the DNS.

Public reports^[2] indicate that there is a pattern of multifaceted attacks utilizing different methodologies. Some of the attacks target the DNS, in which unauthorized changes to the delegation structure of domain names are made, replacing the addresses of intended servers with addresses of machines controlled by the attackers. This particular type of attack, which targets the DNS, only works when DNSSEC is not in use. DNSSEC is a technology developed to protect against such changes by digitally “signing” data to assure its validity. Although DNSSEC cannot solve all forms of attack against the DNS, when it is used, unauthorized modification to DNS information can be detected, and users are blocked from being misdirected.

ICANN has long recognized the importance of DNSSEC and is calling for full deployment of the technology across all domains. Although this will not solve the security problems of the Internet, it aims to assure that Internet users reach their desired online destination by helping to prevent so-called *Man in the Middle* attacks where a user is unknowingly re-directed to a potentially malicious site. DNSSEC complements other technologies, such as *Transport Layer Security* (TLS) (most typically used in HTTPS) that protect the end user/domain communication.

As the coordinator of the top-most level of the DNS, ICANN is in the position to help mitigate and detect DNS-related risks, and to facilitate key discussions together with its partners. The organization believes that all members of the domain name system ecosystem must work together to produce better tools and policies to secure the DNS and other critical operations of the Internet.

- [1] “Alert Regarding Published Reports of Attacks on the Domain Name System,”

<https://www.icann.org/news-announcement-2019-02-15-en>

- [2] “A Deep Dive on the Recent Widespread DNS Hijacking Attacks,” Krebs on Security, February 19, 2019.

<https://krebsonsecurity.com/2019/02/a-deep-dive-on-the-recent-widespread-dns-hijacking-attacks/>

The RIPE NCC has run out of IPv4 Addresses

From the RIPE-NCC Website: “Today, at 15:35 (UTC+1) on 25 November 2019, we made our final /22 IPv4 allocation from the last remaining addresses in our available pool. We have now run out of IPv4 addresses. Our announcement will not come as a surprise for network operators—IPv4 run-out has long been anticipated and planned for by the RIPE community. In fact, it is due to the community’s responsible stewardship of these resources that we have been able to provide many thousands of new networks in our service region with /22 allocations after we reached our last /8 in 2012.

Even though we have run out, we will continue to recover IPv4 addresses in the future. These will come from organisations that have gone out of business or are closed, or from networks that return addresses they no longer need. These addresses will be allocated to our members (*Local Internet Registries* [LIRs]) according to their position on a new waiting list that is now active.

While we therefore expect to be allocating IPv4 for some time, these small amounts will not come close to the many millions of addresses that networks in our region need today. Only LIRs that have never received an IPv4 allocation from the RIPE NCC (of any size) may request addresses from the waiting list, and they are only eligible to receive a single /24 allocation. LIRs that have submitted an IPv4 request can see their position on the waiting list in the LIR Portal. A graph is available at <https://www.ripe.net/> that shows the number of requests on the waiting list and the number of days that the LIR at the front of the queue has been waiting.

This event is another step on the path towards global exhaustion of the remaining IPv4 addressing space. In recent years, we have seen the emergence of an IPv4 transfer market and greater use of *Carrier Grade Network Address Translation* (CGNAT) in our region. There are costs and trade-offs with both approaches and neither one solves the underlying problem, which is that there are not enough IPv4 addresses for everyone.

Without wide-scale IPv6 deployment, we risk heading into a future where the growth of our Internet is unnecessarily limited—not by a lack of skilled network engineers, technical equipment or investment—but by a shortage of unique network identifiers. There is still a long way to go, and we call on all stakeholders to play their role in supporting the IPv6 roll-out.

At the RIPE NCC, we are here to support our membership and the wider RIPE community in this work. Aside from allocating the IPv6 resources that will be required, we will continue to provide advice, training, measurements and tools to help network operators as they put their deployment plans into action. We are optimistic and excited to see what the next chapter will bring. So let’s get to work—and together, let’s shape the future of the Internet.”

Check your Subscription Details!

If you have a print subscription to this journal, you will find an expiration date printed on the back cover. For the last couple of years, we have “auto-renewed” your subscription, but now we ask you to log in to our subscription system and perform this simple task yourself. The subscription portal is here: <https://www.ipjsubscription.org/>. This process will ensure that we have your current contact information as well as delivery preference (print edition or download). For any questions, contact us by e-mail at: ipj@protocoljournal.org

Our Privacy Policy

The *General Data Protection Regulation* (GDPR) is a regulation for data protection and privacy for all individual citizens of the *European Union* (EU) and the *European Economic Area* (EEA). Its implementation in May 2018 led many organizations worldwide to post or update privacy statements regarding how they handle information collected in the course of business. Such statements tend to be long and include carefully crafted legal language. We realize that we may need to provide similar language on our website and in the printed edition, but until such a statement has been developed here is an explanation of how we use any information you have supplied relating to your subscription:

- The mailing list for *The Internet Protocol Journal* (IPJ) is entirely “opt in.” We never have and never will use mailing lists from other organizations for any purpose.
- You may unsubscribe at any time using our online subscription system or by contacting us via e-mail. We will honor any request to remove your name and contact information from our database.
- We will use your contact information *only* to communicate with you about your subscription; for example, to inform you that a new issue is available, that your subscription needs to be renewed, or that your printed copy has been returned to us as undeliverable by the postal authorities.
- We will *never* use your contact information for any other purpose or provide the subscription list to any third party other than for the purpose of distributing IPJ by post or by electronic means.
- If you make a donation in support of the journal, your name will be listed on our website and in print unless you tell us otherwise.

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Fabrizio Accatino	Roberto Canonico	The Flirble Organisation	Kevin Iddles	Warren Kumari
Michael Achola	David Cardwell	Gary Ford	Mika Ilvesmaki	Darrell Lack
Martin Adkins	John Cavanaugh	Jean-Pierre Forcioli	Karsten Iwen	Yan Landriault
Christopher Affleck	Lj Cemerar	Susan Forney	David Jaffe	Markus Langenmair
Scott Aitken	Dave Chapman	Christopher Forsyth	Ashford Jaggernauth	Fred Langham
Jacobus Akkerhuis	Stefanos Charchalakakis	Andrew Fox	Martijn Jansen	Andrew Lamb
Antonio Cuñat Alario	Greg Chisholm	Craig Fox	Jozef Janitor	Richard Lamb
Nicola Altan	David Chosrova	Fausto Franceschini	John Jarvis	Sig Lange
Matteo D'Ambrosio	Marcin Cieslak	Tomislav Futivic	Dennis Jennings	Tracy LaQuey Parker
Jens Andersson	Brad Clark	Edward Gallagher	Edward Jennings	Rick van Leeuwen
Danish Ansari	Narelle Clark	Andrew Gallo	Aart Jochem	Simon Leinen
Tim Armstrong	Steve Corbató	Chris Gamboni	Brian Johnson	Robert Lewis
Richard Artes	Brian Courtney	Xosé Bravo Garcia	Curtis Johnson	Martin Lillepuu
Michael Aschwanden	Dave Crocker	Osvaldo Gazzaniga	Richard Johnson	Roger Lindholm
David Atkins	Kevin Croes	Kevin Gee	Jim Johnston	Sergio Loreti
Jac Backus	John Curran	Greg Giessow	Jonatan Jonasson	Eric Louie
Jaime Badua	André Danthine	John Gilbert	Daniel Jones	Guillermo a Loyola
Eric Baker	Morgan Davis	Serge Van Ginderachter	Gary Jones	Hannes Lubich
Santosh Balagopalan	Jeff Day	Greg Goddard	Jerry Jones	Dan Lynch
David Belson	Julien Dhallenne	Tiago Goncalves	Anders Marius	Miroslav Madić
Hidde Beumer	Freek Dijkstra	Octavio Alfageme	Jørgensen	Alexis Madriz
Pier Paolo Biagi	Geert Van Dijk	Gorostiaga	Amar Joshi	Carl Malamud
John Bigrow	David Dillow	Barry Greene	Merike Kao	Jonathan Maldonado
Orvar Ari Bjarnason	Richard Dodsworth	Richard Gregor	Andrew Kaiser	Michael Malik
Axel Boeger	Ernesto Doelling	Martijn Groenleer	Christos Karayiannis	Yogesh Mangar
Keith Bogart	Michael Dolan	Geert Jan de Groot	David Kekar	Bill Manning
Mirko Bonadei	Eugene Doroniuk	Christopher Guemez	Jithin Kesavan	Harold March
Roberto Bonalumi	Karlheinz Dölger	Gulf Coast Shots	Jubal Kessler	Vincent Marchand
Julie Bottorff	Joshua Dreier	Sheryll de Guzman	Shan Ali Khan	Gabriel Marroquin
Photography	Lutz Drink	Jason Hall	Nabeel Khatri	David Martin
Gerry Boudreaux	Andrew Dul	James Hamilton	Dae Young Kim	Jim Martin
L de Braal	Joan Marc Riera	Stephen Hanna	Russell Kirk	Ruben Tripiana Martin
Kevin Breit	Duocastella	Martin Hannigan	Anthony Klopp	Timothy Martin
Thomas Bridge	Holger Durer	John Hardin	Henry Kluge	Juan Jose Marin
Ilia Bromberg	Mark Eanes	David Harper	Michael Kluk	Martinez
Václav Brožík	Peter Robert Egli	Edward Hauser	Andrew Koch	Carles Mateu
Christophe Brun	George Ehlers	David Hauweele	Ia Kochiashvili	Ioan Maxim
Gareth Bryan	Peter Eisses	Marilyn Hay	Carsten Koempe	David Mazel
Stefan Buckmann	Torbjörn Eklöv	Headcrafts SRLS	Richard Koene	Miles McCredie
Caner Budakoglu	Y Ertur	Hidde van der Heide	Alexader Kogan	Brian McCullough
Darrell Budic	ERNW GmbH	Johan Helsingius	Antonin Kral	Joe McEachern
Scott Burleigh	ESdatCo	Robert Hinden	Mathias Körber	Alexander McKenzie
Jon Harald Bøvre	Steve Esquivel	Asbjorn Hojmark	Robert Krejčí	Jay McMaster
Olivier Cahagne	Jay Etchings	Damien Holloway	John Kristoff	Mark Mc Nicholas
Antoine Camerlo	Mikhail Evstiounin	Alain Van Hoof	Terje Krogdahl	Carsten Melberg
Tracy Camp	Paul Ferguson	Edward Hotard	Bobby Krupczak	Kevin Menezes
Ignacio Soto Campos	Ricardo Ferreira	Bill Huber	Murray Kucherawy	Bart Jan Menkveld
Fabio Caneparo	Kent Fichtner	Hagen Hultzs	Dirk Kurfuerst	William Mills

David Millsom	Rob Pirnie	Yaron Sheffer	Surendran Vangadasalam
Desiree Miloshevic	Marc Vives Piza	Doron Shikmoni	Ramnath Vasudha
Joost van der Minnen	Jorge Ivan Pincay Ponce	Tj Shumway	Philip Venables
Thomas Mino	Victoria Poncini	Jeffrey Sicuranza	Buddy Venne
Rob Minshall	Blahoslav Popela	Thorsten Sideboard	Alejandro Vennera
Wijnand Modderman	Eduard Llull Pou	Greipur Sigurdsson	Luca Ventura
Mohammad Moghaddas	Tim Pozar	Andrew Simmons	Tom Vest
Charles Monson	David Raistrick	Pradeep Singh	Dario Vitali
Andrea Montefusco	Priyan R Rajeevan	Henry Sinnreich	Michael L Wahrman
Fernando Montenegro	Balaji Rajendran	Geoff Sisson	Laurence Walker
Joel Moore	Paul Rathbone	Helge Skrivervik	Randy Watts
Maurizio Moroni	William Rawlings	Darren Sleeth	Andrew Webster
Brian Mort	Bill Reid	Richard Smit	Tim Weil
Soenke Mumm	Petr Rejhon	Bob Smith	Jd Wegner
Tariq Mustafa	Robert Remenyi	Courtney Smith	Westmoreland
Stuart Nadin	Rodrigo Ribeiro	Mark Smith	Engineering Inc.
Michel Nakhla	Glenn Ricart	Job Snijders	Rick Wesson
Mazdak Rajabi Nasab	Justin Richards	Ronald Solano	Peter Whimp
Krishna Natarajan	Mark Risinger	Asit Som	Russ White
Naveen Nathan	Ron Rockrohr	Ignacio Soto Campos	Jurrien Wijnhuizen
Darryl Newman	Carlos Rodrigues	Evandro Sousa	Derick Winkworth
Thomas Nikolajsen	Magnus Romedahl	Peter Spekrijse	Pindar Wong
Paul Nikolich	Lex Van Roon	Thayumanavan Sridhar	Janko Zavernik
Travis Northrup	Alessandra Rosi	Paul Stancik	Muhammad Ziad
Marijana Novakovic	William Ross	Ralf Stempfer	Ziayuddin
David Oates	Boudhayan Roychowdhury	Matthew Stenberg	Romeo Zwart
Ovidiu Obersterescu	Carlos Rubio	Adrian Stevens	Bernd Zeimet
Tim O'Brien	Timo Rüter	Clinton Stevens	廖明沂.
Mike O'Connor	RustedMusic	John Streck	
Mike O'Dell	Babak Saberi	Martin Streule	
Jim Oplotnik	George Sadowsky	Viktor Sudakov	
Carlos Astor Araujo Palmeira	Scott Sandefur	Edward-W. Suor	
Alexis Panagopoulos	Sachin Sapkal	Vincent Surillo	
Gaurav Panwar	Arturas Satkovskis	T2Group	
Manuel Uruena Pascual	PS Saunders	Roman Tarasov	
Ricardo Patara	Richard Savoy	David Theese	
Dipesh Patel	John Sayer	Douglas Thompson	
Alex Parkinson	Phil Scarr	Lorin J Thompson	
Craig Partridge	Elizabeth Scheid	Joseph Toste	
Dan Paynter	Jeroen Van Ingen Schenau	Rey Tucker	
Leif Eric Pedersen	Carsten Scherb	Sandro Tumini	
Rui Sao Pedro	Ernest Schirmer	Angelo Turetta	
Juan Pena	Dan Schrenk	Phil Tweedie	
Chris Perkins	Richard Schultz	Steve Ulrich	
Michael Petry	Roger Schwartz	Unitek Engineering AG	
Alexander Peuchert	SeenThere	John Urbanek	
David Phelan	Scott Seifel	Martin Urwaleck	
Derrell Piper	Yury Shefer	Betsy Vanderpool	



Follow us on Twitter and Facebook

@protocoljournal



<https://www.facebook.com/newipj>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Diamond Sponsors



Ruby Sponsors

Your logo here!

Sapphire Sponsors

Your logo here!

Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
NMS
535 Brennan Street
San Jose, CA 95131

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

David Conrad, Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

May 2020

Volume 23, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

FROM THE EDITOR

In This Issue

From the Editor	1
Network Buffer Sizes	2
Mail Security with DMARC and ARC	21
Letter to the Editor	35
Fragments	37
Thank You!	40
Call for Papers	42
Supporters and Sponsors	43

In mid-February, I traveled to Melbourne, Australia, to attend the *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT). I normally attend around 10 or 12 similar Internet-related events in a year, be they *Network Operator Group* (NOG) conferences, *Regional Internet Registry* (RIR) events, or meetings of *The Internet Engineering Task Force* (IETF). This year, most of these events have either been cancelled or have “gone virtual” as the world tackles the COVID-19 pandemic.

The pandemic has clearly demonstrated the resilience and flexibility of the Internet and the people and organizations that rely on it for work, education, and entertainment. The various lock-downs or shelter-in-place orders have also given many of us an opportunity to take a closer look at some of the underlying technologies of the Internet, as we participate in online events or perhaps read more books and articles. This journal continues to receive many interesting articles on all aspects of networking, and in addition to the normal issues in print (and PDF format), we are also planning to expand our online presence in the near future.

Buffering is a central concept in packet-switched networks. Applications such as streaming audio or video rely on buffers to compensate for the fact that packets do not arrive at a fixed rate or even in a fixed order. Memory buffers are also used within the switches of the network to account for variations in network bandwidth and throughput. In our first article, Geoff Huston discusses network buffers and explains the numerous mechanisms that are used or have been proposed to tackle network congestion.

Previous articles in this journal have discussed various aspects of unsolicited e-mail, commonly referred to as “spam.” This time, John Levine explains recent developments in anti-spam efforts, specifically *Domain-based Message Authentication, Reporting & Conformance* (DMARC) and *Authenticated Received Chain* (ARC).

As always, we welcome your feedback and suggestions on anything you read in this journal. Letters to the Editor may be edited for clarity and length and can be sent to ipj@protocoljournal.org. Please make sure your subscription details are accurate.

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

What's the Right Network Buffer Size?

by Geoff Huston, APNIC

Packet-switched networks need to use memory buffers within the switches of the network. In a simple example, if two packets arrive at a switch at the same time and are destined to the same output port, then one packet needs to wait in a local buffer while the other packet is sent on, assuming that the switch does not want to needlessly discard packets. Not only do these buffers address such timing issues that are associated with multiplexing, they are also useful in smoothing packet bursts and performing rate adaptation that is necessary when packet sources are self-clocked. However, there is a question that has never been answered satisfactorily: What's the "right" size for the memory buffer of a switch? If buffers are generally good and improve data throughput by reducing the incidence of packet drop, then more (or larger) buffers are better, right? Not necessarily, because buffers also add additional delay to packet transit through the network if the packet gets parked into one of more buffers in transit. If you want to provide a low-jitter packet-transit service, then deep buffers in the network are decidedly unfriendly! The result is the rather enigmatic observation that network buffers have to be as big as they need to be, but no bigger.

Buffers in a packet-switched communication network serve at least two purposes. They impose some order on the highly erratic instantaneous packet rates that are inherent when many diverse packet flows are multiplexed into a single common transmission system, and they compensate for the propagation delay inherent in any congestion feedback control signal and the consequent coarseness of response to congestion events by end systems.

The Internet adds an additional dimension to this topic. Most Internet traffic is still controlled by the rate adaptation that various forms of *Transmission Control Protocol* (TCP) congestion-control algorithms use. The overall objective of these rate-control mechanisms is to make *efficient use of the network*, such that no network resource is idle when there is latent demand for more resources, and *fair use of the network*, such that if the network has multiple flows, then each flow will be given a proportionate share of the network resources, relative to the competing resource demands from other flows.

The study of buffer sizing is not one that occurs in isolation. The related areas of study encompass various forms of queueing disciplines that these network elements use, the congestion-control protocols that data senders use, and the mix of traffic on the network. The area also encompasses considerations of hardware design; *Application-Specific Integrated Circuit* (ASIC) chip layouts; and the speed, cost, and power requirements of switch hardware.

The topic of buffer sizing was the subject of a workshop at Stanford University in early December 2019. The workshop drew together academics, researchers, vendors, and operators to look at this topic from their perspectives. It hosted 98 attendees from 12 countries, with 26 from academia and 72 from industry. The following are my notes from this highly stimulating workshop^[0].

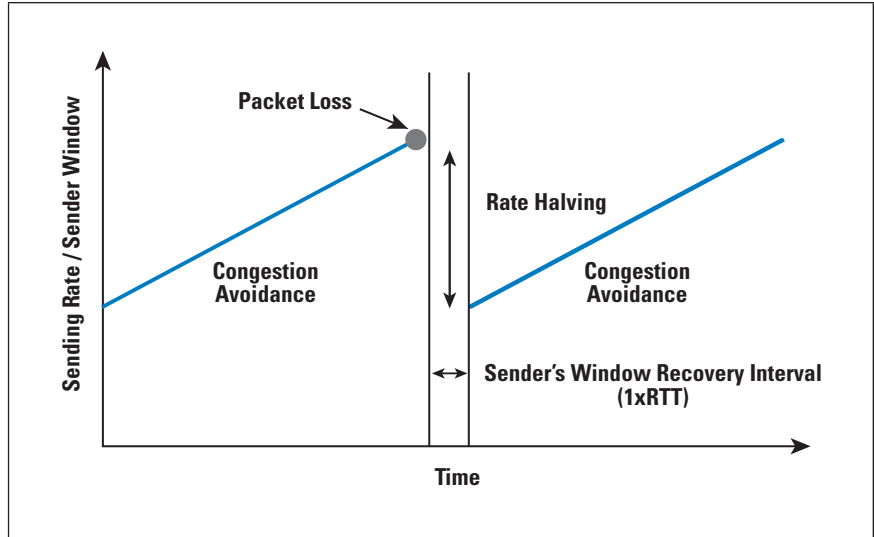
Background

In an autonomously managed packet network, packet senders learn from reflections of data that packet receivers provide, in a similar fashion to the way a radar system “learns” from a reflected signal. In a reliable flow-controlled TCP session the receiver sends an ACK packet to acknowledge the receipt of one or more data packets. Each ACK describes how many in-sequence bytes the receiver removed from the network. The sender can use this ACK signal to guide the injection of additional data into the network. One aim of each packet sender is a position of *stability*, where every packet passed into the network is matched against a packet leaving the network. In the TCP context, this behaviour is termed *ACK Pacing*.

While the sender can use ACK pacing to determine a stable sending rate, it cannot readily determine a *fair* and *efficient* sending rate. The unknown factor in this model is that the sender is not aware of the right amount of network resources to claim for the data transaction that would sustain a *fair* and *efficient* outcome. The TCP approach to solve this problem is to use a process of *dynamic discovery* where the sender probes the network by gently increasing sending rates until it receives an indication that the sending rate is too high. It then backs off its sending rate to a point that it believes is lower than the sustainable maximum rate and resumes the probe activity^[1].

This classic model of TCP flow management is termed *Additive Increase, Multiplicative Decrease* (AIMD). The sending rate is increased by a constant amount over each time interval (usually the time to send a packet to the receiver and the receiver to send an acknowledgement packet back to the sender, or a *Round Trip Time* (RTT) interval). In response to a packet-loss event, indicated by *Duplicate ACKs* that suggest the next in-sequence packet has been lost and the receiver considers successive packets to be out of order, the sender decreases the sending rate by a multiplicative ratio. The classic model of TCP uses an additive factor of 1 TCP *Message Segment Size* (MSS)^[9] of data per RTT and a rate halving (divide by 2) in response to a packet loss. The result is a “sawtooth” TCP behaviour^[2] (Figure 1). This control is determined by the sender maintaining a *Congestion Window* value, which is the maximum amount of unacknowledged data that the sender can have. Increasing the sending rate is achieved by increasing this value, and a decrease is achieved by reducing this value. When the value is reduced, then the sender must wait for the amount of unacknowledged data to drop below the new value before sending new data into the network.

Figure 1: TCP AIMD Congestion-Control Behaviour



We can mathematically model this behaviour of rate halving in response to packet loss and linear increase otherwise. If the packet-loss function is assumed to be a random loss function with a probability p , then the data-flow rate is proportional to the inverse square root of the packet-loss probability, as given in following equation (0)^[2]:

$$BW = \frac{MSS}{RTT} \cdot \frac{C}{\sqrt{p}} \quad \text{where } C = \sqrt{\frac{3}{2}} \quad (0)$$

This result implies that the achievable capacity of an AIMD TCP flow is inversely proportional to the square root of the packet-loss probability.

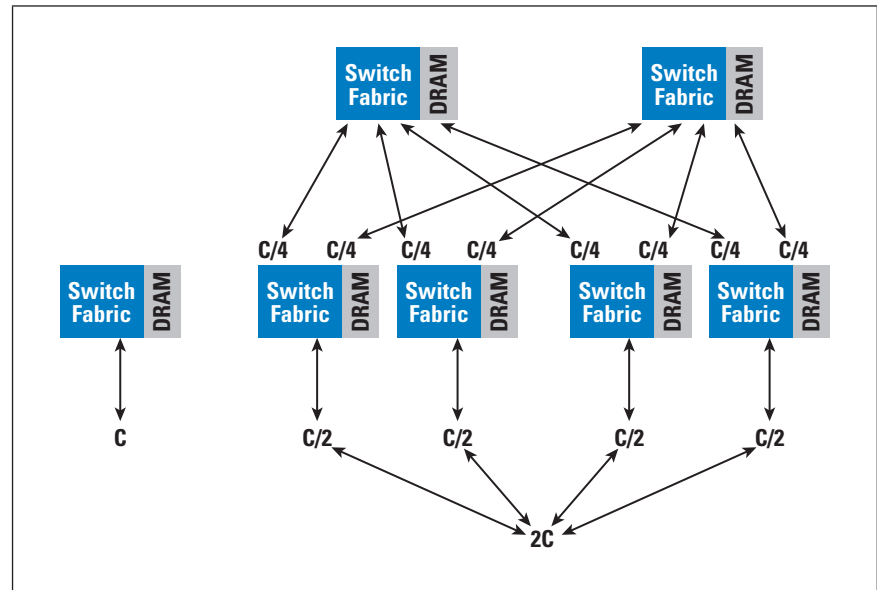
But packet loss is not a random event. If we assume that packet loss is the result of buffer overflow, then we also need to consider buffers and buffer depth in more detail. An informal standard for the Internet is that the buffer size should be equal to the delay-bandwidth product of the link (the derivation of this “rule of thumb” result is explained in the next section).

$$Size = MSS \cdot RTT \quad (1)$$

As network link speeds increase, the associated buffers similarly need to increase in size, based on this engineering rule of thumb. The rapid progression of transmission systems from megabits per second to gigabits per second and the prospect of moving to terabit systems in the near future pose, particular scaling issues for silicon-based switching and buffer systems. As networks increase in scale, the switching scaling factors tend to show multiplicative properties.

For example, if we have a single switch of capacity C and we want to double the effective switching capacity but cannot increase the capacity of the switching chip, then how many switching chips will we need to produce a composite switch of capacity $2C$? The answer is not 2 but 6, as shown in Figure 2. A packet will also need to traverse up to three switch fabrics, so the aggregate buffer size of the path through the switch fabric may triple in size.

Figure 2: Doubling Switch Capacity



Self-clocking packet sources imply that congestion events within the network are inevitable, and any control mechanism that is imposed on these sources requires some form of *feedback* that allows the source to craft an efficient response to congestion events. However, this feedback is constrained by propagation delays and this lag creates some coarseness in the response mechanisms. If the response is too extreme, the sources will over-react to congestion and the network will head into instability with oscillations between periods of intense use and high packet loss and periods of idle operation. If the response is too small, the congestion events will extend over time, leading to protracted periods of operation with full buffers, high lag, and high packet loss.

Robust control algorithms need to be stable for general topologies with multiple constrained resources, and ways of achieving this stability are still the subject of investigation and experimentation. If feedback based on rate mismatch is available, then feedback based on queue size is not all that useful for stabilising long-lived flows. Feedback based on queue size is, however, important for clearing transient overloads.

Over more than three decades of experience with congestion-management systems, we have seen many theories, papers, and experiments.

Clearly, there is no general agreement on a preferred path to take with congestion-control systems. However, a consistent factor here is the network buffer size, and the sizing of these buffers in relation to network capacity. One view, possibly extreme, says that buffers are at the root of all performance issues.

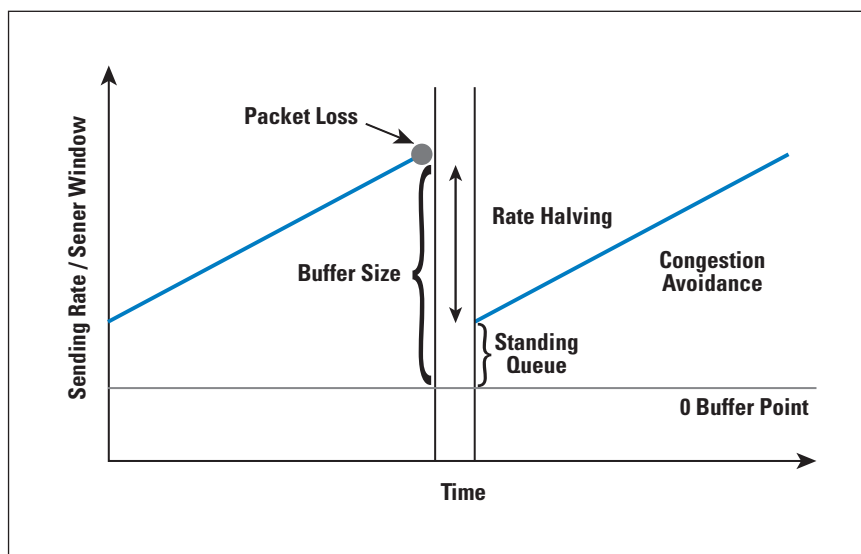
The task of dimensioning buffers in a switching system has implications right down to the design of the ASIC that implements the switch fabric. On-chip memory can be fast, but it is limited in capacity to some 100 MB or less. Larger memory buffers need to be provisioned off-chip, requiring I/O logic to interface to the memory bank, and the speed of the off-chip system is typically slower than on-chip memory. Hybrid systems have to compromise between devoting chip capacity to switching, memory, and external memory interfaces. And layered on top of these design trade-offs is the continuing need to switch at higher speed across larger numbers of ports.

One objective of the Buffer Size Workshop was to continue the conversation about buffers, determine their relationship to congestion control, and improve our understanding about the interdependence among buffer size, queuing control, self-clocking algorithms, network dimensioning, and traffic profiles.

How Big Should a Buffer Be in the Internet?

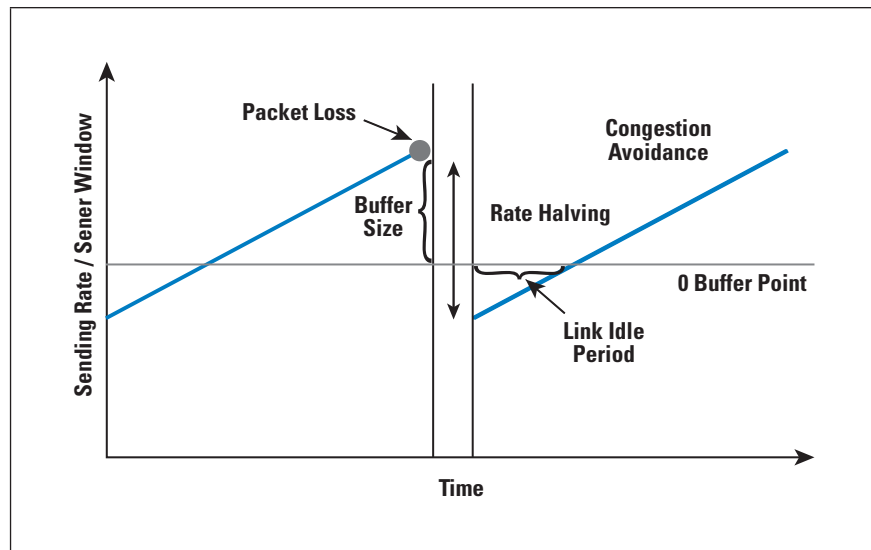
The single AIMD flow model predicts poor outcomes for flows operating across buffers that are too deep or too shallow. Too deep and the flow's loss response of halving the congestion window does not clear the buffer, and a standing queue forms that contributes to an increased latency imposed in the flow (Figure 3).

Figure 3: Deep Buffers and Rate Halving
Halving Congestion Response



If the buffer is too shallow, then rate halving drops the sending rate below the bottleneck capacity, and the link will be under-used until the additive increase brings the rate up to the link capacity (Figure 4).

Figure 4: Shallow Buffers and Rate Halving Congestion Response



The delay-bandwidth product rule of thumb generates some extremely large queue-capacity requirements in medium-delay, high-capacity systems. A 10-Gbps system using a 100-ms RTT link requires a 125-MB memory pool per port that can read and write at 10 Gbps. A 1-Tbps system would require a 12.5-GB memory pool per port that can read and write at 1 Tbps. A 16-port switch would require 200 GB of high-speed buffer memory using this same design guideline.

Such numbers are challenging for switch designers, and it is reasonable to review the original work to understand the derivation of this provisioning rule.

This model of provisioning the queue to the bandwidth-delay product is derived from a AIMD control algorithm of a single flow using an additive value of 1 segment per RTT and halving the congestion window on packet loss, coupled with the objective of using the buffer to keep the link busy during the period of window deflation.

However, something a little deeper in the oscillation of the AIMD flow-control process affects the selection of the buffer size. In the purely hypothetical situation of a single flow operating across a single switch with a lossless transmission medium, the only source of packet loss is buffer exhaustion. If the switch has no buffer at all, then the AIMD algorithm will operate in a steady state between half capacity and full capacity, leaving approximately one-fourth of the capacity unused by the flow.

The objective is to use a buffer as a reservoir to fill the transmission link while the sender pauses, waiting for the receiver's count of unacknowledged data to fall below the new congestion-window value. If the buffer size is set to the link bandwidth times the link RTT, then the buffer will be drained at the point when the sender's unacknowledged data reaches the congestion-window value and the sender can resume sending.

While this result is from a theoretical analysis of a single flow through a single link, experiments by Villamizar and Song in 1994^[3] pointed to a more general use of this dimensioning guideline in the case of multiple flows across multiple links. The rationale for this experimental observation was a supposition that synchronisation occurs across the dominant TCP flows, and the aggregate behaviour of the elemental flows was similar to a single large flow. This work was the foundation of today's common assumption that buffers in the network should be provisioned at a size equal to the round-trip delay multiplied by the capacity in order to ensure efficient loading of the link; see equation (1).

This supposition has been subsequently questioned. The scenario of a link loaded with a diversity of flows in RTT, duration, and burst profiles implies that synchronisation across such flows is highly unlikely, obviously having implications for buffer-size calculation. If there are two concurrent TCP flows, they have the same RTT, and they resonate in the increase and decrease events, then the buffer requirement will be the same for an efficient use of the network and a fair sharing of the available bandwidth. But if the increase and decrease of the two sessions are exactly out of phase, then a fair and efficient outcome would be created by a buffer size that is three-quarters of the original single flow. The real world typically sees a number of concurrent flows where both the RTT and the phase of the TCP duty cycle all vary. A Stanford TCP research group study in 2004^[4] used the central-limit theorem to point to a radically smaller model of buffer size. You can maintain link efficiency for N desynchronised flows with a buffer that is dimensioned to the size of:

$$Size = \frac{BW \cdot RTT}{\sqrt{N}} \quad (2)$$

This result is radical for high-speed extended latency links in a busy network. The consequences on router design are enormous: “For example, a 1 Tb/s ISP router carrying one TCP flow with an RTT_{min} of 100ms would require 12.5 GB of buffer and off-chip buffering. If it carries 100,000 flows, then the buffer can be safely reduced to less than 40MB, reducing the buffering and worst-case latency by 99.7%. With small buffers, the buffer would comfortably fit on a single chip switch ASIC.”^[5]

Queue Management

The default operation of a queue within a switch is to accept new packets while there is still space in the queue and discard all subsequently arriving packets until the output process has cleared space in the queue. If an incoming packet burst arrives at a switch and the queue capacity is insufficient to hold the burst, then the tail of the burst will be discarded. This tail-drop behaviour can compromise the performance of the flow, because the clocking information for the tail end of the burst has been lost.

One mitigation of this behaviour is *Active Queue Management* (AQM), where the process of queue formation triggers “early” drop. In other words, a packet drop will occur even when there is space in the queue to accept the packet. The ideal outcome of AQM is that packet drop in a large burst will occur inside the burst and the trailing packets following the dropped packet (which are not dropped as there is still space in the queue) will carry a coherent clocking signal in the ACK packet train that allows the flow to repair the loss quickly without losing the implicit clocking signal. Loss-based congestion-control algorithms will react to this packet drop by dropping their congestion-control window size, reducing their sending rate without collapsing the sending rate back to zero.

Drop-based TCP control algorithms react predictably to packet loss. However, the Internet is not entirely homogenous with respect to flow-control algorithms, and we are seeing increasing interest in flow systems that account for variance of the RTT measurements in a flow, or so-called *delay-based* TCP control systems. Delay-based paced control algorithms react differently to queue drop, and a “pure” delay-based flow-control system is indifferent to a loss signal. The question is: Are there AQM functions that can support a mix of congestion-control algorithms? Indeed, is the question of what form of AQM to use a more important question than the size of the underlying buffer?

Explicit Network Feedback

For many years there has been considerable debate between an end-system approach that uses *only* the received ACK stream to infer the network congestion state in the data forwarding direction from a packet loss signal (Figure 5), and an approach that uses some form of *explicit* signalling from the network that can directly inform the source of the network state. Very early efforts in such direct signalling through *Internet Control Message Protocol* (ICMP) *Source Quench* messages were quickly discounted because of the various issues related to its potential for *Denial-of-Service* (DoS) attacks and its inability to authenticate the messages.

The *Explicit Congestion Notification* (ECN) proposal^[6] tried to address the most obvious failing of the earlier approach by placing the congestion signal inside the end-to-end IP packet exchange. Switching elements that were experiencing the onset of local congestion load in their buffers were expected to set a *Congestion Experienced* bit in the IP packet header of packets that were contributing to this load condition. Receivers were expected to translate this bit into the ACK packet header, so that the sender received an explicit congestion signal rather than having to infer congestion from an ACK signal that reflects packet loss (Figure 6).

Figure 5: Loss-Based Congestion-Control Behaviour

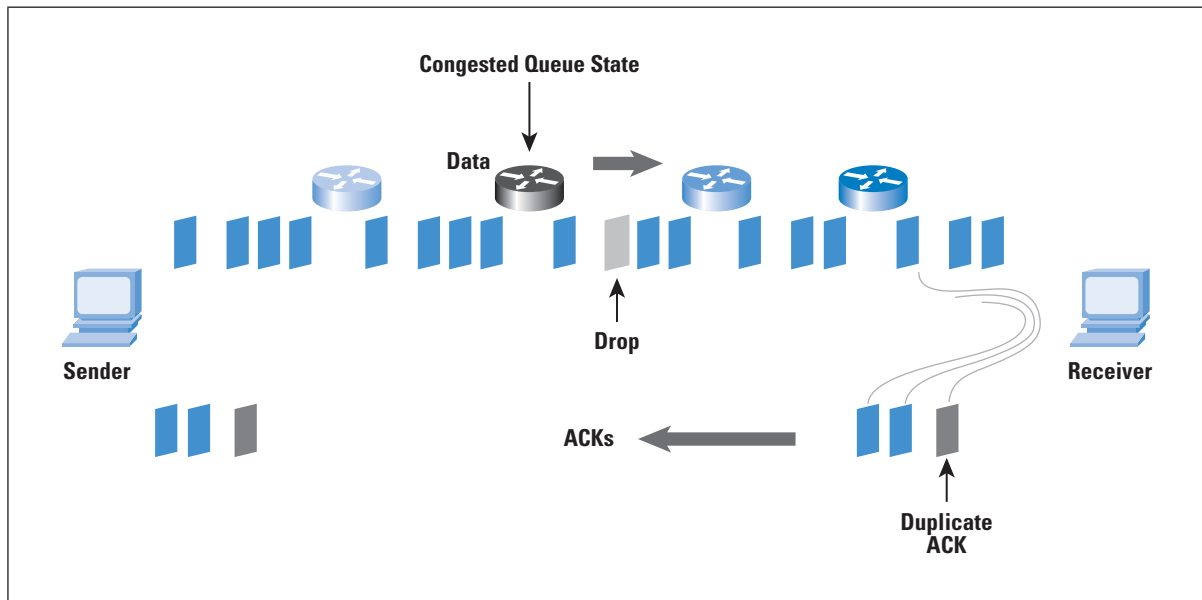
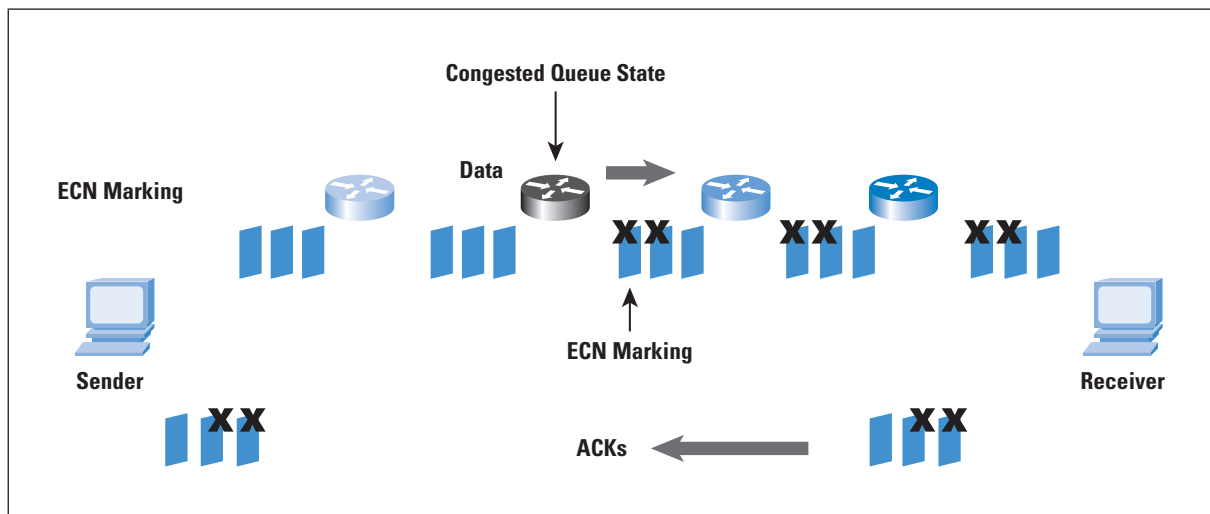


Figure 6: ECN Marking



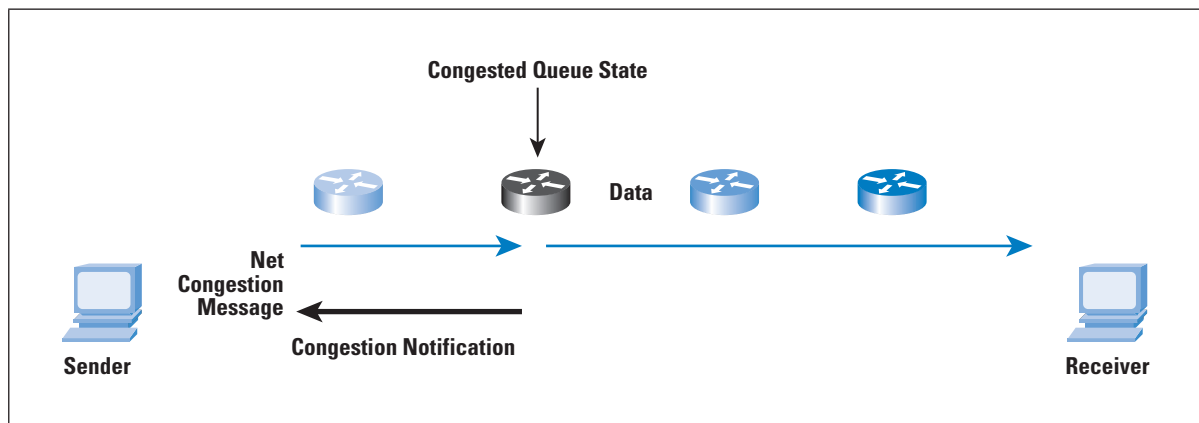
The advantage of ECN is that the sender is not placed in the position of being informed of a congestion condition well after the condition has occurred. Explicit notification allows the sender to be informed of a condition as it is forming, so that it can take action while there is still a coherent ACK pacing signal coming back from the receiver (that is, before packet loss occurs). This measure mimics the intention of delay-based flow systems, but with increased precision assuming that all switches were to perform this congestion marking.

However, ECN is only a single bit marking. Is that enough? Would a richer marking framework facilitate a more precise sender response? What if we had a marking regime that marks based on the distance from the current rate to a desired fair-efficient rate? Or use a larger vector to record the congestion state in multiple queues on the path?

The conclusion from one presentation is that the single-bit marking, while coarse and non-specific, is probably sufficient to moderate self-clocking TCP flows such that they do not place pressure on network buffers, leaving the buffers to deal with short-term bursts from unconstrained sources.

Another presentation at the workshop explored a network-level direct-feedback message, analogous to the ICMP *Packet Too Big* messages in *Path MTU Discovery* (PMTUD). To short-circuit the delays associated with completing the entire round trip, this approach envisages the switch experiencing the onset of congestion to explicitly message the source of this congestion condition (Figure 7).

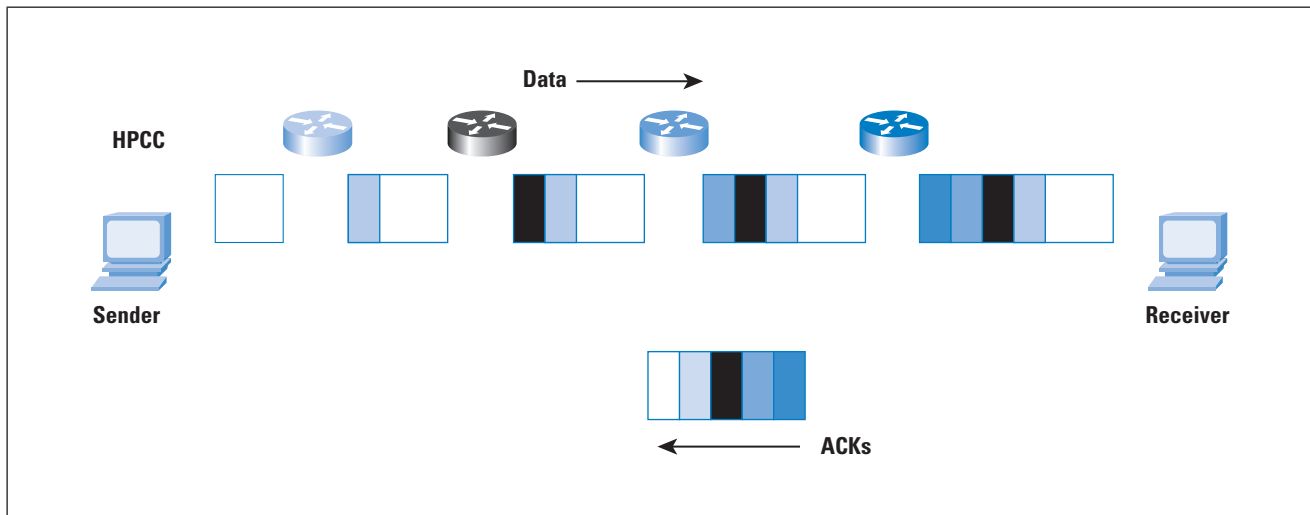
Figure 7: Network-Congestion Signalling



Another presentation looked at the attachment of a detailed telemetry log to each packet in a data-centre application. In the *High-Precision Congestion Control* (HPCC) framework each switch attaches the time, queue length, byte count, and link bandwidth to the data packet. The receiver takes this data and attaches it to the corresponding ACK, so that the sender can form a detailed model of the recent state of path capability. HPCC allows the sender to calculate a fair sending rate and then rapidly converge to this rate, while at the same time bounding the formatting of queues and bounding queuing delays. The domain of application of this approach appears to be the data centre, and the objective is to achieve high speed with bounded delay for *Remote Direct Memory Access* (RDMA)-style applications (Figure 8).

There is a degree of debate between *congestion-based* TCP control and delay-based mechanisms. On the one hand, we hear that delay-based mechanisms can operate the flows at the onset of queue formation in the network. On the other hand, we hear that attempting to set the flow to a fixed delay and operating with fairness to other flows is intrinsically impossible and that we need to operate flows with congestion moderation.

Figure 8: High-Precision Congestion Control



Near and Far Buffers

What is the cost/power trade-off of buffers on-chip and off-chip? And if we are considering off-chip, what do we actually mean, because there are different implementation approaches to off-chip memory. As a general observation, the performance of off-chip memory is not remotely close to what is required by a high-capacity, high-speed switch. This performance is not improving over time because memory speed is not scaling at the same rate as transmission or switch speeds, so the gap in performance between transmission and switching and memory speed is only getting larger over time.

One switch chip fabricator, Broadcom, implements both deep and shallow buffers. On its switch fabric chip Broadcom uses small, fast buffers and wraps the switch fabric with everything it can to reduce the dependency on deep buffers.

Recent operational data at Intel suggests that shallow buffers may be “good enough,” but because of limitations in instrumenting technologies there is insufficient confidence in these results to allow switch chip designs to completely discount external memory interfaces and a local cache and use on-chip memory exclusively. Current switch designs use between 10% and 50% of chip area on memory management. This observation applies to high-capacity, high-speed switches, because at lower capacity and lower speeds there is no such constraint and you can use large pools of off-chip memory (relative to transmission speeds), although some constraint in the amount of memory will likely produce a better outcome in these contexts as well.

The question of future requirements is always present in chip design, given the long times between phrasing requirements and deployment into networks. Where is this situation heading? Memory buffers are not growing as fast as chip bandwidth. Clock speeds are not increasing, and scaling chip bandwidth is currently achieving parallelism rather than increasing the chip clock speed.

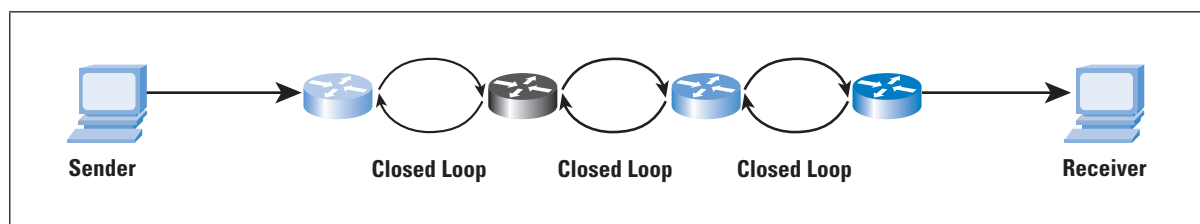
While doubling switch capacity may be feasible, contemplating an increase in capacity by factors of 20, 50, or even 100 seem like particularly tough challenges.

Today packet rates are typically achieved with multiple parallel pipelines, and orchestrating such highly parallel mechanisms creates its own complexity in design. No one is yet prepared to call an end to the prodigious outcomes of *Moore's Law* in the semiconductor realm, but it looks like clock speeds are not keeping up, and pin density and even increasing gate density are becoming challenging. Is doubling the number of ports on a switch chip good enough? If a chip has twice the switch bandwidth, does it need twice the on-chip memory capacity? Or less? The answer lies in external factors such as congestion-control algorithms, queue-management disciplines, and delay management.

Hop-by-Hop Flow Control

Hop-by-Hop Flow Control represents a revival of a very early attribute of packet-switched networks, where an end-to-end path is composed of a sequence of flow-controlled hops. Each switching element sends at its line rate into the buffer of the next switch. When a queue forms at the receiver, the hop flow control can pause the flow coming from the adjacent switch and resume it when the queue is cleared. Yes, this process sounds very reminiscent of X.25 and the *Digital Data Communications Message Protocol* (DDCMP) component of DECnet, and it's the opposite of the intent of the end-to-end approach. However, the approach can produce direct back pressure on a bursting source with no packet drop and yield highly efficient use without extensive buffer-induced delays. Essentially the self-clocking nature of the flow is replaced with a network clocking function (Figure 9). Admittedly, this approach is not universally applicable, and it appears to offer a potential match to the intra-data-centre environment where traffic patterns are highly bursty, propagation times are low, paths are short, and volumes and speeds are intense.

Figure 9: Hop-by-Hop Flow Control



Flow-Aware Buffer Management

It appears that the move towards shorter buffers relative to the link speed is inevitable. But how to manage the feedback systems to allow self-clocking data sources to adjust to the shortened buffer space is still an outstanding issue.

One approach starts with a basic traffic characterisation of a relatively small proportion of “elephant flows” (high volume, long duration) mixed with a far higher count of “mice” flows (low volume, short duration). While elephant flows are highly susceptible to congestion signalling, mice flows are not.

If the network could classify all currently active flows into either elephants or mice, then the network could use different queuing regimes for each traffic class. This sorting adds to the cost and complexity of packet switches, and if scaling pressures are a factor in switch design, then it’s not clear that the additional cost of switch complexity would be offset by a far superior efficiency outcome in the switching function.

Assuming that such a flow classification could be achieved dynamically, we can consider differential responses. For short flows, there is little benefit to be gained by any form of explicit congestion control other than placing all such flows into their own queuing regime. For long-lived large flows, we could contemplate an explicit network-congestion signal. It could take the form of an explicit packet back to the network-generated source. The advantage of this approach is that the feedback of excessive sending rate is faster than a full RTT interval, allowing the sender to give a timely response. However, this idea does seem like a reprise of the ill-fated ICMP Source Quench message, and all that was problematical with ICMP Source Quench is probably still an issue in this form of network-congestion notification.

We can exploit this concept of the use of various queue regimes for different flow types in a different way by using a short buffer for long flows in the expectation that the implicit congestion signal of packet drop would allow the long-duration flow to stabilise into the available network resource, while short unregulated bursts could have access to a deeper buffer, allowing effective use of the buffer as a rate-adaptation tool to mitigate the burst.

This concept is taken even further in one project, which used the observation that if a buffer is too deep, then the flow-rate reduction following packet drop will leave a standing queue in the buffer, and if the buffer is too shallow, then the rate reduction will leave a period of an empty queue and an idle transmission system. This observation means that a flow-aware buffer manager could adjust its buffer size following observation of the post-reduction behaviour, reducing the buffer if standing queues form and increasing it if the queue is idle. It’s an interesting approach to fair-queuing flow management, treating the per-flow buffer as an elastic resource that can resize itself to adapt to the congestion-management discipline of the flow.

ISP Network Buffer Profile

P4 is a language used to program the data plane of network devices. The language can express how a switch should process packets (“P4” itself comes from the original paper that introduced the language, *Programming Protocol-independent Packet Processors*^[7].)

Barefoot's *Tofino* is an example of a new class of programmable Ethernet packet switches that are controlled through P4 constructs, and these units can currently handle aggregate capacity of some 12.8 Tbps of data-plane capacity. This capability allows for a measurement regime that can expose packet characteristics at a nanosecond level of granularity. By tapping the packet flow of a high-speed trunk transmission system into and out of a switching element in the network and attaching the taps to a P4 switch unit, it is possible to match the times of ingress and egress of individual packets and generate a per-packet record of queuing delay within the switching element at a nanosecond level of granularity.

This capability provides a new level of insight into burst behaviour in high-speed carriage systems *Internet Service Providers* (ISPs) use. The major observation from an exercise conducted on a large ISP network was that network buffers are lightly used except for “microbursts,” bursts of some 100 microseconds or so, where the queue adds a delay element of more than 10 ms on a 10 GigE port. Further analysis reveals an estimate of packet drop rates if the network buffers were reduced in size, and for this case the analysis revealed that an 18-msec buffer would be able to sustain a packet drop rate of less than 0.005%.

If buffer-congestion behaviours in such ISP networks are, in fact, microbursts, then network measurement tools that operate at the per-minute or even at the per-second level of granularity are simply too crude. P4-based measurements that can resolve behaviours at the nanosecond level offer new insights into buffer behaviours in networks that carry a large volume of diverse flows. Even though the per-flow control cycle of the data-plane flows is of the scale of some milliseconds and longer, the microburst behaviour is that of a load model that exhibits sub-millisecond burstiness. The timescale of end-to-end congestion control operates at a far coarser level than the observed behaviour of congestion within a switch running a conventional traffic load.

This discussion leads to the observation that large-scale systems are creating extremely rapid queue size fluctuations, and it is unrealistic to expect that end-to-end control algorithms can control the queue size. It might be that at best these control algorithms can contribute to influencing the distribution of queue sizes.

Sender Pacing

The Internet can be seen as a process of statistical multiplexing of a collection of self-clocked packet flows where the flows exhibit a high degree of variance and a low level of stability. The reaction to this unconstrained input condition so far is to use large buffers that can absorb the variations in traffic. How large is “large enough” becomes the critical question in such an environment. The work on buffer sizing as being in proportion to the bandwidth-delay product of the transmission elements is an outcome of a process that measures the properties of the control algorithm for traffic flows.

It then derives estimates of buffer sizes that should be capable of carrying such a volume of traffic that it will efficiently and fairly load the transmission system.

The exercise assumes that the buffer dimension is a free parameter in network design, and control algorithms are fixed. Buffer speed inside the network has to double at a cycle of some 2 years, and the buffer size has to double in a similar timeframe. The product of size and speed is a quadrupling every 2 years. The current tactical response to this escalation of buffer requirements due to transmission capacity increases has been to reduce the size of the buffers relative to the transmission capacity. However, this response is not a long-term sustainable solution because such under-provisioned network buffers will impair overall network efficiency in these self-clocking flow regimes.

The future prospects for self-clocked traffic flows are not looking all that bright given that the growth demands for network buffer-based mitigation of unconstrained sender behaviours appears to be in excess of what can be satisfied within constraints of constant unit cost of network infrastructure. Without overall economies of scale where larger service-delivery systems achieve lower unit costs of service delivery, the management of traffic and content assumes a different trajectory that tends to drive towards greater distribution and dispersal rather than continued aggregation and amalgamation. For the large hyper-scaled content enterprises in today's Internet, this outcome is certainly not optimal.

It is a potentially fruitful thought process to consider this topic from an inverted perspective and look at the desirable control-algorithm behaviour that efficiently uses the network transmission resources when the available buffering is highly restricted. This thought process leads to the consideration of "pacing," where the server uses high-precision timers to smooth data flows as they leave the server, attempting to create a stable traffic flow that matches bottleneck capacity on the path. The more accurate this estimation of bottleneck bandwidth, the lower the demand for buffer capacity due to burst adaptation. Residual buffer demand is presumably based on the demands of statistical multiplexing of disjoint flows. Given that the senders are under the control of the service-delivery platforms and there are orders of magnitude fewer high-volume senders than receivers, this form of change is actually far less than the change required by, say, the IPv6 transition.

It is this thinking that lies behind the *Bottleneck Bandwidth and Round-trip Propagation Time* (BBR) protocol work. The sender's flow-control algorithm generates an estimate of the bottleneck bandwidth and the minimum RTT interval, and then paces packet delivery so as to feed traffic into the bottleneck at exactly the bottleneck capacity, which should not involve the formation of a queue at the bottleneck.

The BBR control algorithm periodically probes up to revise its previous bandwidth estimate, and probes down to revise its previous minimum RTT, and accounts for other congestion-formation signals, such as ECN. This probing up and down, or dithering, is not precisely specified in the core BBR algorithm, and these parameters are being revised in the light of deployment experience to determine dithering settings that are both efficient and fair. The expectation is that BBR will not drive the formation of standing queues in the network and will pace the flow at the maximal rate that the network path can fairly sustain.

However, BBR is not the only way to perform flow pacing, and a large number of outstanding questions remain. How does pacing at the sender affect the queue management at the edge close to the client? What are the cross impacts of burst traffic with pacing? How should a pacing-control algorithm react to packet loss? Or to out-of-order packet delivery? Can strict flow pinning still be required for *Equal Cost Multiple Path* (ECMP) routing or does pacing relax such requirements for strict path pinning? Are pacing or self-clocking the only options, or are there other approaches?

One perspective is that we are sitting between two constraint sets. Escalating volume and speed in the core parts of the network implies that bandwidth-delay product model buffer sizing is an unsustainable approach. The scaling back of buffer sizes in the network means that self-clocked protocols will potentially become more unstable and compromise achievable network efficiency and fairness. From this perspective sender pacing looks to be a promising direction to pursue.

Is There a Buffer Sizing Problem?

In the Internet we are currently seeing a diversity of responses to network provisioning. Some network operators use equipment with generous buffers. These buffers are overly generous according to the buffer-bloat argument. Other network operators field equipment with scant buffers that run the risk of starving data sources while leaving idle network capacity.

There is a mix of congestion-control algorithms (CUBIC, NewReno, BBR, *Low Extra Delay Background Transport* [LEDBAT], etc.) and a mix of queue-management regimes (*Controlled Delay* [CoDel], *Random Early Detection* [RED], *Weighted Random Early Detection* [WRED])^[8]. A diversity of deployment environments exists, including mobile networks, Wi-Fi, wired access systems, LANs, and data centres. And there is a mix of parameters of the desired objective here, whether it is some form of fairness, loss, jitter, start-up speed, steady-state throughput, stability, efficiency, or any combination of these factors. It is little wonder that it's challenging to formulate a clear picture of common objectives and to determine what actions are needed to achieve whatever we might want!

There is the assumption that large network buffers absorb imprecision in clocking (timing “slop”) and allow simpler coarse rate-control algorithms to operate effectively without needing high-precision tuning. Small network buffers provide little leeway and tolerance for such approximate approaches. This mistrust of the level of precision of control that end systems exercise is a pervasive view within the networking community, and it could even be characterized as an entrenched view. So entrenched is this view that probably no experimental result could convince the community as a whole that network buffers can be far smaller than they are today, all other factors being equal. This fact is true despite the overwhelming evidence that overly large buffers compromise network performance, a position that has been described as “buffer bloat.” There are other reasons why large buffers are a problem for networks and users. As we scale up the size and speed of the network, large very-high-speed buffers are also increasing in cost. If we are going to admit compromises and trade-offs in network design, is reducing the relative size of the buffer an acceptable trade-off?

And if we want to reduce buffer size and maintain efficient and fair performance, how can we achieve it? One view is that sender pacing can remove much of the pressure on buffers, and self-clocking flows can stabilise without emitting transient bursts that buffers will need to absorb. Another view, one that does not necessarily contradict the first, is that the self-clocking algorithm can operate with higher precision if there were some form of feedback from the network on the state of the network path. This feedback can be as simple as a single bit (ECN) or a complete trace of path element queue state (HPCC).

This topic remains a rich area of unanswered questions. What does it imply when the timescale of buffer-congestion events are orders of magnitude smaller than the timescale of self-clocking flows? Are flows overly reliant on loss signals and too insensitive to delay variation? Can paced delay-based algorithms like BBR coexist with loss-based oscillating algorithms such as CUBIC and NewReno? Would the general adoption of sender pacing change the picture of buffer sizing in the Internet?

How big should buffers be in the network? Or perhaps the opposite is the more practical question: How small can we provision buffers in an increasingly faster and larger network and still achieve efficient and fair outcomes in a variety of deployment environments?

All of these questions are good and legitimate for further research, experimentation, and measurement.

References and Further Reading

- [0] Workshop on Buffer Sizing, Stanford University, December 2–3, 2019. <https://buffer-workshop.stanford.edu>
- [1] Van Jacobson and Mike Karels, “Congestion Avoidance and Control,” *ACM SIGCOMM Computer Communications Review*, Volume 18, Issue 4, August 1988.
This foundational paper is frequently cited by many papers on TCP behaviour.
- [2a] Matt Mathis, Jeffrey Semke, Jamshid Mahdavi, and Teunis J. Ott, “The Macroscopic Behaviour of the TCP Congestion Avoidance Algorithm,” *ACM SIGCOMM Computer Communications Review*, Volume 27, Issue 3, July 1997.
Matt and Jamshid published a second paper of this topic in October 2019, where they argue that this macroscopic model will soon be completely obsolete:
- [2b] Matt Mathis and Jamshid Mahdavi, “Deprecating the TCP Macroscopic Model,” *ACM SIGCOMM Computer Communications Review*, Volume 49, Issue 5, October 2019.
- [3] Curtis Villamizar and Cheng Song, “High Performance TCP in ANSNET,” *ACM SIGCOMM Computer Communications Review*, Volume 24, No. 5, October 1994.
An effort to generalise the buffer sizing theory into observed practice in networks. It has been commonly acknowledged as the rationale for using bandwidth-delay product as the buffer sizing model for network equipment. This small-scale study was within a single network, and the results have been applied in a far more diverse set of deployment scenarios than the single setup that was analysed in this paper.
- [4] Guido Appenzeller, Isaac Keslassy, and Nick McKeown, “Sizing Router Buffers,” *ACM SIGCOMM Computer Communications Review*, Volume 34, Issue 4, September 2004.
A widely cited paper that provides an analysis of multiple diverse flows over a single common buffer, concluding that an efficient and fair buffer size model is related to the inverse of the square root of the number of active flows that traverse this common link (and buffer).
- [5] Nick McKeown, Guido Appenzeller, and Isaac Keslassy, “Sizing Router Buffers (Redux),” *ACM SIGCOMM Computer Communications Review*, Volume 49, No. 5, October 2019.
A very recent review of the buffer sizing conversation and highlighting some of the significant experiments with small buffers in large networks since the 2004 paper. The paper includes numerous questions about future requirements for buffer sizes.

- [6] Sally Floyd, K. K. Ramakrishnan, and David L. Black, “The Addition of Explicit Congestion Notification (ECN) to IP,” RFC 3168, September 2001.

The specification of re-purposing two bits in the IPv4 packet header for routers to use to mark congestion events into active flows.

- [7] Pat Bossharty, Dan Daly, Glen Gibb, Martin Izzard, Nick McKeown, Jennifer Rexford, Cole Schlesinger, Dan Talayco, Amin Vahdat, George Varghese, and David Walker, “Programming Protocol-independent Packet Processors,” ACM SIGCOMM *Computer Communications Review*, Volume 44, No. 3, July 2014.

The specification of a programming language for packet processors that has been used in recent very-high-speed packet-switch processors.

- [8] “TCP Congestion Control,” Wikipedia,
https://en.wikipedia.org/wiki/TCP_congestion_control

- [9] Geoff Huston, “MSS Values of TCP,” *The Internet Protocol Journal*, Volume 22, No. 3, December 2019.

- [10] Geoff Huston, “Buffers and Protocols,” Presentation at RIPE 80 Meeting, May 12, 2020. Slides and video available at:
<https://ripe80.ripe.net/archives/video/316/>

GEOFF HUSTON, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

Mail Security with DMARC and ARC

by John Levine

Electronic mail is both one of the most useful services of the Internet and the most frustrating. The best thing about mail is that anyone can send a message to anyone else without prearrangement, while the worst thing about mail is that anyone can send a message to anyone else without prearrangement. As mail became ever more ubiquitous in the 1990s and 2000s, an increasing fraction of it was mail the recipients didn't want. In 2005, Dave Crocker^[1] and John Klensin^[2] wrote articles in this journal about the spam problem. Since then, several of the anti-spam techniques described in Crocker's article have become ubiquitous as the spam problem has become worse.

One can distinguish between *spam*, unsolicited mail sent in bulk, and *phishing*, mail sent to trick the recipient into revealing account credentials or other private information. (Some phishes are sent in bulk, some are sent to specific victims, known as *Spear Phishing*.) Starting in 2007, PayPal, which had long been among the biggest phishing targets, started working with some large consumer mail systems to keep PayPal phishes out of recipient mailboxes. The idea was that the recipient systems could identify genuine e-mail from PayPal, and reject anything else purporting to be from PayPal. In 2012 an industry group started the DMARC project to generalize this technique. In 2015 the *Domain-based Message Authentication, Reporting & Conformance* (DMARC) specification was published as an independent track RFC^[3], and now it is ubiquitous in large mail systems.

DMARC works by tying the address in the RFC 5322^[4] **From:** header of a message to mail authentication, and letting a domain offer policy advice to mail recipients. If a message is successfully validated by *Sender Policy Framework* (SPF) or *DomainKeys Identified Mail* (DKIM), and the domain in that validation matches the one in the **From:** header, the message is DMARC "aligned." Sending mail systems can publish DMARC policy records in the *Domain Name System* (DNS) requesting recipient systems to quarantine (send to the spam folder) or reject unaligned mail. This system works quite well for the original intended application of DMARC, business-to-consumer mail, where the sending organization generally has full control over all mail sent from its domain. It works particularly well for PayPal, where all the mail is some variation of "log into your account to see what's new," so if a few messages are accidentally lost because DMARC miscategorizes a legitimate message it's not a big problem.

Underlying and Previous Work

DMARC depends on two existing mail authentication schemes, *Sender Policy Framework* (SPF)^[5] and DKIM^[6]. SPF does path validation of the domain in the RFC 5321^[7] **MAIL FROM:** address. A domain can publish an SPF record that uses a complex syntax to specify a set of IP addresses.

If the message was sent from one of those addresses, SPF validation passes. (This description is oversimplified; see [5] for the full details.) SPF has the virtue of being easy to implement because it requires no changes to outgoing messages and a single DNS record to implement, but it can describe only a limited subset of the ways that mail is delivered. For the most part it can handle only mail sent directly from sender to recipient, without any forwarding or remailing, and does not deal well with mail sent by third parties on the sender's behalf. While SPF provides a **-all** code that advises recipients to reject mail from the domain if SPF validation fails, most mail systems disregard the advice because the false positive rate is so high.

DKIM does message content validation by adding cryptographic signature headers to a message that the recipient can check using a key in the DNS. Each signature is stored in a **DKIM-Signature** header field that contains several subfields, including the name of the domain that added the signature. If a DKIM signature validates, it means both that the message hasn't been modified since it was signed and that the domain in the signature takes responsibility for the message. Since DKIM validates the contents of the message rather than the path, it is unaffected by forwarding.

DKIM is considerably harder to implement than SPF because it requires modifications to mail software to add the signature headers to each outgoing message. It also requires the signing system to create public/private key pairs, publish the public key in the DNS, and configure the private key into the signing software. DKIM validation fails when a message is modified in transit, such as when a mailing list adds a subject tag or a message footer, and sometimes simply because a *Message Transfer Agent* (MTA) hasn't been configured to add the signature in the first place. (In large enterprises it can be remarkably hard to track down all of the computers sending mail. DMARC helps address this problem, as we will see later.)

DKIM had an optional add-on called *Author Domain Signing Practices* (ADSP)^[8], which was sort of a proto-DMARC. A domain could publish an ADSP record in the DNS saying that if a message with the domain in the **From:** header didn't have a valid DKIM signature from the same domain, recipients should discard the message. ADSP was never deployed beyond experiments, and the *Internet Engineering Task Force* (IETF) has since made it an historic specification.

DMARC Deployment

One of the reasons that previous sender policy approaches like SPF **-all** and ADSP failed is that there is no way to test them other than turning them on to see what happens. For a small domain with one or two mail servers that might be possible, but for a large organization the risk is impossibly high since they rarely have complete knowledge of all of the systems sending their mail and how those systems are configured.

DMARC offers a variety of features to check the alignment of the mail of a domain before publishing any policy advice. It has powerful reporting features that let a domain owner ask other systems to send reports about the mail purporting to be from the domain. Domains invariably ask for reports before publishing any policies, so they can see what mail they send is and isn't aligned. This information lets them fix alignment issues before they do publish policies.

DMARC Validation

When a message arrives, DMARC validation involves first finding a DMARC policy record for the **From:** header domain, then validating SPF and the DKIM signature(s) on the message, and then perhaps doing something to the message. The first step is to find the policy record for the **From:** domain of the message, a DNS TXT record. If that domain is **marketing.mybiz.example**, the first place to look is **_dmarc.marketing.mybiz.example**. If there is a TXT record in DMARC syntax, for example, it starts with **v=DMARC1**; that's the policy record. If not, it looks for a policy record in the "organizational domain."

The DMARC specification is deliberately vague about how to find the organizational domain, but in practice everyone uses the Mozilla *Public Suffix List* (PSL)^[9] where the organizational domain is the superdomain just below the public suffix. In this case if the domain were a typical *Top-Level Domain* (TLD) that accepts registrations at the second level, the organization domain would be **mybiz.example**, so it would look for a TXT record at **_dmarc.mybiz.example**. If there is a TXT record in DMARC syntax, that's the policy record; otherwise there is no policy record for this domain.

The DMARC record is a list of key=value pairs, with rules for checking alignment, what to do with unaligned mail, and where to send aggregate and failure reports. A typical record might be:

```
v=DMARC1; p=none; rua=mailto:dmarc-a@example.net;  
ri=3600; ruf=mailto:dmarc-f@example.net
```

There is no policy (**none**), abuse and failure reports are mailed to the given addresses (**rua** and **ruf**), and the requested report interval is an hour (**ri** is 3,600 seconds.) The point of the second check for the organizational domain is twofold. First, the second check makes it easier to deploy DMARC across a large enterprise, since one DMARC organizational record can cover all of an organization's subdomains. The other is that it covers non-existent subdomains of the organizational domain, for when hostile or buggy mailers send mail purporting to be from such a subdomain.

The next step in validation is to check whether the **From:** header domain is aligned with the SPF identity of the message. The SPF validation process can produce a result of *None*, *Neutral*, *Pass*, *Fail*, *Softfail*, *Temperror*, or *Permerror*. For DMARC alignment, only a *Pass* result is acceptable.

The DMARC policy record can require strict SPF alignment, meaning the **From:** domain and SPF identity have to be the same, or relaxed SPF alignment, meaning they need only be in the same organizational domain. In the previous example, if the **From:** domain were **marketing.mybiz.example**, an SPF identity of **mail.mybiz.example** or just **mybiz.example** would be sufficient for relaxed SPF alignment. Relaxed alignment is the default.

Next, the validator checks for DKIM alignment. For each valid DKIM signature on the message, the validator compares the **From:** domain to the **d=** domain of the signature. The policy record can specify strict or relaxed DKIM alignment, again requiring either an identical signature domain or just one in the same organizational domain. If at least one valid DKIM signature is aligned, the message is DKIM aligned. If the message is either SPF or DKIM aligned, it is DMARC aligned.

If the message is aligned, we're done other than perhaps saving some statistics for later reporting. If it's not aligned, the situation is potentially much more complex if the recipient system opts to follow the policy advice, as most mail systems (at least by mail volume) now do.

The policy record can specify policy advice of **none**, **quarantine**, or **reject**. It can also specify an optional percentage of how often to apply the policy. Advice of **none** means to do whatever the recipient would have done with the message anyway. Advice of **quarantine** means to treat the message extra skeptically, perhaps by filing it in a spam folder or marking it as suspicious. Advice of **reject** asks the recipient to **reject** the message at the end of the *Simple Mail Transfer Protocol* (SMTP) session and not handle it further. If the percentage is less than 100, the advice is to treat that percentage of unaligned mail from the domain according to the advice, and the rest one step less. For example, if the advice were **reject** and the percent was 25, one-fourth of unaligned mail would be rejected and the other three-quarters would be quarantined. (The percent has no effect if the policy is **none**.)

As noted previously, the point of the percentage is to allow domain owners to enable policies gradually, see what happens, and limit the damage from misconfigurations.

DMARC Reporting

DMARC has two powerful reporting features. A domain can ask for daily aggregate reports of what IP addresses have sent mail with the domain in the **From:** header, with details about DMARC alignment and DKIM and SPF validation. Many large mail systems including Google, Yahoo/AOL, Comcast, and Fastmail send aggregate reports.

It is also possible to request copies of messages that fail DMARC validation, but for privacy reasons very few systems do. The only significant mail system in the U.S. that sends failure reports is LinkedIn.

Even for a site that has no plans to publish a DMARC policy, the reports are useful and interesting. They can provide insight into where your mail is actually going, and who else might be sending mail purporting to be from you.

To request each kind of report, the domain policy record includes a tag with a list of *mailto: Uniform Resource Identifiers* (URIs), each with an optional size limit of the maximum message report size the system can handle. The default aggregate report interval is once a day.

Aggregate reports constitute an *Extensible Markup Language* (XML) file attached to an e-mail message in gzip or ZIP compressed form. The XML file includes a section (a “record”) for each sending IP address, with subsections (a “row”) for each combination of authentication results. For example, here’s a section of a report Google sent to my *Smail* system describing mail it received from two IP addresses:

```
<record>
  <row>
    <source_ip>2001:470:1f07:1126:0:43:6f73:7461</source_ip>
    <count>1</count>
    <policy_evaluated>
      <disposition>none</disposition>
      <dkim>pass</dkim>
      <spf>pass</spf>
    </policy_evaluated>
  </row>
  <identifiers>
    <header_from>taugh.com</header_from>
  </identifiers>
  <auth_results>
    <dkim>
      <domain>iecc.com</domain>
      <result>pass</result>
      <selector>k1912</selector>
    </dkim>
    <dkim>
      <domain>taugh.com</domain>
      <result>pass</result>
      <selector>k1912</selector>
    </dkim>
    <spf>
      <domain>taugh.com</domain>
      <result>pass</result>
    </spf>
  </auth_results>
</record>
```

```

<record>
  <row>
    <source_ip>209.85.220.55</source_ip>
    <count>4</count>
    <policy_evaluated>
      <disposition>none</disposition>
      <dkim>fail</dkim>
      <spf>fail</spf>
    </policy_evaluated>
  </row>
  <identifiers>
    <header_from>taugh.com</header_from>
  </identifiers>
  <auth_results>
    <dkim>
      <domain>googlegroups.com</domain>
      <result>pass</result>
      <selector>20161025</selector>
    </dkim>
    <spf>
      <domain>googlegroups.com</domain>
      <result>pass</result>
    </spf>
  </auth_results>
</record>

```

The first record for the IPv6 address reports on a message sent from my mail server. It has a valid SPF and two valid DKIM signatures, one with the **From:** header domain and one for the server domain, so it was DMARC aligned. The second record describes four messages with valid SPF and DKIM signatures, but with SPF and DKIM domains that don't match the **From:** header, so they wouldn't have been DMARC aligned. Since the second group of messages have **googlegroups.com** authentication identifiers, they're probably the same message, modified and remailed to a Google Groups mailing list. [Since I know I sent only one message to the list that day, this message leaks the number of *Gmail* subscribers to the list. I've seen similar leakage for much larger lists; for example, like the one operated by the *North American Network Operators' Group* (NANOG).]

Larger mail systems receive reports with larger numbers of messages and more report sections. The reports are intended to be mechanically handled. Some open source software is available to analyze reports and put summaries in a database^[10]. More often the reports are sent directly to specialist services like *Dmarcian*^[11] or *Agari*^[12] that offer freemium report analysis services, simple analysis for free, or more sophisticated analysis and remediation advice for a fee.

The other kind of report is a failure report. When a message arrives that has the domain address in the **From:** header and fails DMARC validation, the recipient system may (but usually doesn't) send the message back in a failure report. The report is a multipart/report e-mail message containing a structured report section and a full or partial copy of the failing message.

A typical report section follows:

```
Feedback-Type: auth-failure
User-Agent: Lua/1.0
Version: 1.0
Original-Mail-From: nanog-bounces@nanog.org
Original-Rcpt-To: xxx@linkedin.com
Arrival-Date: Thu, 26 Dec 2019 19:22:54 +0000
Message-ID: <20191226191849.6BBF111BA67D@ary.qy>
Authentication-Results: dmarc=fail (p=none; dis=none)header.from=iecc.com
Source-IP: 50.31.151.76
Delivery-Result: delivered
Auth-Failure: dmarc
Reported-Domain: iecc.com
```

The message in a failure report might be a legitimate one that was unaligned when sent, or modified on the way to become unaligned. Or it might be a fraudulent one, either an attempted phish, or just a random spam message where the spamware happened to pick your domain for the fake return address. For this particular report, it's obviously a real message relayed through the NANOG mailing list.

The original failure report included the full address of the recipient, meaning that by looking at the failure reports, anyone who posts to NANOG can see who subscribes to LinkedIn. This kind of data leakage explains why most sites don't send failure reports at all, and most of the ones that do limit what they send, typically including only the headers of a failing message and redacting recipient address details.

Using DMARC Reporting to Prepare for Policy Publication

Before publishing a DMARC policy of **quarantine** or **reject**, domain operators should be confident that as close as possible to 100% of the mail they send is DMARC aligned. They might send unaligned mail if SPF records of a domain do not cover all of the IP addresses that send valid mail, causing SPF validation to fail. Some outgoing *Mail Transfer/Transfer Agents* (MTAs) might have DKIM configured incorrectly or not at all, so there's no aligned DKIM signature. Large organizations often can have MTAs sending mail that the network managers didn't know about; for example, if a department set up its own local server, or contracted with a third-party mail sender.

The data from DMARC reports tells the operator what IP addresses are sending unaligned mail, and generally makes it straightforward to figure out why it's unaligned. Mitigation might involve updating the domain SPF records to include missing MTAs, fixing the DKIM signing configuration in MTAs, or enforcing rules about unapproved mail servers or third-party mail senders. (Many third parties can do DKIM signing with a client's domain, but that requires either sharing the private signing keys(s) or delegating a DNS subtree that the third party can manage.)

After the operator has the mail sufficiently under control, it can gradually turn on sending policies. DMARC provides the quarantine policy as an intermediate step between no policy and reject so there is a chance for recipients to retrieve miscategorized mail. It can also use the percentage parameter in the policy record to apply policies gradually and limit the damage if mistakes occur.

DMARC vs. Discussion Lists

DMARC was originally intended for domains at organizations like banks that send primarily business-to-business and business-to-customer mail, and little or no person-to-person mail. When the organization considers when to publish a DMARC policy, and what policy to publish, it should remember that some fraction of its legitimate mail will arrive unaligned because of intermediate processing that DMARC cannot describe. Since the organization presumably knows what mail it sends, it can weigh the benefits of less phishing versus the cost of lost mail and make a decision that is reasonable for the organization.

In 2014, AOL and Yahoo, two large consumer mail systems, had separate security breaches in which intruders stole copies of millions of their users' address books. The stolen data was quickly sold to spammers, who used it to send spam to AOL and Yahoo users that appeared to be from the recipients' friends. This situation caused an expensive support problem at AOL and Yahoo as users complained about the spam and asked why their friends were spamming them. First AOL, and then Yahoo, "solved" the problem by quickly publishing DMARC **p=reject** policies that told every mail system that implements DMARC to reject any AOL or Yahoo mail that didn't come directly from AOL or Yahoo. This decision was very different from the ones made by organizations described previously. In this case the benefit of the policy was to mitigate the cost of an operational failure, with little if any benefit for most of their users, while creating major problems for their discussion list users.

A small but important part of the mail from users of any consumer mail system is unaligned yet legitimate mail that recipients want. That happens typically because the mail is routed indirectly from the sending user to the ultimate recipients. A particular point of contention is e-mail discussion lists where the normal actions of list managers make most of the mail unaligned. This situation can cause non-receipt of mail sent to subscribers on mail systems that enforce DMARC policies on incoming mail, and it can also cause removal of subscribers from lists because of bounces caused by DMARC failures. (Yahoo was aware of the mailing list issues but decided to publish **p=reject** anyway, according to someone who was there at the time.) Another source of unaligned mail is third-party mailing services. A small organization like an athletic club or scout troop often has an announcement list where the return address on the announcements is the personal address of the organization's secretary, who may use AOL or Yahoo.

A variety of proposed workarounds have been made for the problems that DMARC causes to mailing lists, none of which are very satisfactory. Initially, the easiest approach was to tell people sending mail from addresses with DMARC policies to subscribe from another address. That approach stopped being practical when AOL and Yahoo flipped the switch.

Since then, mailing list software has taken a range of approaches to ensure that the messages the list sends out are aligned. In a few cases, lists tried to turn off any features that would modify messages in ways that would invalidate DKIM signatures, hoping that DKIM signatures on incoming messages would remain valid when resent from the list. This idea didn't work very well, partly because remailed messages weren't SPF aligned (the list uses its own envelope address for bounce management), and users want the changes that lists make, such as adding subject line tags to identify the list.

Mailing lists have settled on two general anti-DMARC approaches^[13]. The most common is to put the list address into the **From:** header so the list can add a DKIM signature with its own domain and make the message DMARC aligned. For example, if the incoming message included:

```
From: Steve C <steve@aol.com>
To: nodule@lists.example.com
```

The list might rewrite it as:

```
From: Steve C via the nodule list <nodule@lists.example.com>
To: somelist@lists.example.com
Reply-To: Steve <steve@aol.com>
```

The rewritten **From:** header usually includes the author's address comment and the list name. The author's actual address is placed in a **Reply-To:** header, or occasionally a **Cc:** header. This approach allows DMARC alignment, since the list can add a lists.example.com DKIM signature, but makes mail from the list harder to handle. Mail user agents treat **Reply-To:** in different ways, leading to confusion about whether someone is replying to the author of a message, or to the list, or both. Adding to the confusion, some lists only rewrite the headers for messages in author domains that publish a DMARC policy, so messages from the same list have different headers.

Another approach is to rewrite the **From:** header to replace the problematic author address with one that is DMARC aligned but still represents the author. For example, my mailing lists would rewrite the headers in the previous example like this, changing the author address only by adding a local domain suffix:

```
From: Steve C <steve@aol.com.dmarc.fail>
To: nodule@lists.example.com
```

The domain **dmarc.fail** is a real domain I registered. (It was available.) I publish an MX record for ***.dmarc.fail**, to catch any mail sent to rewritten addresses. The rewritten message as sent by the mailing lists has a **dmarc.fail** DKIM signature, so it's properly DMARC aligned. When the list software rewrites an address, it creates a forwarding entry for the rewritten address that redirects back to the original address. The forwarding entries are deleted after a few days so that replies to the author sent shortly after the original message go back to the author, but the forwarding is limited, so it's not a useful vector to relay third-party spam.

This technique works fairly well. Since only the **From:** header is changed, there's no effect on **Reply-To:** or other mail behavior, and the author's identity is easy to recognize. Other systems have implemented the same idea in perhaps less passive-aggressive ways. The IETF's working mailing lists rewrite the address into the local part; for example:

From: Steve C <steve=40aol.com@lists.ietf.org>

The commercial LISTSERV mailing list service rewrites the address into an opaque local address and puts the real address in **Reply-To:**

From: Steve C <00000006b01fa96f-dmarc-request@lists.example.com>
Reply-To: Steve C <steve@aol.com>, Nodule list <nodule@lists.example.com>

The primary disadvantage of the address rewriting is that it requires access to the local mail system of the list to manage the set of temporary forwarding addresses, rather than doing it entirely inside the list software.

The other anti-DMARC approach that some lists take is message wrapping, enclosing the message as a *Multi-Purpose Internet Mail Extensions* (MIME) part within an outer message from the list. Most mailing lists have a MIME digest option, to send the day's messages as a set of MIME parts within a single daily message, so this process in effect turns each message into a one-message digest. The outer message typically has the list address in the **From:** header, while the inner message is unmodified.

Technically, this approach should work well, because it uses existing well-standardized features of RFC 5322 mail. Having done some experiments to see how workable it is, I found that in practice it works very badly because mail user agents treat MIME attached messages as an afterthought. While the inner message is typically displayed legibly, it is often not possible to reply to the inner message without clumsy extra steps, or in some cases at all, and multipart messages or those with attachments are handled inconsistently. The IETF experimented with several varieties of MIME wrapping before deciding that rewriting the **From:** header was the best of a bad lot.

While all of these approaches allow mailing lists to send DMARC-aligned mail, none of them are very satisfactory, and none let mailing lists work as well as they did before DMARC.

ARC

While the amount of mail that large providers get from mailing lists is small, on the order of 1% to 2% of the non-spam total, it is mail that the recipients care about deeply. After years of complaints, several large mail providers developed *Authenticated Received Chain* (ARC) to help them handle wanted but unaligned user mail.

An obvious way to handle unaligned mail from mailing lists would be to whitelist them. Large mail systems have a pretty good idea of where the lists are (the number of mailing list hosts worldwide is probably only about 10,000), so they could just accept the mail from the lists that they know their users want. The problem with this concept is that mailing lists don't do a very good job of spam filtering, and spam leaks through them all the time.

In particular, most lists check only that the address in the **From:** header is subscribed to the list before forwarding a message. If a subscriber's account is compromised and starts sending spam, any message sent to a list will generally get forwarded to the list. Even without an account being compromised, if a stolen address book happens to contain your address and the address of a list to which you subscribe, spamware can forge mail from you to the list, and again the list will forward it. I've seen this happen multiple times, and it is quite frustrating since the person whose address is being forged can't do anything about it.

The goal of ARC is to add a "chain of custody" to a message that shows what happened to it each time it was forwarded. This technique lets the ultimate recipient system retroactively make spam filtering decisions based on what happened to the message at the forwarding systems.

ARC builds on existing mail technology. It adapts the *Authentication-Results* (A-R) header^[14] that many mail systems apply to incoming messages that records the authentication status of the message at the time an MTA received the message. Here is a typical A-R header that my MTA applied to an incoming message from Apple's **me.com**:

```
Authentication-Results: iecc.com; spf=pass spf.mailfrom=xxx@me.com
spf.helo=mr85p00im-hyfv06011401.me.com smtp.remote-ip="17.58.23.191";
dkim=pass header.d=me.com header.s=1alhai header.a=rsa-sha256;
dmarc=pass header.from=me.com (p=quarantine, pct=100)
```

The first field is the name of the system that added the header, followed by groups of authentication results, in this case for SPF, DKIM, and DMARC. Each group includes the result and relevant items like the envelope **MAIL FROM** and sending IP for SPF. All of the fields are optional other than the system name, and they are added only for the kinds of authentication the system checked. ARC combines a modified A-R header and two DKIM-like signature headers into an *ARC seal*, which is intended to describe the passage of a message through a system such as a mailing list manager. A single message may have multiple ARC seals if it has passed through multiple forwarding systems.

Each seal is numbered, starting with 1 for the first seal applied. Each header in an ARC seal has an `i=` clause to indicate which seal it's part of.

The headers in an ARC seal look like this:

```
ARC-Message-Signature: i=1; a=rsa-sha256; d=microsoft.com; s=abcd; h=From:Date:...
ARC-Authentication-Results: i=1; mx.microsoft.com 1; spf=pass ...; dkim=pass ...
ARC-Seal: i=1; a=rsa-sha256; s=abcd; d=microsoft.com; cv=none; b=j7M/jt9eVP...
```

The *ARC-Message-Signature* (AMS) is almost identical to a DKIM signature, with the added `i=` field. It is intended to cover the usual headers and body of the message, at the time the message was sent from the signing system. If the system makes changes to the message, the AMS is applied after those changes. When a message is received, the most recent AMS signature will be valid unless an intermediate system has modified the message since the ARC seal was applied and not added a seal of its own.

The *ARC-Authentication-Results* (AAR) header reports the authentication status at the time the sealing system received the message; that is, before any modifications reflected in the AMS.

The *ARC-Seal* header is a DKIM-like signature that covers only the three headers in the ARC seal, to validate the seal itself. It also indicates whether the chain of ARC seals in the message was intact when the message was sealed, using the `cv=` (chain value) field. If this seal is the first one, the chain value is “none” for no previous seal. For any subsequent seal, the chain value is “pass” if the previous seal was valid (the DKIM-like signatures validated) and the previous seal had `cv=none` or `cv=pass`. Otherwise the chain value is “fail.”

If a mail system receives a message from a trustworthy source with a valid ARC chain, it can use the information in the ARC seals to make exceptions to its DMARC policy. As a simple example, assume a message that is not DMARC aligned arrives, but it has a valid chain of ARC seals. In one of the seals, an AAR header shows that the message was DMARC aligned (`dmARC=pass`) and the **header.from** domain was the same as the one currently in the message. That means the lack of alignment is due to changes made by the forwarding system. If the forwarding system is considered trustworthy, for example, a host that hosts discussion lists, the receiving system can decide to deliver the message. More complicated analysis is possible, but I expect this sort of analysis looking for typical mailing list operations is likely to be the most common. Since malicious systems can add fake ARC seals, this analysis makes sense only for mail from trust-worthy sources. Identifying sources trustworthy enough to apply ARC exceptions may be a problem for mail systems too small to develop reliable data on hosts that send mail to them. There are some efforts to provide shared lists of reputable mailing list hosts that will likely be good enough, since the number of active list hosts is not large and changes slowly.

At this point the implementation of ARC has started, but it is not yet common enough to let mailing lists stop doing anti-DMARC header munging. *Python* and *Perl* libraries for DKIM have both added ARC support^[15]. The Sympa 6.2 mailing list manager has ARC support, as does GNU Mailman 3.1, but not Mailman 2.x.

Large mail systems including Google's *Gmail* and Microsoft's *outlook.com* have some ARC support, and both *Gmail* and *outlook.com* put ARC seals on forwarded and mailing list mail, but neither is yet using it for mail filtering other than experimentally. Few mailing lists yet add ARC seals, partly because of the lack of ARC support in the list software they currently use, and partly because the list managers are unaware of ARC.

Conclusions

DMARC started as a relatively simple technique to deter phishing of high-profile commercial domains such as those of banks and payment providers. Consumer mail systems AOL and Yahoo then repurposed it to deal with spam forging their users' addresses. While this repurposing largely solved the spam forgery problem of mail systems, it caused severe collateral damage to e-mail discussion lists. While many lists have tried to work around the DMARC problems, all of the workarounds have drawbacks that make them ultimately unsatisfactory. To help undo the DMARC damage, a group of large mail providers invented ARC, which makes it somewhat possible to examine the history of a message and see how a message that is not DMARC aligned got that way.

The ongoing evolution of DMARC, mailing lists, and ARC is yet another round of security measures with unexpected consequences. With any luck, ARC will be the end of this sequence of effect, side-effect, and counter-effect, but we won't know until ARC is more widely deployed, hopefully in a few years.

References and Further Reading

- [1] Dave Crocker, "Challenges in Anti-Spam Efforts," *The Internet Protocol Journal*, Volume 8, No. 4, December 2005.
- [2] John Klensin, "Another Look at Spam," *The Internet Protocol Journal*, Volume 8, No. 4, December 2005.
- [3] Murray Kucherawy and Elizabeth Zwicky, Eds., "Domain-based Message Authentication, Reporting, and Conformance (DMARC)," RFC 7489, March 2015.
- [4] Peter W. Resnick, "Internet Message Format," RFC 5322, October 2008.
- [5] Scott Kitterman, "Sender Policy Framework (SPF) for Authorizing Use of Domains in Email, Version 1," RFC 7208, April 2014.

- [6] Murray Kucherawy, David Crocker, and Tony Hansen, “DomainKeys Identified Mail (DKIM) Signatures,” RFC 6376, September 2011.
- [7] John C. Klensin, “Simple Mail Transfer Protocol,” RFC 5321, October 2008.
- [8] John Levine, Mark Delany, Eric Allman, and Jim Fenton, “DomainKeys Identified Mail (DKIM) Author Domain Signing Practices (ADSP),” RFC 5617, August 2009.
- [9] See <https://publicsuffix.org/> for the PSL, and https://wiki.mozilla.org/Public_Suffix_List for a description of its use and history.
- [10] See <https://www.taugh.com/rddmarc/>
- [11] Dmarcian: www.dmarcian.com
- [12] Agari: www.agari.com
- [13] Mailman and DMARC, <https://wiki.list.org/DEV/DMARC>
- [14] Murray Kucherawy, “Message Header Field for Indicating Message Authentication Status,” RFC 8601, May 2019.
- [15] See <https://pypi.org/project/dkimpy/> for the *Python* library, and <https://metacpan.org/release/Mail-DKIM> for the *Perl* library.

JOHN R. LEVINE writes, speaks, and consults on the Internet, electronic mail, cybersecurity, and related topics. He speaks to many trade, policy, and general groups. He has testified at the *Federal Trade Commission Spam Forum* on the mechanics of spam, to the *Senate Commerce Committee on Spyware*, and is part of the *Industry Canada Task Force on Spam*. He has spoken at the Internet Law and Policy Forum and at many conferences. He is frequently interviewed in the print and electronic media and has extensive working relationships with reporters. John consults and provides advice and expertise on e-mail and Internet systems, security, and software. He co-founded the *Domain Assurance Council*, a non-profit industry consortium that establishes standards for e-mail certification and security. Levine has served as an expert witness on a variety of computer topics including e-mail spam, compiler software, and graphic image file formats. He has written many books on the Internet and other computer topics. His books range from the best-selling *Internet for Dummies*, with over seven million copies of eleven editions in print in dozens of languages, *Fighting Spam for Dummies*, and *Windows Vista: The Complete Reference*, to books on computer language tools and graphics programming. E-mail: john1@taugh.com

Letter to the Editor

Ole,

I enjoyed reading Geoff Huston's article "MSS Values of TCP," (*The Internet Protocol Journal*, Volume 22, No. 3, December 2019). I had not been familiar with the variation of *Maximum Segment Size* (MSS) values that are used in the broad Internet.

I have never run into a client or server device that has been unable to operate with greater than a 576-byte *Maximum Transmission Unit* (MTU). Even the early Intel 8088-based MS DOS PCs in the 1980s with 3Com 3c501 *Network Interface Cards* (NICs) could handle 1500-byte MTU. In modern times only a tiny fraction of our tens of thousands of connected hosts are capable of Ethernet Jumbo packets (SAN nodes replicating data between data centers; our routers along just those paths are configured for 9,216-byte MTU).

Geoff touched only briefly on encapsulation influencing resulting TCP MSS values. Our *Wide-Area Network* (WAN) connections between our office locations use *Internet Protocol Security* (IPsec) tunnels, and we also use IPsec tunnels between our office locations and our virtual routers inside the Amazon *Virtual Private Cloud* (VPC). In addition to those, we have some *Generic Routing Encapsulation* (GRE) tunnel connections with some information providers and partners (Zscaler is an example). With the recent availability of low-cost high-speed Layer 2 connections between some sites, we have been implementing *Media Access Control security* [MACsec] (lower router resource consumption than IPsec). We also are beginning to use *Virtual Extensible LAN* (VXLAN) between our two data centers. With all of these encapsulations in our network, we've been avoiding IP fragmentation of TCP packets by configuring TCP `adjust-mss` on our Cisco routers. We also use `adjust-mss` on our Wireless LAN Controllers.

We still face plenty of *User Datagram Protocol* (UDP) packets from video devices that would need to be fragmented if attached with default configuration. When our IT group controls such devices, we configure their MTU to fit within our IPsec tunnels.

Regards, —Richard Berke, Richard.Berke@troweprice.com

The author responds: Richard,

Thanks for your comments about MSS sizes and the related topic of MTU selection.

We need to go back to the 1970s to find some variation from the current ubiquity of the 1,500-octet MTU that dominates today's communications, and very little of that early network environment has survived. However, we can piece together some of the thinking behind the original design of the Internet Protocol and the selection of MTU and MSS values.

Smaller packet sizes made packets less susceptible to bit-error-rate corruption and could reduce jitter (which was a major consideration behind the design of *Asynchronous Transfer Mode* (ATM) cells), but at the same time smaller packets had a reduced payload efficiency. Various mainframe vendors tuned their network products to match their intended deployment environment, including the choice of supported packet sizes.

As a “network of networks,” the Internet was envisaged to work across various permutations of networks, all with differing MTU sizes. The fragmentation model of IP Version 4 came from this approach, where an IP router was permitted to fragment a packet if it was too large for the next network.

IP Version 4 permitted IP packets between 20 and 65,535 octets in size. While in a strict sense the minimum MTU is 20 octets, that is without any payload at all. A 21-octet MTU would make some level of progress in sending a payload, albeit extremely inefficiently.

Where does 576 octets come from? IP hosts were not required to accept the entire protocol-permitted range of packet sizes. The specification required IP hosts to reassemble and accept IP packets up to 576 octets in size. Why 576? It is such an odd number. I could understand a value of 532, 542, 572, or even 592 octets, all based on a 512-octet payload and various permutations of minimum or maximum IP headers and optionally including a TCP header. However, I can’t get to 576 octets that way, so I don’t have any credible explanation as to why this value was chosen.

By the time we were designing IPv6 in the early 1990s the thinking had changed, and fragmentation was frowned upon. It was slow and insecure, and the experts advised its avoidance wherever possible. What should the minimum unfragmented MTU be for IPv6? Ethernet framing was ubiquitous by this time, so a 1,500-octet MTU size seemed like a good first answer. But the Internet had a new aspect by then: it was no longer a “network of networks” but was the base substrate network upon which other networks were overlaid. Various other headers, including IP-in-IP, were being added. So, we needed to specify a universal minimum unfragmented IPv6 packet size that would be relevant in many kinds of IP-in-IP contexts. The value of 1,280 octets as the new minimum unfragmentable packet size was chosen for IPv6. Why 1,280? I understand that this number was chosen because it’s the sum of 1,024 and 256.

My view is that the marginal loss of payload efficiency is small enough that for the public Internet a 1,220-octet MTU can be used with some confidence that it will not encounter MTU mismatch issues.

Regards,

—Geoff Huston, gih@apnic.net

Keeping the DNS Secure During the Coronavirus Pandemic

The Internet's value in bringing people together has never been more apparent than it is now. While most of us are under some form of "stay at home" order in an effort to slow the spread of the coronavirus, the Internet provides us with a lifeline. It brings us information and entertainment, allows some of us to continue our work and education, and brings us what we need most at times of isolation—social connections.

The role of the *Internet Corporation for Assigned Names and Numbers* (ICANN) community, Board, and organization in maintaining a secure, stable, and unified Internet has always been important, but at this time, when reliance on the Internet has skyrocketed, our collective role has become all the more vital. ICANN's mission frames our concern about cybercriminals who are exploiting the pandemic by perpetrating scams and victimizing Internet users. Some are selling phony cures, treatments, and vaccines. Some are using domain names as part of their efforts to prey on people at this time when many are experiencing anxiety, fear, and loneliness. The U.S. *Federal Trade Commission* reports that it has fielded more than 7,800 coronavirus-related complaints. The agency noted that U.S. consumers alone have collectively lost more than U.S. \$5 million.

Of course, ICANN cannot involve itself in content issues, both because of our Bylaws as well as practically, but that does not mean we are unaware or unconcerned about those who are using the *Domain Name System* (DNS) to victimize others. It is this concern that prompted me to contact the registries and registrars thanking them for their efforts and actions aimed at helping to mitigate and minimize the abusive domain names being used to maliciously take advantage of the coronavirus pandemic. For example, the *Registrar Stakeholder Group*^[1] has posted a useful guide, entitled "Registrar approaches to the COVID-19 Crisis" that provides a number of steps and resources the registrar community can use in their efforts.

Many of our contracted parties already support a *Framework to Address Abuse*,^[2] which deals with DNS abuse and website content abuse. I continue to commend them for making this commitment to protect the DNS from those who would maliciously exploit domain names. In my correspondence to the registries and registrars, I expressed ICANN org's appreciation for their work during the pandemic.

Additionally, I'm pleased to tell you that ICANN org has joined registries, registrars, security experts, law enforcement, Internet engineers, and others, in the COVID-19 *Cyber Threat Coalition* (CTC)^[3]. The CTC's mission is to, "operate the largest professional-quality threat lab in the history of cybersecurity out of donated cloud infrastructure and with rapidly assembled teams of diverse, cross-geography, cross-industry threat researchers."

I am proud that so many in the Internet ecosystem are joining together during this crisis to stop those who prey on the desperate. We will continue to keep you advised of our engagement efforts to mitigate the misuse of domain names during these critical times.

—Göran Marby, *President and Chief Executive Officer, ICANN*

[1] <https://rrsg.org/>

[2] <http://dnsabuseframework.org>

[3] <https://www.cyberthreatcoalition.org/>

Global Encryption Coalition Formed

Encryption is a critical technology that helps keep people, their information, and communications private and secure. However, some governments and organisations are pushing to weaken encryption, which would create a dangerous precedent that compromises the security of billions of people around the world. Actions in one country that undermine encryption threaten us all.

As a global coalition, we call on governments and the private sector to reject efforts to undermine encryption and pursue policies that enhance, strengthen and promote use of strong encryption to protect people everywhere. We also support and encourage the efforts of companies to protect their customers by deploying strong encryption on their services and on their platforms.

The mission of the *Global Encryption Coalition* is to promote and defend encryption in key countries and multilateral gatherings where it is under threat. It also supports efforts by companies to offer encrypted services to their users.

With a steering committee led by the *Center for Democracy and Technology* (CDT), *Global Partners Digital* (GPD) and the *Internet Society* (ISOC), the Global Encryption Coalition is composed of national coalitions, civil society groups, corporations, academics, and technologists around the world who agree to support its founding statement.

For more information, visit: <https://www.globalencryption.org/>

APNIC Launches Networking from Home Events

With most Asia Pacific economies forced into various states of lockdown to minimize COVID-19 infections, *Network Operator Group* (NOG) meetings and other technical events in the region have either been cancelled or postponed. NOGs are a great forum for network engineers to share experience with their peers, work out solutions to common technical problems, and build the strong relationships that help the Internet operate. There are 22 NOGs in the APNIC region, but sadly that means a lot of events have been cancelled in 2020.

Networking from Home is a new virtual event initiative to provide a place for technical folk in the region to share their experience and expertise with their peers, just like they would at a NOG event.

There will be four free online events—one each held in the time zones of South East Asia, South Asia, East Asia, and Oceania—and they will be a digestible 2.5 hours long. Presentations will be short and punchy and interaction is encouraged! APNIC’s Geoff Huston will deliver a different keynote at each event, and he will be supported by a range of great speakers suggested by the NOG communities.

If you have a great presentation in mind, get in touch with the Program Committees for the events. For more information, visit: <https://nfh.apnic.net/>

Check your Subscription Details!

If you have a print subscription to this journal, you will find an expiration date printed on the back cover. For the last couple of years, we have “auto-renewed” your subscription, but now we ask you to log in to our subscription system and perform this simple task yourself. The subscription portal is here: <https://www.ipjsubscription.org/> This process will ensure that we have your current contact information, as well as delivery preference (print edition or download). For any questions, contact us by e-mail at: ipj@protocoljournal.org

Our Privacy Policy

The *General Data Protection Regulation* (GDPR) is a regulation for data protection and privacy for all individual citizens of the *European Union* (EU) and the *European Economic Area* (EEA). Its implementation in May 2018 led many organizations worldwide to post or update privacy statements regarding how they handle information collected in the course of business. Such statements tend to be long and include carefully crafted legal language. We realize that we may need to provide similar language on our website and in the printed edition, but until such a statement has been developed here is an explanation of how we use any information you have supplied relating to your subscription:

- The mailing list for *The Internet Protocol Journal* (IPJ) is entirely “opt in.” We never have and never will use mailing lists from other organizations for any purpose.
- You may unsubscribe at any time using our online subscription system or by contacting us via e-mail. We will honor any request to remove your name and contact information from our database.
- We will use your contact information only to communicate with you about your subscription; for example, to inform you that a new issue is available, that your subscription needs to be renewed, or that your printed copy has been returned to us as undeliverable by the postal authorities.
- We will never use your contact information for any other purpose or provide the subscription list to any third party other than for the purpose of distributing IPJ by post or by electronic means.
- If you make a donation in support of the journal, your name will be listed on our website and in print unless you tell us otherwise.

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Fabrizio Accatino	Tracy Camp	Mikhail Evstiounin	Johan Helsingius	Andrew Koch
Michael Achola	Ignacio Soto Campos	Paul Ferguson	Robert Hinden	Ia Kochiashvili
Martin Adkins	Fabio Caneparo	Ricardo Ferreira	Asbjorn Hojmark	Carsten Koempe
Christopher Affleck	Roberto Canonico	Kent Fichtner	Damien Holloway	Richard Koene
Scott Aitken	David Cardwell	Michael Fiumano	Alain Van Hoof	Alexader Kogan
Jacobus Akkerhuis	John Cavanaugh	The Flirble Organisation	Edward Hotard	Antonin Kral
Antonio Cuñat Alario	Lj Cemerar	Gary Ford	Bill Huber	Robert Krejčí
Nicola Altan	Dave Chapman	Jean-Pierre Forcioli	Hagen Hultzs	Mathias Körber
Matteo D'Ambrosio	Stefanos Charchalak	Susan Forney	Kevin Iddles	John Kristoff
Jens Andersson	Greg Chisholm	Christopher Forsyth	Mika Ilvesmaki	Terje Krogdahl
Danish Ansari	David Chosrova	Andrew Fox	Karsten Iwen	Bobby Krupczak
Finn Arildsen	Marcin Cieslak	Craig Fox	David Jaffe	Murray Kucherawy
Tim Armstrong	Brad Clark	Fausto Franceschini	Ashford Jaggernaut	Warren Kumari
Richard Artes	Narelle Clark	Valerie Fronczak	Martijn Jansen	George Kuo
Michael Aschwanden	Joseph Connolly	Tomislav Futiv	Jozef Janitor	Dirk Kurfuerst
David Atkins	Steve Corbató	Edward Gallagher	John Jarvis	Darrell Lack
Jac Backus	Brian Courtney	Andrew Gallo	Dennis Jennings	Yan Landriault
Jaime Badua	Dave Crocker	Chris Gamboni	Edward Jennings	Markus Langenmair
Bent Bagger	Kevin Croes	Xosé Bravo Garcia	Aart Jochem	Fred Langham
Eric Baker	John Curran	Osvaldo Gazzaniga	Brian Johnson	Andrew Lamb
Santosh Balagopalan	André Danthine	Kevin Gee	Curtis Johnson	Richard Lamb
Michael Bazarewsky	Morgan Davis	Greg Giessow	Richard Johnson	Sig Lange
David Belson	Jeff Day	John Gilbert	Jim Johnston	Tracy LaQuey Parker
Hidde Beumer	Julien Dhallenne	Serge Van Ginderachter	Jonatan Jonasson	Rick van Leeuwen
Pier Paolo Biagi	Freek Dijkstra	Greg Goddard	Daniel Jones	Simon Leinen
John Bigrow	Geert Van Dijk	Tiago Goncalves	Gary Jones	Robert Lewis
Orvar Ari Bjarnason	David Dillow	Octavio Alfageme	Jerry Jones	Christian Liberale
Axel Boeger	Richard Dodsworth	Gorostiaga	Anders Marius	Martin Lillepuu
Keith Bogart	Ernesto Doelling	Barry Greene	Jørgensen	Roger Lindholm
Mirko Bonadei	Michael Dolan	Jeffrey Greene	Amar Joshi	Sergio Loreti
Roberto Bonalumi	Eugene Doroniuk	Richard Gregor	David Jump	Eric Louie
Julie Bottorff	Karlheinz Dölger	Martijn Groenleer	Merike Kao	Guillermo a Loyola
Photography	Joshua Dreier	Geert Jan de Groot	Andrew Kaiser	Hannes Lubich
Gerry Boudreaux	Lutz Drink	Christopher Guemez	Christos Karayiannis	Dan Lynch
L de Braal	Andrew Dul	Gulf Coast Shots	David Kekar	Sanya Madan
Kevin Breit	Joan Marc Riera	Sheryll de Guzman	Jithin Kesavan	Miroslav Madić
Thomas Bridge	Duocastella	Rex Hale	Jubal Kessler	Alexis Madriz
Ilia Bromberg	Holger Durer	Jason Hall	Shan Ali Khan	Carl Malamud
Václav Brožík	Mark Eanes	James Hamilton	Nabeel Khatri	Jonathan Maldonado
Christophe Brun	Peter Robert Egli	Stephen Hanna	Dae Young Kim	Michael Malik
Gareth Bryan	George Ehlers	Martin Hannigan	William W. H.	Yogesh Mangar
Stefan Buckmann	Peter Eisses	John Hardin	Kimandu	Bill Manning
Caner Budakoglu	Torbjörn Eklöv	David Harper	John King	Harold March
Darrell Budic	Y Ertur	Edward Hauser	Russell Kirk	Vincent Marchand
Scott Burleigh	ERNW GmbH	David Hauweele	Gary Klesk	Gabriel Marroquin
Jon Harald Bøvre	ESdatCo	Marilyn Hay	Anthony Klopp	David Martin
Olivier Cahagne	Steve Esquivel	Headcrafts SRLS	Henry Kluge	Jim Martin
Antoine Camerlo	Jay Etchings	Hidde van der Heide	Michael Kluk	Ruben Tripana Martin

Timothy Martin	John O'Neill	David Ross	Job Snijders	Tim Weil
Carles Mateu	Jim Oplotnik	William Ross	Ronald Solano	Jd Wegner
Juan Jose Marin	Packet Consulting	Boudhayan	Asit Som	Westmoreland
Martinez	Limited	Roychowdhury	Ignacio Soto	Engineering Inc.
Ioan Maxim	Carlos Astor Araujo	Carlos Rubio	Campos	Rick Wesson
David Mazel	Palmeira	Timo Ruiter	Evandro Sousa	Peter Whimp
Miles McCredie	Alexis Panagopoulos	RustedMusic	Peter Spekreijse	Russ White
Brian McCullough	Gaurav Panwar	Babak Saberi	Thayumanavan	Jurrien Wijlhuizen
Joe McEachern	Manuel Uruena Pascual	George Sadowsky	Sridhar	Derick Winkworth
Alexander McKenzie	Ricardo Patara	Scott Sandefur	Paul Stancik	Pindar Wong
Jay McMaster	Dipesh Patel	Sachin Sapkal	Ralf Stempffer	Phillip Yaleloglou
Mark Mc Nicholas	Alex Parkinson	Arturas Satkovskis	Matthew Stenberg	Janko Zavernik
Carsten Melberg	Craig Partridge	PS Saunders	Adrian Stevens	Muhammad Ziad
Kevin Menezes	Dan Paynter	Richard Savoy	Clinton Stevens	Ziayuddin
Bart Jan Menkveld	Leif Eric Pedersen	John Sayer	John Streck	Jose Zumalave
William Mills	Rui Sao Pedro	Phil Scarr	Martin Streule	Romeo Zwart
David Millsom	Juan Pena	Elizabeth Scheid	Viktor Sudakov	Bernd Zeimetz
Desiree Miloshevic	Chris Perkins	Jeroen Van Ingen	Edward-W. Suor	廖 明沂.
Joost van der Minnen	Michael Petry	Schenau	Vincent Surillo	
Thomas Mino	Alexander Peuchert	Carsten Scherb	T2Group	
Rob Minshall	David Phelan	Ernest Schirmer	Roman Tarasov	
Wijnand Modderman	Derrell Piper	Philip Schneck	David Theese	
Mohammad Moghaddas	Rob Pirnie	Dan Schrenk	Douglas Thompson	
Charles Monson	Marc Vives Piza	Richard Schultz	Lorin J Thompson	
Andrea Montefusco	Jorge Ivan Pincay Ponce	Timothy Schwab	Joseph Toste	
Fernando Montenegro	Victoria Poncini	Roger Schwartz	Rey Tucker	
Joel Moore	Blahoslav Popela	SeenThere	Sandro Tumini	
John More	Eduard Llull Pou	Scott Seifel	Angelo Turetta	
Maurizio Moroni	Tim Pozar	Yury Shefer	Phil Tweedie	
Brian Mort	David Raistrick	Yaron Sheffer	Steve Ulrich	
Soenke Mumm	Priyan R Rajeevan	Doron Shikmoni	Unitek Engineering AG	
Tariq Mustafa	Balaji Rajendran	Tj Shumway	John Urbanek	
Stuart Nadin	Paul Rathbone	Jeffrey Sicuranza	Martin Urwaleck	
Michel Nakhla	William Rawlings	Thorsten Sideboard	Betsy Vanderpool	
Mazdak Rajabi Nasab	Bill Reid	Greipur Sigurdsson	Surendran	
Krishna Natarajan	Petr Rejhon	Andrew Simmons	Vangadasalam	
Naveen Nathan	Robert Remenyi	Pradeep Singh	Ramnath Vasudha	
Darryl Newman	Rodrigo Ribeiro	Henry Sinnreich	Philip Venables	
Thomas Nikolajsen	Glenn Ricart	Geoff Sisson	Buddy Venne	
Paul Nikolich	Justin Richards	Helge Skrivervik	Alejandro Vennera	
Travis Northrup	Mark Risinger	Darren Sleeth	Luca Ventura	
Marijana Novakovic	Gregory Robinson	Richard Smit	Tom Vest	
David Oates	Ron Rockrohr	Bob Smith	Dario Vitali	
Ovidiu Obersterescu	Carlos Rodrigues	Courtney Smith	Michael L Wahrman	
Tim O'Brien	Magnus Romedahl	Eric Smith	Laurence Walker	
Mike O'Connor	Lex Van Roon	Mark Smith	Randy Watts	
Mike O'Dell	Alessandra Rosi	Craig Snell	Andrew Webster	



Follow us on Twitter and Facebook

@protocoljournal



<https://www.facebook.com/newipj>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Diamond Sponsors



Ruby Sponsors



Sapphire Sponsors

Your logo here!

Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

David Conrad, Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

September 2020

Volume 23, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

FROM THE EDITOR

In This Issue

From the Editor	1
Link-State Protocols on Data-Center Fabrics.....	2
The IPv4 Marketplace	20
In Memoriam	34
Book Review.....	36
Fragments	38
Thank You!	40
Call for Papers.....	42
Supporters and Sponsors	43

In June 2013 we published an article entitled “Optimizing Link-State Protocols for Data Center Networks.” In that article, Alvaro Retana and Russ White wrote: “With the advent of cloud computing, the pendulum has swung from focusing on wide-area or global network design toward a focus on *Data Center* network design. Many of the lessons we have learned in the global design space will be relearned in the data center space before the pendulum returns and wide-area design comes back to the fore.” In this issue, Russ White and Melchior Aelmans examine the use of link-state alternatives to the *Border Gateway Protocol* (BGP) in data center designs. One such alternative, *Routing in Fat Trees* (RIFT), will be explored further in an upcoming article in this journal, so please make sure your subscription details are up-to-date.

The depletion of the IPv4 address space and transition to IPv6 has been covered in numerous articles in IPJ over more than two decades. It was initially believed that “everyone” would implement IPv6 by the time the *Regional Internet Registries* (RIRs) ran out of addresses, but such predictions have proven to be too optimistic for a variety of reasons. The demand for public IPv4 address space has led to an “aftermarket,” whereby blocks of addresses can be purchased (or leased) through the use of address brokers. We asked David Strom to explore this market in more detail, and he approached this assignment by deciding to sell his own Class C address block.

We are excited to bring you another book review, this time on the topic of information security. Please send us your suggestions for networking-related books that we should have reviewed.

Publication of *The Internet Protocol Journal* is made possible by the generous support of numerous individuals and organizations. Please consider making a donation or getting your company to sign up for a sponsorship.

As always, we welcome your feedback and suggestions on anything you read in this journal. Letters to the Editor may be edited for clarity and length and can be sent to ipj@protocoljournal.org

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

Recent Developments in Link State on Data-Center Fabrics

by Russ White and Melchior Aelmans, Juniper Networks

Since the initial publication of the drafts resulting in RFC 7938^[0], the *Border Gateway Protocol* (BGP)^[10] has been the default choice for *Data-Center* (DC) fabrics, assumed by most controllers, intent-based systems, training courses in DC fabrics, and implementers. Recent activity in the *Internet Engineering Task Force* (IETF) and implementers suggests using link-state protocols in DC fabrics. This article explores why this move towards link-state protocols on DC fabrics is taking place, and then considers three specific avenues to link state on DC fabrics: *Distributed* (or localized) *Optimized Flooding in Intermediate System-to-Intermediate System Protocol* (IS-IS)^[1], centralized calculation of optimal flooding trees^[2, 3], and *Routing in Fat Trees* (RIFT)^[4].

The arguments presented in this article are legitimate reasons not to use BGP for the DC fabric underlay and show that options other than BGP are available. Readers might (incorrectly) conclude the authors believe BGP should never be used as the routing protocol for a DC fabric overlay—but that is not true. To make a case for link-state protocols in DC fabric underlays, an extensive examination of the positive and negative aspects of BGP—and the other available protocols—is essential. Ultimately, it is up to individual operators to decide which protocol is “the best” for their application, a decision based on business and operational—as well as technical—reasons.

Defining Terms

Defining terms is often considered pedantic, and therefore often overlooked. But as the networking world spreads to wholly virtual environments, definitions quickly become blurry and local. In this article, two distinct terms that might be used differently in other contexts are used. The first is the *underlay*. For this article, the underlay is the physical infrastructure, control plane, and telemetry; it provides basic connectivity, including IPv4, IPv6, and *Multiprotocol Label Switching* (MPLS), edge-to-edge in the fabric. The *overlay*, on the other hand, provides virtual topologies which tunnel traffic edge-to-edge through fabric-side interfaces and devices. In other words, the overlay consists of tunnels with head- and tail-ends on *Top of Rack* (ToR, or *leaf*) switches or servers attached to the fabric and the control planes that provide reachability through those tunnels.

In other words, underlay control planes do not carry overlay reachability, overlay control planes do not carry underlay reachability, and underlay devices (other than where they terminate an overlay tunnel) do not switch based on overlay destinations. If this explanation sounds vaguely like the *Exterior Gateway Protocol* (EGP) versus *Interior Gateway Protocol* (IGP) split in a traditional transit provider network, that is because it is—just like the EGP/IGP split in a conventional transit network.

The underlay/overlay distinction might be confusing in some discussions because people who work entirely within cloud services may well consider the set of tunnels built between virtual machines or containers the overlay, and everything under these tunnels the underlay—even if the network has two layers of tunnels. There is no set of standard terms for the situation where a bottom layer provides connectivity based on the physical fabric topology, a collection of virtual networks on top of that used to create logical topologies, and another set of virtual networks within those logical topologies formed by the applications running over the network. Overlay tends to end up being used for both the “middle layer” and the “upper layer,” hence the importance of defining the terms as they are used here.

Two other terms of importance here are *autonomic* and *automatable*. Confusion around these terms arises from the use of *Zero Touch Provisioning* (ZTP), used to mean the configuration of a device that does not need manual configuration to deploy. While both automated and autonomic networks are ZTP, there is still a difference between the two concepts. The closer a protocol comes to not needing any configuration at all, whether that configuration is automated or not, the closer the protocol is to being autonomic. While autonomic and automatable protocols appear similar, there are differences. The automation system must still be managed and maintained, there are still interfaces to integrate and manage, etc. Autonomic control planes may (or may not) be more complex, at least under the surface, than automatable ones; regardless, they are not the same thing.

It is rare for a protocol to be fully automated or fully autonomic; these two are a continuum rather than a binary space. For instance, BGP is not autonomic by design but can be modified to allow BGP speakers to discover one another and form a peering relationship automatically. In other cases, it might be possible to derive information, such as the IS-IS system ID, automatically, but doing so might make troubleshooting and maintenance more difficult—so automatic assignment might be possible, but not always desirable. Individual operators may have different optimal positions along the automatable-to-autonomic continuum.

Given autonomic to automatable is a spectrum of options, why would you choose to move towards autonomic operation? After an automation system is put in place to support the network, it may not seem to make much difference.

In a sense, moving from automatable to autonomic is simply shifting complexity from one place in the network to another. Moving configuration from the automation system to the protocol moves complexity from the automation system to the protocol as well. The one vital difference is each piece of complexity moved from the automation system to the protocol is one less interface to manage, one less piece of state the automation system must build and keep track of, etc.

Complex automation systems can be difficult to create and manage. Even in fully automated networks, research shows a major portion of network failures are caused by configuration failures.^[5] Deducing the amount of state the automation system, and the humans supporting the automation system, is managing can be justified by reducing the number of places mistakes can happen. The more the protocol can figure out on its own, the less you must figure out how to configure.

This article assumes *spine-and-leaf* (or leaf-and-spine!) fabrics. A *Clos*^[14] is a three-stage fabric, while a five-stage fabric wired in the most common way is called a *butterfly*. Butterfly fabrics are illustrated in two ways; “as-wired” with the leaves on both the top and bottom of the diagram and “folded” with the leaves arrayed at the bottom of the diagram.^[6]

The number of stages in a spine-and-leaf fabric denotes the maximum number of switches a packet is forwarded through when crossing from edge to edge. A spine-and-leaf fabric may have multiple spines; a three-stage fabric has one spine, while a five-stage fabric has three spine stages. The inner or top-most tier (depending on how the fabric is drawn) is considered the *Top of Fabric* (ToF) or the fabric layer. In a five-stage fabric, the “middle tier” is called either the *Top of Pod* (ToP) or the *spine*. The level or tier denotes the distance from the edge, with the ToR or leaf nodes being considered T0, the “middle” stage T1, and the fabric or ToF stage T2.

BGP in the Underlay

For the last 10 to 15 years, BGP has been the “underlay protocol of choice” for DC fabrics. While the reasons for using BGP in the underlay have been outlined in several places through the years, including RFC 7938^[9], it is useful to recap and explore some of these reasons as background.

First, BGP is widely implemented; virtually every routing vendor and every open-source routing stack such as *Free Range Routing* (FRRouting) has a fairly complete and well-tested BGP implementation. You can be confident that no matter whose hardware and software you choose, BGP will be supported—and the implementation is likely to be mature, interoperable with other implementations, and running in production in a lot of networks.

Second, BGP was—at least at one time—conceived of as one of the most straightforward routing protocols to understand and implement well. The logic of path-vector is reasonably easy to implement correctly, and the underlying transport mechanism, the *Transmission Control Protocol* (TCP), is built into every operating system already.

Third, BGP is widely deployed, and hence well understood by operators. Operators consider it easier to hire someone who knows BGP than one who knows any other protocol, and it is easy to find tooling for operating BGP in the open-source community.

There is a bit of irony in this point as 10 years ago it was almost impossible to find engineers with solid BGP experience; the advent of BGP on large-scale data-center fabrics has become something of a “self-fulfilling prophecy” in this regard.

Fourth, where scale is of concern, the perception is BGP outshines every other protocol. After all, “BGP runs the global Internet,” and you cannot ask for a better proof point of scalability than that. The initial implementations of BGP on large-scale DC fabrics originally tried various IGP, and found they could not scale to the size required.

Fifth, BGP has extensive prefix-filtering, route-tagging, and traffic-engineering capabilities. No other protocol, other than perhaps *Enhanced Interior Gateway Routing Protocol* (EIGRP) (!), can match the ability of BGP to control route flow.

Sixth, you can use BGP for both the underlay and the overlay in a single network. In theory, this possibility makes the configuration simpler. The normal configuration when using BGP for both is to configure the underlay using *External BGP* (eBGP) peering and the overlay as *Interior BGP* (iBGP) peering.

With all these advantages, why would you decide to move away from using BGP in both the underlay and overlay?

Challenging BGP in the Underlay

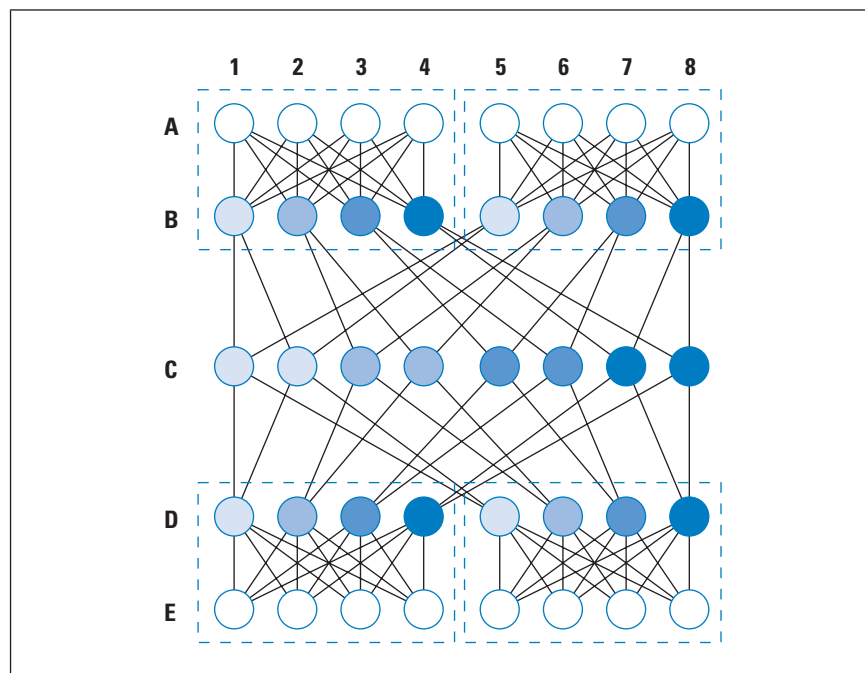
There are, however, counterpoints to many of the advantages of using BGP as the underlay protocol listed previously. Beginning with the second one—BGP is one of the simplest routing stacks to implement. With the advent of multiple address families, the *Resource Public Key Infrastructure* (RPKI), *Ethernet VPN* (EVPN), *Virtual Private LAN Service* (VPLS), MPLS traffic engineering, *BGP Link-State* (BGP-LS), and the many other features that have been “piled into” BGP across the last 20 years, BGP implementations have exploded in complexity. BGP may be the most complex protocol to implement among all the routed control planes today.

Using BGP as a singular DC fabric protocol, both overlay and underlay, is one factor causing the increasing complexity of BGP implementations. The ability to peer on unnumbered interfaces, the ability to accept any peer with any *Autonomous System* (AS) number, the ability to accept routes without any filters implemented, and many other changes must be made to make BGP work correctly in a DC fabric. It is easy enough to create a single knob that turns on a group of features at once. It is not so easy to hide the increased complexity—and the higher chance of a defect in the code or a misconfiguration of some kind—resulting from these kinds of changes. BGP is strongly automatable, but it will take massive work to make it autonomic. Is pushing that work into code used at critical points throughout the Internet a good idea?

At some point, the routing community needs to decide if it is wise to make one protocol the “protocol to end all protocols.” Is a single solution the right answer for all problems? Or is it better to move back towards developing multiple parallel protocols to support different purposes? This criticism may not apply to operators building their private implementation of BGP for use on their DC fabrics—but these kinds of implementations are uncommon.

A second related issue is the amount of specialized configuration required to allow BGP to converge quickly on the kinds of dense topologies used for DC fabrics. Figure 1 illustrates the design.

Figure 1: A Small Butterfly Fabric



Note that in this diagram, A and E are ToR switches or leaf nodes, B and D are spines, and C is either the superspine or fabric. Dashed boxes around a set of devices indicate the pods. How BGP converges depends on the kind of topology change. In the case of a single router or link failure, BGP can converge almost as quickly as an IGP, given the failure timers are tuned correctly, BFD and other underlying mechanisms are in use, etc. The case of a withdrawal from the edge of the network, however, is much different.

In the case of a withdrawal, BGP converges by hunting across available paths, starting from the shortest and ending in the longest. This hunt does not happen because of the way BGP is designed, but rather because of the timing of processing and forwarding updates. To prevent loops, a BGP speaker must process an update locally, modifying the routing table before it can forward the update to its peers. Longer paths just take longer than shorter ones for withdrawals to traverse. This withdrawal behavior can be a problem in at least two situations: when a workload is moved from one location on the fabric to another, and when an anycast address representing a service instance is removed from the fabric.

In these cases, the slow convergence time of BGP can impact applications running on the fabric.

Controlling the impact of the hunt is fairly easy. The key is to reduce the length of the paths through which BGP must “search.” The easiest way to do it is to block the reflection of updates and withdrawals through the network. For instance, E1 in Figure 1 should not reflect any withdrawals or updates to any of its peers in row D, and D1 should not reflect any updates or withdrawals to any of its peers in row C. There are many ways to accomplish these stipulations, but a common method is to create filters on the routers at rows A and E, the leaf nodes or ToR switches, so only BGP updates with an empty AS path (^\$) are permitted, and to place all the routers at the spine routers (such as B and D) within a single pod in the same AS.

With these changes, BGP is essentially converted into “Fancy *Routing Information Protocol* (RIP),” and you can reduce the time required to withdraw a route (or move it from one place to another in the fabric) to about 1 minute in large-scale fabrics. It is possible to modify BGP to converge more quickly, but doing that returns the discussion to the first argument discussed previously—is creating a single protocol to solve all problems really the right answer? When is the complexity of the BGP code “complex enough” to start considering other options?

Let’s examine two other considerations before moving on to examining link-state protocols in DC fabric underlays. One of the advantages listed for BGP is that it has many different policy options, such as route filtering and tagging. If the underlay is really designed to provide undifferentiated IP connectivity, these policies do not seem like much of a real advantage. Policy, such as route tagging and filtering, should be moved to the overlay—which is most likely going to be BGP anyway.

A final point is that transit providers separate infrastructure and customer routes to split these two kinds of information into different failure domains. One misunderstanding about failure domains is they must be “absolute” and “complete,” where the two failure domains are completely decoupled at every point, if they are effective. They are not, however, always effective because it is likely impossible to build networks out of completely decoupled failure domains. Instead, it is a matter of tradeoffs. How much gain is there in separating these two kinds of information in this way, versus how difficult is it to separate these two kinds of information, and how much deoptimization is likely to occur?

In a DC fabric, separating the infrastructure routes of the underlay from the “customer” routes in the overlay is a legitimate way to form two different failure domains. These two failure domains might be somewhat tightly coupled, but they are still two different failure domains.

Separating the routes this way also creates multiple administrative domains, leaving open the possibility of allowing “customers,” or workload processes, to control some aspects of the reachability information in the overlay without the risk of causing problems in the basic IP connectivity the underlay provides.^[7]

Link State in the Underlay

Link-state protocols, like BGP, are also widely implemented and understood. Every commercial routing stack and many open-source routing stacks—including an implementation of *Open Shortest Path First* (OSPF) or IS-IS—are mature, well tested, and widely deployed. However, most of these implementations are not optimized for use on DC fabrics. This section considers the positive aspects of using a link-state protocol on a DC fabric, some of the challenges operators face when deploying standard link-state protocols on DC fabrics, and realistic expectations for scale when using these unmodified implementations. The following sections address modified link-state protocols currently being designed and implemented, and the probable scaling characteristics of these implementations.

The first advantage link-state protocols have over BGP in DC fabrics is convergence speed—but the irony is link-state protocols are at their fastest where BGP is at its slowest, and vice versa. Link-state protocols are most challenged at scale during initial convergence because of the density of the topology through which flooding must take place. Considering the network in Figure 1, shown previously; when E1 originates a new *Link State Update* (LSU)—whether a *Label Switched Path* (LSP) fragment in IS-IS or a *Link-State Advertisement* (LSA) in *Open Shortest Path First* (OSPF)^[11,12], it sends the update to every router in row D. Every router in row D, in turn, sends the LSU to every router in row E, which then sends the LSU to every router in row D. The number of copies each fabric device receives depends primarily on timing, but in topologies of around 2,600 fabric devices, each one was observed receiving more than 40 copies of each LSU. Nonetheless, unmodified link-state protocols converge at their worst as fast or faster than BGP up to some scale, where scale includes both the number of devices (nodes in the *Shortest Path Tree*, or SPT) and the number of reachable destinations. To what scale? The number will vary, but 1,000 (or more) fabric devices with a 100,000 reachable destinations are not unreasonable within a single flooding domain (or area in OSPF terms) based on prior large-scale deployments. Optimizations will increase these numbers somewhat—though to what degree depends on many factors.

Where link-state protocols converge much faster than BGP is when a reachable destination either moves from one place on the fabric to another or is disconnected from the fabric entirely. From the perspective of IS-IS, any reachable destination changes are just changes in leaf connectivity, meaning the destination can just be removed from the SPT without running *Shortest Path First* (SPF). This process is called a *partial SPF*; it is extremely fast and requires minimal processing on each of the fabric devices.

The second advantage link-state protocols have over BGP in DC fabrics is topology *visibility*. Link-state protocols require each device to maintain a full view of the topology, which must be synchronized with every other router in the network (or rather flooding domain); this process is called the *Link State Database* (LSDB). To obtain a copy of the LSDB, you need only to connect to one (or two, if you are concerned with resilience) router connected to the fabric. This kind of information is useful for traffic engineering and traffic steering. Further, periodic snapshots of the network topology from the perspective of the control plane can be a useful mine of telemetry information.

The first challenge for link-state protocols in the DC fabric is scaling, mainly related to flooding. We will consider several ways to reduce the number of LSUs each device receives in the following sections, so we don't consider them here. Another problem often cited in this area is the impression that link-state protocols can drop or fail to deliver LSUs—that flooding is periodic, rather than reliable, and the period is long enough to allow significant problems to develop. All link-state protocols, however, use reliable transport to deliver flooded packets. For instance, IS-IS tracks whether each neighbor has received an LSU through acknowledgments and retransmits LSUs until they are acknowledged. IS-IS can also send a description of the entire database periodically to ensure a neighbor's LSDB is correctly synchronized. OSPF has similar mechanisms.

Two other challenges link-state protocols face are scaling the number of reachable destinations and the time required to run the SPF algorithm used to calculate the set of loop-free paths. Faster processors combined with well-designed and well-tested implementations of SPF, along with optimizations such as partial SPF, have largely mitigated these concerns up to much larger scales than many engineers realize. Link-state protocols will never scale to the same levels as BGP, but they will scale enough to support a large proportion of the DC fabrics operators will build.

This article considers three proposed methods to control flooding designed to allow link-state protocols to support dense large-scale topologies. The first is *Distributed Optimized Flooding* (distopt-flood), arguably the least complex of the three options. The second is a centralized flooding controller, and the third is *Routing in Fat Trees* (RIFT), which is essentially a modified link-state protocol designed specifically for spine-and-leaf fabrics.

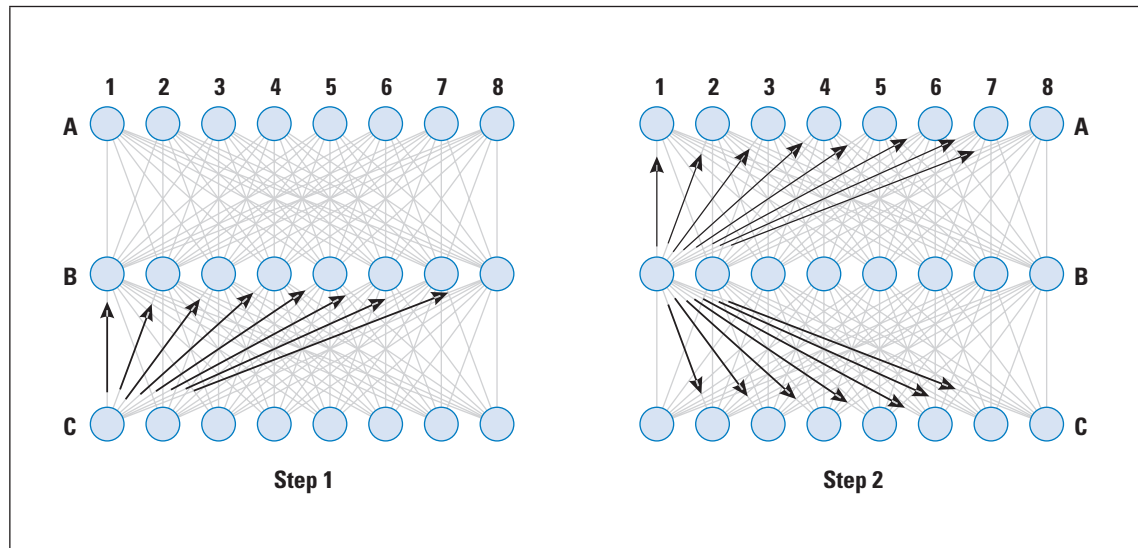
Distributed Optimized Flooding

Distoptflood outlines two optimizations to flooding, both of which work across all topologies and do not require a centralized controller of any kind. The first optimization is selecting a reduced set of reflooders^[8] when flooding an LSP (or fragment—LSP is used interchangeably with LSPF fragment in these explanations) by doing the following:

- Set all link metrics to 1.
- Calculate the shortest path tree.
- Group nodes with a cost of 2 by directly connected neighbors (nodes reachable with a cost of 1) through which they are reachable.
- Select a set of directly connected neighbors that can reach all nodes with a cost of 2.
- Remove any directly connected neighbors that are on the shortest path towards the origin of the change.

Figure 2 illustrates the flooding optimizations in a Clos fabric, while Figure 3 illustrates flooding optimizations in a Butterfly fabric.

Figure 2: Distributed Optimized Flooding in a Clos Fabric



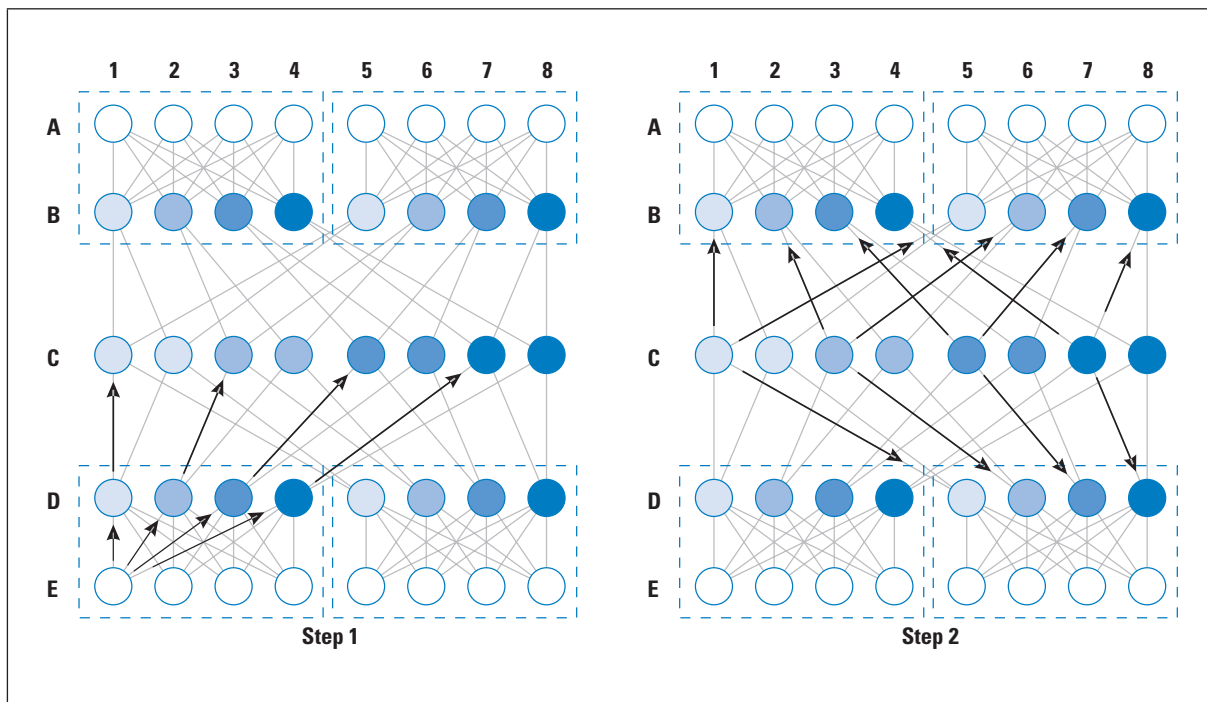
In the three-stage Clos (Figure 2), some change happens at C1; for instance, some network that was connected to C1 is disconnected. C1 calculates its two-hop neighborhood and determines it needs only to designate one of its neighbors, B1 through B8, as a reflooder. Let's assume it chooses B1 as the reflooder. C1 will send the LSP to C1 normally, and send the LSP to B2 through B8 using a link-local packet; these neighbors will receive the LSP and process it, but will not flood the changed LSP to their neighbors (they will not, in IS-IS terms, set their *Send-Receive* flag).

B1 will discover all of its neighbors can reach the same set of neighbors, and hence will select one connected neighbor as a reflooder; say B1 selects A1 as its reflooder. B1 will send the updated LSP to A1 through A8 and C2 through C8 using a link-local packet, so each of these routers will receive and process the change, but not reflood it. A1 will receive and process the update, but building its optimized flooding set will discover every one of its connected neighbors is on the shortest path towards the origin of the change, which is C1, so it will not reflood the update to any neighbors.

After about a second, each of the reflooders will send a *Complete Sequence Number Packet* (CSNP), which contains a description (or digest) of the local LSDB. If an IS notices a mismatch between its local LSDB and a neighbor's LSDB, it can send a *Partial Sequence Number Packet* (PSNP) requesting the retransmission of the missing information.

Figure 3 illustrates a slightly more complex five-stage spine-and-leaf fabric; while there is no “official” name for this configuration, it is often called a *Butterfly*.

Figure 3: Distributed Optimized Flooding in a Butterfly Fabric



Once again, let's assume a change occurs at E1, such as losing connectivity to a network (or reachable destination). In stage 1, E1 will build an LSP and calculate a set of reflooders that can reach its entire two-hop neighborhood—which is all of its neighbors in this case (D1 through D4). D1 through D4 will build a set of reflooders, which will include one of the two routers they are connected to in row C (C1 through C8). In stage 2, the selected reflooders in row C (C1, C3, C5, and C7) will determine a set of reflooders, which will be one spine router in each pod (such as A1, A5, and D5 for C1). The result: A1 through A8 and E5 through E8 will receive four copies of the changed LSP. None of the row A or row E routers will reflood the change because all their neighbors are on the shortest path back to E1, which originated the change. Depending on the timing of flooding, the number of copies of the changed LSP routers in rows A and E will likely be less than four.

A virtual testbed of around 2,600 routers configured as a butterfly showed this optimization decreased the number of LSPs each router received by a factor of 10 and doubled the initial convergence speed with more than 100,000 routes.

Because IS-IS runs over Ethernet natively and you can calculate the local system ID from an attached *Media Access Control* (MAC) (or Ethernet ID) address, you can run a modified IS-IS fabric with virtually no configuration on the fabric devices. You can assign locally calculated IPv6 addresses to the loopback address of each device, and use link-local IPv6 addresses to forward IPv6 traffic across the fabric. Forwarding IPv4 traffic would, of course, require an address plan and some form of automated configuration for loopback addresses. Fully autonomic configuration of this kind, however, can make troubleshooting issues and tracing flows through the fabric difficult. Therefore, current implementations do not include fully autonomic operations, so you must configure the system ID and the loopback address on each device.

A more controversial point is that using a control plane that runs natively at Layer 2 could improve security somewhat. A host that has been taken over or “pwned” by an attacker could not use the IP capabilities of the host to attack the operation of the fabric itself. Whether this feature results in an improvement in security is left to the reader (and operator!) to consider more deeply.

Centrally Calculated Optimal Flooding Trees

Centralized flooding management requires several modifications to link-state protocols, explained in “Dynamic Flooding on Dense Graphs”^[2]—a framework describing the changes required rather than a specific implementation. Rather than approaching the problem of optimally flooding information through a dense topology using local calculations, you can calculate a *flooding leader*, which then distributes an optimal flooding tree to all nodes in the fabric. Individual nodes would normally flood only along the designated tree, and then “by request” to resolve any flooding issues or to add links temporarily without impacting the flooding topology.

Each flooding domain (area in OSPF terms) must have a flooding leader; the draft suggests OSPF can elect this leader in much the same way as the *Designated Router* (DR) or a *Designated Intermediate System* (DIS) in IS-IS, both of which are used to reduce the amount of flooding required to synchronize a set of routers connected to a single broadcast link. While a single leader per flooding domain is required, the draft suggests each flooding domain should have multiple candidates, so the failure of the flooding leader does not cause an outage. This setup would be similar to the way OSPF elects a *Backup DR* (BDR) first, promotes the BDR to the DR role, and then elects a new BDR. In this way, a new flooding leader can “listen in” and be ready to take over the role of flooding leader if failure occurs.

A new *Type Length Value* (TLV) is added to IS-IS to enable the election of an area leader. An IS advertising this TLV is considered in the area leader election on all devices in the flooding domain. Rather than specifying the algorithm used to elect the flooding leader, an algorithm field is used to indicate how the flooding leader should be elected. Perhaps the simplest algorithm would be to elect the device with the highest (or lowest) priority, as advertised in the new TLV, and select among multiple advertisers with the same priority using the system ID (in the case of IS-IS).

When elected, the flooding leader calculates an optimal flooding topology. The flooding leader does not need any special information here; it already has a full view of the topology of the flooding domain through the synchronization of LSDBs required by normal link-state protocol operation. The precise calculation used is not specific in the draft, but a simple one might be to just use the shortest path tree as calculated by the flooding leader as the optimal flooding tree. The flooding tree does not need to be optimal from every point in the topology; it is not used to forward traffic, only to reduce flooding. The flooding tree also does not need to be perfect. A single device receiving two copies of a flooded link-state change might be less than optimal, but it will not cause routing loops or other significant network problems. In the same way, if a device fails to receive some new link-state information, the result might be suboptimal traffic flow. The normal flooding processes in OSPF and IS-IS will eventually catch the error (generally on the order of seconds) and fix the problem. Some optimizations, such as choosing only one link from a set of parallel links, and handling multiple nodes connected to a shared multi-access link, are considered in some detail.

After the topology is calculated, it must be advertised to the network devices in the flooding domain. IS-IS advertises it using a new TLV that is similar to the way link-states are already advertised. Each TLV contains a series of system IDs through which the flooding path passes. Similar additions to the OSPF protocol are described as well.

What “Dynamic Flooding on Dense Graphs”^[2] provides is a framework for a solution to flooding inefficiencies in link-state protocols, rather than a solution. In fact, you can advertise the use of distributed optimized flooding within a network by using the mechanisms in this draft. One specific algorithm for computing a dynamic flooding topology is described in “An Algorithm for Computing Dynamic Flooding Topologies”^[3] in this way:

“The proposed algorithm constructs a subgraph composed of small overlapping cycles. The base graph is denoted by $G(V, E)$, where V is the set of all reachable nodes in this area, and E is the set of edges. The subgraph to be computed is denoted by $G'(\{\}, \{\})$, which starts with an empty set of nodes and an empty set of edges.”

It is beyond the scope of this article to describe the precise way this algorithm operates or proposed alternatives.

RIFT

Routing in Fat Trees (RIFT) is a recent addition to this list, combining link-state and distance-vector concepts. Link-state-like operation is retained as information is transmitted up the fabric towards the ToF, while distance-vector-like operation carries reachability and topology information towards the edges of the fabric, the leaves.

Work on this new protocol started in IETF when the RIFT working group charter was approved in February 2018. The charter states:

“The Routing in Fat Trees (RIFT) protocol addresses the demands of routing in Clos and Fat-Tree networks via a mixture of both link-state and distance-vector techniques colloquially described as ‘link-state towards the spine and distance vector towards the leaves.’ RIFT uses this hybrid approach to focus on networks with regular topologies with a high degree of connectivity, a defined directionality, and large scale.”

The working group was chartered to create a protocol that will:

- Deal with automatic construction of fat-tree topologies based on detection of links.
- Minimize the amount of routing state held at each topology level.
- Automatically prune topology distribution exchanges to a sufficient subset of links.
- Support automatic disaggregation of prefixes on link and node failures to prevent black-holing and suboptimal routing.
- Allow traffic steering and rerouting policies.
- Provide mechanisms to synchronize a limited key-value data-store that can be used after protocol convergence.

According to the charter: “It is important that nodes participating in the protocol should need only very light configuration and should be able to join a network as leaf nodes simply by connecting to the network using the default configuration. The protocol must support IPv6 and should also support IPv4.”

Basic Operations

As briefly described earlier, RIFT combines concepts from both link-state and distance-vector protocols. A *Topology Information Element* (TIE) is used to carry topology and reachability information; it is like an OSPF LSA or IS-IS LSP. Figure 4 illustrates the advertisement of reachability and topology information in RIFT.

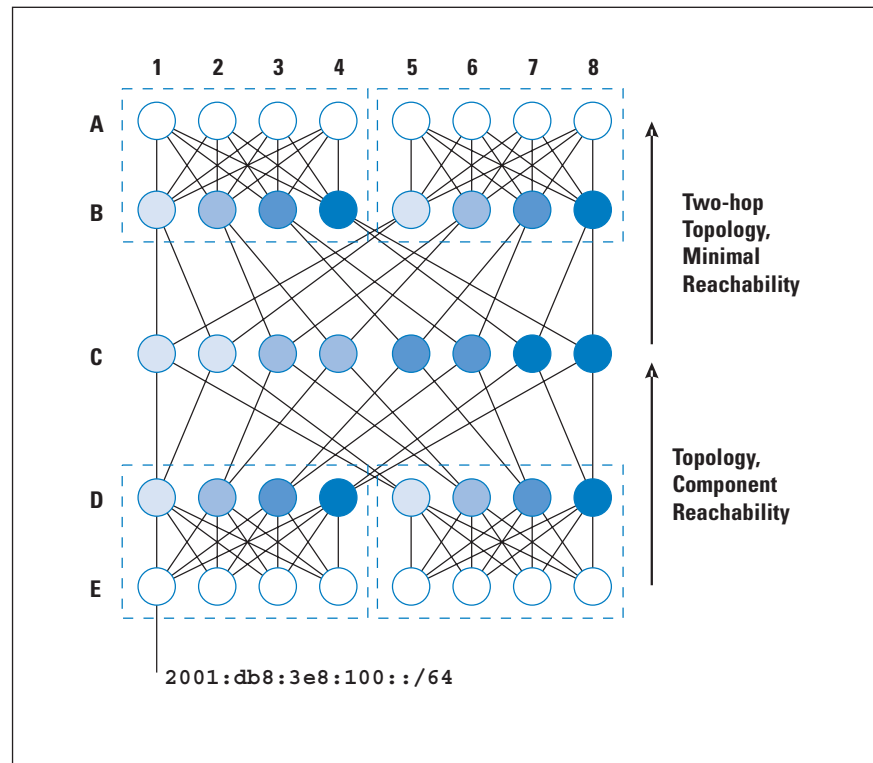
In Figure 4, E1, a ToR, advertises `2001:db8:3e8:100::/64` and its connections to D1–4 to D1–4. D1–4, in turn, refloods this information towards C1–8.

Rather than flooding the TIEs received from E1 back down the fabric, however, C1–8 advertises the minimal amount of reachability possible (normally this route would be a single default route) and their full set of neighbors to D5–8 and B1–8. When this process is completed:

- E2 will know about any locally connected destinations, four potential default routes from D1–4, and all the neighbors of D1–4 (the two-hop neighborhood for each connected neighbor).
- D5 will know about the neighbors and destinations connected to E5–8, the C1 two-hop neighborhood, the C2 two-hop neighborhood, and two possible default routes from C1 and C2.
- C4 will know about the destinations and neighbors reachable from all devices in rows A and E, and the neighbors connected to devices in rows B and D.

Using this information, the routers in row C can calculate a full SPT to discover the best path to each destination in the network. Routers in rows B and D will forward any traffic destined within the pod based on local information learned from the ToR switches, and all other traffic towards the default route learned from the ToF (row C). ToR switches will forward traffic using the default route learned from the spine stage above them, rows B and D.

Figure 4: RIFT Operation in a Butterfly Fabric



Automatic Disaggregation

In normal circumstances, devices other than those in the ToF stage rely on the default route to forward packets towards the ToF, while more specific routes are available to forward packets towards the ToR switches. What happens, however, if the C3→D2 link fails? B2, B6, and D6 will continue forwarding traffic towards C3 based on the default route being advertised in the TIE flooded by C3, but C3 will no longer have a route by which it can reach **2001:db8:3e8:100::/64**—so this traffic will be dropped at C3.

To resolve this problem, RIFT can use the two-hop neighbors advertised to all routers to automatically determine when there is a failed link and push the required routes along with the default route down the fabric. In this case, C4 can determine D3 should have an adjacency with D2, but the adjacency does not exist. Because of the failed adjacency, C4 can flood reachability to **2001:db8:3e8:100::/64** alongside the default route it is already sending to all its neighbors. This more specific route will draw any traffic destined to the **100::/64** route, so C3 no longer receives this traffic. The default route will continue to draw traffic towards C3 for the other destinations it can still reach.

Other RIFT Features

When ToF fabric switches are configured, fabric devices running RIFT can compute their fabric location and largely self-configure (there are exceptions for devices requiring Layer 2 support and leaf nodes in the topology). This self-configuration includes the use of IPv6 *Neighbor Discovery* (ND)^[13] to determine local IPv6 addresses, so no addressing plan or distribution protocols are required for pure IPv6 operation. If native IPv4 forwarding is required in the underlay, those addresses must be managed and configured in some way.

RIFT also offers the ability to perform unequal-cost load balancing from the ToR towards the ToF. Since each node has only the default route, and the stages closer to the ToF have more complete routing information, it is not possible for the ToR to cause a routing loop by choosing one possible path over another, or unequally sharing traffic along its available links.

Conclusion

BGP has been and will continue to be an important option for DC fabric underlays for many years to come. BGP may eventually offer some of the interesting features link-state protocols already offer, such as faster convergence and closer-to-autonomic deployment. On the other hand, some features of a link-state protocol, such as the ability to get a complete view of the entire topology from a single place—pulling a copy of the LSDB—are going to be very difficult to replicate in BGP, and the BGP convergence speed is always likely to lag behind a link-state protocol.

Table 1 summarizes many of the differences between the options outlined here.

Table 1: Differences Between Modified BGP, Modified IS-IS, and RIFT

Feature	BGP (Modified for DC Fabrics)	IS-IS (Modified for DC Fabrics)	RIFT
Peer Discovery	Partial	Yes	Yes
Automatic Tier Calculation	No	Potentially	Yes
Mis-Cabling Detection	No	Capability in progress	Yes
Fabric Addressing	Loopback address, peering; can be reduced with protocol modifications; can be automated	System ID, loopback address; can be automated or locally calculated	ToF state and others; can be automated
Aggregation; Default only on ToR and Below	Manually configured	No	Yes
Scales to Underlay Routing on Host	Yes	Depends on fabric size and implementation	Yes
High <i>Equal-Cost Multi-Path</i> (ECMP) Fanout Support	Yes	Yes	Yes
Unequal-Cost Load Sharing	Yes (in some implementations)	No	Yes
Full View of Topology	No	Yes	Yes (in the ToF)
Carry Opaque Configuration Data	No (can carry opaque information through Communities)	No (can carry opaque information through Tags)	Yes
Drain Node without Disruption	Yes	Yes	Yes
Automatic Disaggregation	No	No	Yes
Fast Convergence Speed	Partial (Depends on event type)	Yes	Yes
Security Includes Origin Validation and Replay Protection	Origin validation could be implemented, but heavy weight; no replay protection	No	Yes
Initial Implementation	Simple	Moderate	Complex
Overlay Support	Assumes single protocol (eBGP underlay, iBGP/eVPN overlay)	Assumes eVPN overlay	Supports eVPN overlay, can operate pure Layer 3 fabric with no overlay to the workload
Support for General Topologies (not just DC fabrics)	Yes	Yes	No

Link-state protocols offer a different set of tradeoffs than BGP does; operators would do well to consider the link-state options described here as strong alternatives to using BGP for DC fabrics underlays.

References and Further Reading

- [0] Petr Lapukhov, Ariff Premji, and Jon Mitchell, “Use of BGP for Routing in Large-Scale Data Centers,” RFC 7938, August 2016.
- [1] Russ White, Shraddha Hegde, and Shawn Zandi, “IS-IS Optimal Distributed Flooding for Dense Topologies,” Internet Draft, Work-in-Progress, September 2019, **draft-white-distoptflood-01**.
- [2] Tony Li, Peter Psenak, Les Ginsberg, Huaimo Chen, Tony Przygienda, Dave Cooper, Luay Jalil, and Srinath Dontula, “Dynamic Flooding on Dense Graphs,” Internet Draft, Work-in-Progress, November 2019, **draft-ietf-lsr-dynamic-flooding-04**.
- [3] Sarah Chen and Tony Li, “An Algorithm for Computing Dynamic Flooding Topologies,” Internet Draft, Work-in-Progress, March 2020, **draft-chen-lsr-dynamic-flooding-algorithm-00**.
- [4] Alankar Sharma, Dmitry Afanasiev, Tony Przygienda, Bruno Rijsman, and Pascal Thubert, “RIFT: Routing in Fat Trees,” Internet Draft, Work-in-Progress, March 2020, **draft-ietf-rift-rift-11**.
- [5] Justin Meza, Tianyin Xu, Kaushik Veeraraghavan, and Onur Mutlu, “A Large Scale Study of Data Center Network Reliability.” In *Proceedings of the Internet Measurement Conference 2018*, Association for Computing Machinery, 2018. <https://doi.org/10.1145/3278532.3278566>.
- [6] “Folded” unfortunately has two distinct meanings in spine-and-leaf networks. The original spine-and-leaf design, the Clos, was considered unidirectional, in that circuit setup proceeded in one direction through the fabric, and the resulting fabrics were nonblocking. Using a spine-and-leaf for bidirectional packet-switched traffic “folds” the fabric. However, folding also means drawing the fabric with the topmost tier at the top of the diagram and all the leaves along the bottom of the diagram.
- [7] Note that in some implementations of BGP, the iBGP and eBGP I/O paths are handled separately, making iBGP and eBGP either closer to, or fully, two separate failure domains. You should consider this point when determining which implementation to deploy in a DC fabric when BGP is used as the underlay protocol.
- [8] The “first flood” to build an initial LSDB on which the two-hop neighborhood can be calculated is performed in the normal way; there is no optimization of this initial flood of topology information.

- [9] Alvaro Retana and Russ White, “Optimizing Link-State Protocols for Data Center Networks,” *The Internet Protocol Journal*, Volume 16, No. 2, June 2013.
- [10] Yakov Rekhter, Susan Hares, and Tony Li, “A Border Gateway Protocol 4 (BGP-4),” RFC 4271, January 2006.
- [11] John Moy, “OSPF Version 2,” RFC 2328, April 1998.
- [12] Dennis Ferguson, Acee Lindem, and John Moy, “OSPF for IPv6,” RFC 5340, July 2008.
- [13] William Allen Simpson, Thomas Narten, Erik Nordmark, and Hesham Soliman, “Neighbor Discovery for IP version 6 (IPv6),” RFC 4861, September 2007.
- [14] Wikipedia entry for “Clos network”:
https://en.wikipedia.org/wiki/Clos_network

RUSS WHITE began working with computers in the mid-1980s, and computer networks in 1990. He has co-authored more than forty software patents, participated in the development of several Internet standards, helped develop the CCDE and the CCAR, and worked in Internet governance with the Internet Society. Russ has a background covering a broad spectrum of topics, including radio frequency engineering and graphic design, and is an active student of philosophy and culture. Russ is a co-host of the *History of Networking* podcast, hosts the *Hedge* podcast, serves on the Routing Area Directorate at the IETF, co-chairs the BABEL working group, and serves on the Technical Services Council as a maintainer on the open-source *Free Range Routing* project. His most recent works are the book *Computer Networking Problems and Solutions*, *Network Disaggregation Fundamentals* video training, and *Abstraction in Computer Networks* video training. E-mail: russ@riw.us

MELCHIOR AELMANS is Lead Engineer Cloud Providers at Juniper Networks, where he has been working with many operators on the design, security, and evolution of their networks. He has over 15 years of experience in various operations, engineering, and sales engineering positions with cloud providers, data centers, and service providers. Before joining Juniper Networks, he worked with eBay, LGI, KPN, etc. Melchior enjoys evangelizing and discussing routing protocols, routing security, and internet routing and peering. He also participates in IETF and RIPE and is a board member at the NLNOG foundation. E-mail: melchior@aelmans.eu

So You Want to Sell Your IPv4 Address Block?

by David Strom

If your company owns a block of IPv4 addresses and is interested in selling it, or if your company wants to purchase additional addresses, now may be the best time to do so. As readers of *The Internet Protocol Journal* (IPJ) are well aware, the number of available IPv4 addresses has been steadily dwindling, to the point now that many of the *Regional Internet Registries* (RIRs) are no longer assigning them. It may be a good time to look at the used-address marketplace. This arena could be a new corner of the Internet for you, so this article can help you understand what is going on and prepare you to do business in it.

For sellers, a good reason to sell address blocks is to make money and get some use out of an old corporate asset. If your company has acquired other businesses, particularly ones that have assets from the early Internet pioneers, chances are you might already have at least one range that is gathering dust, or is underused. Think of this idea of selling blocks as similar to how your company might decide to sell or release its unused real estate. “Many companies have millions of unused IP addresses,” said Vincentas Grinius of Heficed, an address leasing vendor. “They have been holding on to them for future growth or to save as a strategic asset.” Now might also be a good time to sell since prices are starting to level off, according to several brokers that I spoke to (of course, they have a vested self-interest), and the practice is becoming more accepted.

If you’re a buyer, it is also a good time for you, as a way to extend the life of your enterprise IPv4 equipment for a few more years. It is particularly true if your business has resisted a full IPv6 deployment or you can’t easily upgrade your legacy endpoints.

Until recently, the used-address marketplace hasn’t had the best of reputations. Many of us imagine that getting a used-address block from a broker is like buying a cheap used car. Grinius told me that used addresses used to be thought of “as akin to *Hustler* magazine, something folks were ashamed of having in their possession.” But things have gotten more legitimate: in addition to the used-car metaphor, you should also add the digital equivalence of a title insurance company and an accident reporting service like Carfax to establish more of a trusted exchange among buyer, broker, and seller.

Certainly the used-address market is thriving and quite competitive: now we have dozens of block brokers and at least three block lessors (IPv4 Market Group, Prefix Broker, and Heficed) that have solid business operations to help match buyers and sellers.

I have owned my own Class C block since 1993, and it seemed like an opportune time to sell it when IPJ’s editor asked me to write about the used IPv4 marketplace.

So let's first review the history of the IPv4 address depletion and how RIRs work in terms of address allocation before we get into the specifics of how the broker/resale space works. Along the way I will offer my own comments about my experience in selling my own block, and what I learned that can help you decide whether you want to become a buyer, a seller, or a lessor of your own block.

Address Transfer Reference Library

Perhaps the best source of information about IPv4 address depletion, myths (such as changing out customer routers is easy and ISPs still have plenty of IPv4 addresses) and tools about IPv6 transition are available in the back issues of IPJ itself, including articles that Geoff Huston of the *Asia Pacific Network Information Centre* (APNIC) wrote. Following is a guide to the most useful pieces, in IPJ and elsewhere. Note that I shared my thoughts with Huston prior to publication, and have woven in some of his remarks in this summary.

- An IPJ June 2003^[1] article reviews the early stages of IPv6 and includes some early myth busting by Huston, such as IPv6 has innately better security, *Quality of Service* (QoS), and mobility support. In one article, Huston says that “With a continuation of current policies it would appear that IPv4 address space will be available for many years yet.” He was right, just perhaps not in the way that he originally intended. In a recent e-mail, Huston told me, “At the time there was a common expectation that the adoption of IPv6 was meant to complete before the IPv4 pool had exhausted itself.”
- Four years later in an IPJ September 2007^[2] article, Huston talks directly about the state of IPv4 address depletion, and has models that (accurately as it turned out) predicted full depletion by 2011. At that time, address exhaustion was pretty much inevitable. In the article, Huston stated that IPv6 didn't have a very compelling business case and that the use of *Network Address Translation* (NAT) in IPv4 is far easier, a claim you could still make today.
- An IPJ March 2011 special issue on the IPv6 transition^[3] includes commentary on *World IPv6 Day* held June 2011 and a history of the address exhaustion of IPv4. “The stock of [IPv4] addresses is facing imminent depletion,” Huston wrote in that issue. By then, APNIC had exhausted its IPv4 address pool. “Most of the actors in the Internet are unsure about what needs to be done [to make the v6 transition], from the largest of the service providers down to individual end users,” he wrote. That issue is worth reviewing because it has a lot more helpful information about making the IPv6 transition.
- A December 2019 presentation by Huston about IPv6 is also worth reading^[4] He discusses current pricing trends on block sales and predicts that by the time we run out of IPv4 addresses we will have outgrown IPv6: “We didn't need it back when it was first proposed and we still don't need it now.”

Huston mentioned in a recent e-mail to me that it has been “nine years after the initial exhaustion point and IPv6 is still used by less than a quarter of the Internet and IPv4 remains the mainstay of the Internet. It was easier for the industry to change the entire architecture of the Internet than it was to universally adopt a new IP protocol.”

- A nice historical review of the development of RIRs is available in the December 2001 IPJ.^[5] It covers *Classless Inter-Domain Routing* (CIDR), subnetting, and supernetting.
- A white paper from Eric Bais^[6] has loads of practical advice for address transfer, written from the perspective of a broker in the *Réseaux IP Européens* (RIPE) region who both sells and leases blocks.

Historical Review

I first wrote about the depletion of the IPv4 address space when I was editor-in-chief of *Network Computing* magazine back in the early 1990s. Alas, that article is no longer accessible online. I remember it vividly because it got an amusing comment from my father, who never really understood technology but thought it would be funny if I were to leave my job and become an address broker. Needless to say, it was just a passing but prescient thought.

Back in these early days of the commercial Internet, Jon Postel personally and manually assigned IP address ranges. Usually he did it within moments of receiving an e-mail request, and that is how I got my /24 block. Obviously it didn’t scale after the Internet caught on. One of the first to sound the alarm was Frank Solensky, who published his predictions for various run-out dates in 1990 during the 18th meeting of the *Internet Engineering Task Force* (IETF).^[7] See Figure 1.

Figure 1: Solensky’s Original Estimates of Address Exhaustion

<u>Depletion Dates</u>	
• Assigned Class “B” network numbers	Mar. 11, 1994
• NIC “connected” Class B network numbers	Apr. 26, 1996
• NSFnet address space *	Oct. 19, 1997
• Assigned Class “A-B” network numbers	Feb. 17, 1998
• NIC “connected” Class A-B network numbers	Mar. 27, 2000
• BBN snapshots *	May 4, 2002
* all types: may be earlier if network class address consumption is not equal	

The basic “Goldilocks” issue is that for the average business looking to get online, 250 addresses for a class C block is too little and 65,000 addresses for a class B block is too many. Numerous technical approaches have been proposed, including classless addressing (RFC 1918^[8]), NAT, elimination of assigning static addresses to dial-up users, and changes to routing protocols. But the real solution was inventing IPv6 to increase the overall address space. During the early 1990s, the larger blocks of A and B ranges were already being rationed, given that Postel by then had previously assigned many of these blocks.

While the IPv4 addresses were being depleted, three of the RIRs were created through RFC 1366^[9], modified by RFC 1466 in 1993^[10], and further refined a few years later in RFC 2050^[11]. Now there are five of them:

- *African Network Information Centre* (AFRINIC) serving Africa
- *Asia Pacific Network Information Centre* (APNIC) serving parts of Asia and the Pacific region
- *American Registry for Internet Numbers* (ARIN) serving North America and parts of the Caribbean
- *Latin America and Caribbean Network Information Centre* (LACNIC) serving Latin America and parts of the Caribbean
- *Réseaux IP Européens Network Coordination Centre* (RIPE NCC) serving Europe, parts of central Asia, and the Middle East

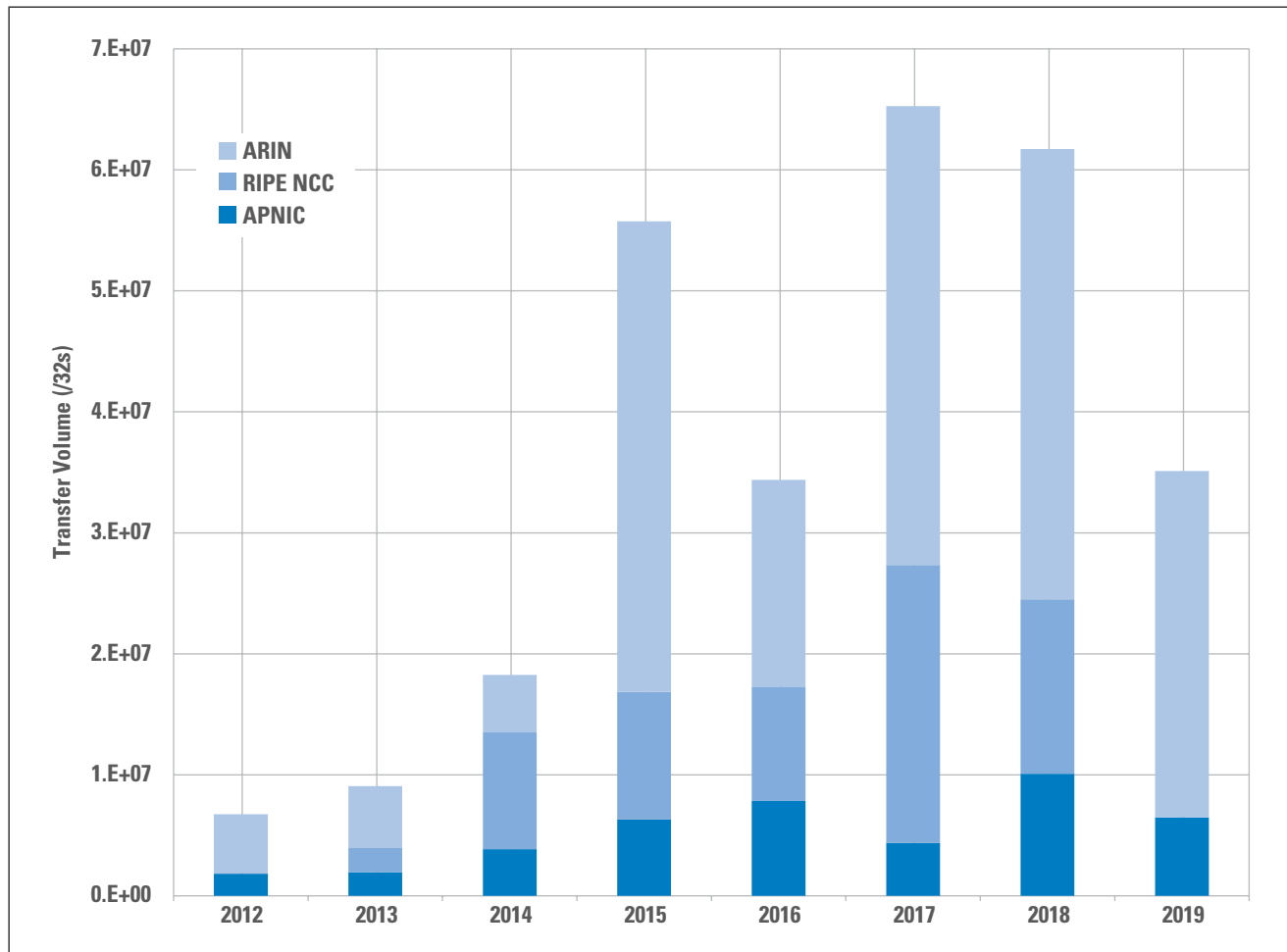
By February 2011, the last remaining common blocks of IPv4 addresses were fully allocated to the RIRs. In an article in IPJ, Raúl Echeberría, Chairman of the *Number Resource Organization* (NRO), the umbrella organization of the five RIRs, was quoted as saying, “It’s only a matter of time before the RIRs and ISPs must start denying requests for IPv4 address space.”^[3] Today almost every block is assigned to some entity. AFRINIC has the most available, and APNIC has a few smaller blocks left. RIPE made its last /22 block assignment in November 2019.^[12]

The Rise of the Used-Address Marketplace and RIR Supervision

Perhaps the origin event for the used-address market was when Microsoft purchased Nortel’s inventory of more than 600,000 individual IPv4 addresses for US\$7.5M in 2011. (Well, that isn’t a completely accurate statement, but it does appear that Nortel’s address pool was the main corporate asset.) Since then, tens of millions of addresses have been transferred^[13] per year, as you can see in Figure 2.

In the last decade, the RIRs have played an increasing role in these transfers. In the references I have the direct links to the current transfer policies for each registry.^[14] Note that some RIRs have more precision and transparency about their process, along with higher thresholds, than others to prove existing ownership of an address block.

Figure 2: Address transfers from 2020 statistics compiled by Geoff Huston



But this system wasn't perfect by any means: block ownership questions weren't easily resolved within a single registry, organization records were full of stale data or listed businesses that were no longer operating entities, and spammers could pollute address blocks by clouding any resale opportunities. Also, many address blocks (such as the one that I owned) pre-date the establishment of RIRs, what they call "legacy resources." How the RIRs deal with these assignments is a challenge, particularly as businesses are no longer around, and tracing the lineage from the original Postel assignment to a current stakeholder can involve some detective work. The question is, who should do the detecting? That isn't a simple question to answer, as you'll see.

Part of the problem was WHOIS itself, the primary domain and block ownership query tool. However, WHOIS is far from perfect. First, its responses differ depending on the data being queried, the RIR in charge of that block, and whether the block owner has provided accurate and up-to-date information or deliberately hidden these details.

But another part of the problem is that the Internet community has made changes to the display of information from WHOIS queries. Changes were necessary because of privacy concerns (from various changes to regulations around the world) and from spammers abusing WHOIS to drive legitimate business owners into hiding their details. I have placed the links to the different RIR WHOIS pages in the reference section if you want to compare them.^[15]

If you were to examine my own /24 block before I began writing this article, you would see:

Organization: David Strom, Inc. (DAVIDS-3)
RegDate: 1993-05-21
Updated: 1996-04-18

The address used for my DAVIDS-3 organization is a New York corporation that is no longer in business. And the point of contact listed is an engineer at an ISP that I used to register the block that is also no longer in business. My challenge: I had to prove that the David Strom Inc. that did business in New York was the same David Strom Inc. that is now doing business in Missouri. Other than finding the plane ticket that I used in my move, I wasn't sure what else I could do to document the "asset transfer" that ARIN was going to eventually ask for.

Thus began my own journey to correct this information and get it ready for resale. The process involved spending a lot of time studying the various transfer webpages at ARIN, calling their transfer hotline several times for clarifications on their process, and paying a \$300 transfer fee to start the process. ARIN staff promises a 48-hour turnaround to answer e-mails, and that can stretch out the time to prepare your block if you have a lot of back-and-forth interactions, as I did.

Enter the Block Broker

This discussion brings us to the modern era (say after 2012) and the IP block-broker marketplace. The goal was to try to make it easier for these transfers, and at the same time improve trust among all parties. As I said earlier, we now have many block brokers doing business. The broker's service (for either selling or leasing a block) is somewhat similar and involves these basic steps:

1. You need to register your business with the broker, a process that involves just answering a few basic questions and creating a login ID so you can interact with them via their various web-based forms and forums and e-mail.
2. Next, if you are a seller, you sign a mutual *Non-Disclosure Agreement* and then list your block that you want to sell. Some brokers have a variety of sales methods, including open and closed auctions and the ability to "buy now." If you are a buyer, you can start browsing the blocks that are available on the open auctions, and participate in the auction. If you have ever bought or sold any physical object via an online auction, you should be familiar with this process.

3. After you select a buyer for a particular block, you request the funds and place them in escrow, and then close the auction.
4. The broker's support team arranges for the transfer with the relevant RIR(s). As a buyer, you will then pay the fees directly to the RIR(s) for the transfer. Each RIR has a different way to calculate fees, ranging from free for RIPE to thousands of dollars, depending on the size of the block. (See Figure 3.)

Figure 3: The Different Fees Each RIR Charges for Transferring Addresses

RIR	Transfer Fee Amount
ARIN	\$300 USD
RIPE	\$0 USD
APNIC	20% of the annual fee for the # of IPv4 addresses being transferred
LACNIC	Initial payment of \$200 USD Smaller than a /19 - \$1,000 USD /19 and larger - \$1,500 USD
AFRINIC	Smaller than /22 - \$0 USD /22 to /20 - \$1,750 USD /20 to /18 - \$2,000 USD

5. Finally, the transaction closes and the block control and the funds released from escrow, minus any commission from the broker, are transferred to the buyer or leaser. Here is where things get interesting. The commissions aren't transparent: you have to get far enough down the process before you can find out what they are; the brokers set up the process this way deliberately so you can't shop around for lower fees. Still, there is a place for brokers, since "nothing is more frustrating than trying to get paid in a country of which you don't know the legal system nor have local representation. Using an escrow makes things easier for all parties involved," says Eric Bais of Prefix Broker.^[6]

One other caveat for block leases: the lessor and lessee have a more intimate and longer-term relationship than if you are buying and selling the block outright, because ultimately the "landlord" business is still responsible for the reputation of the folks who are using your IP addresses. In other words, renting out your space also carries a certain risk to the lessor: just like rentals in the physical space, owners (or landlords) are responsible for their property. If you have a bad tenant who trashes your space, your reputation will suffer. This reality places a bigger burden of trust on the broker to ensure a proper tenant.

Three RIRs list brokers on their websites. They all have somewhat different contact information and number of brokers:

- APNIC has 22 listings^[16], with contact names and phone and skype numbers.
- RIPE has 76 listings^[17], with links to their contact webpages.
- ARIN has 29 listings^[18], with contact names and phones and the date the broker registered with ARIN.

All of these RIRs try hard to indicate that their listings are not a recommendation, just awareness of their businesses. RIPE says its listing, for example, is of brokers who have agreed to conduct their business honorably, but no one checks on the brokers after they are listed to see if they have actually lived up to their promise. That is worth remembering. As the old saying goes, “on the Internet, no one knows if you are a dog.”

If you are starting out in the used marketplace as I was, I recommend that you examine these RIR webpages carefully. Just having these lists of brokers is nice, but if you are going to sell or buy a used block you will find it frustrating to find the right broker for your situation. The biggest issue is that there are no fixed rules for buying, selling, and leasing used addresses. Unlike the used-car industry, there is no overall supervision or agreement on what constitutes the *quality* of an asset. As you can see from the five-step process cited previously, uncertainties and potential problems can arise at every step.

One other thing worth mentioning should be obvious but isn't: The only entities that can play are businesses. If you own a block as an individual, you will first have to transfer ownership to a business to proceed. You notice my ownership is my S-corporation (with my name; that helped me in the transfer process from my New York corporation to my Missouri corporation). Had I initially registered for my block as an individual, I might have had to work harder to prove my identity.

Important Caveats for the Transfer Process

Following are some of the complicating factors to watch out for as you begin your own transfer journey:

First, choosing whether to buy or lease a block can be tricky, and it depends on how many addresses you need and for what purpose. More details will follow, but you need to make this very basic decision before doing anything else, and often you won't have as much data as you might like.

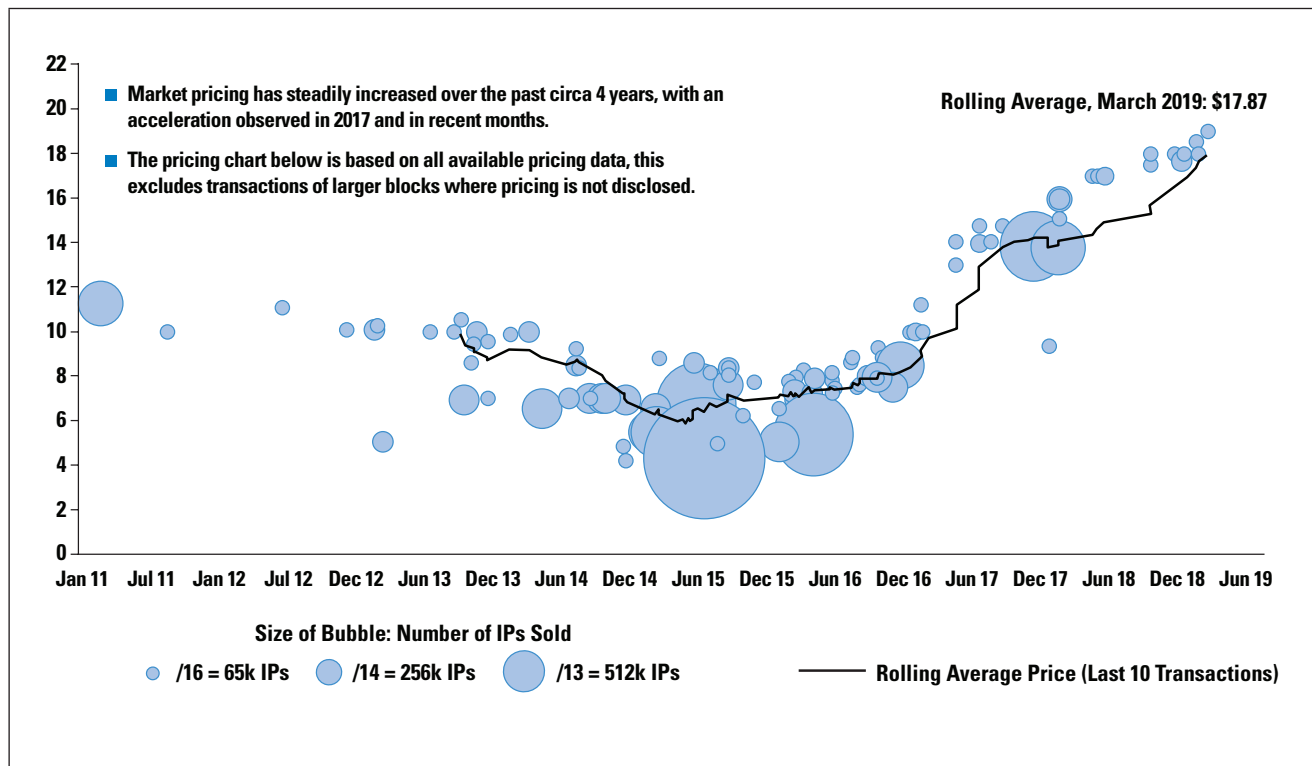
Unlike the used-car industry, there aren't any generally accepted practices or guides to making a tradeoff between buying and leasing. Of course, one aspect is the overall cost, and to estimate it you have to know your time horizon. If you are a buyer, do you need the block for a few years or a few months? Can you eventually migrate the endpoints to IPv6 using these addresses?

If you are a seller, do you want to dispose of the block and make a quick addition to boost your current year's balance sheet, or do you want to invest in a steady rental income over time? As a renter, you are also betting on a particular price curve over the terms of the lease that may or may not materialize. Now imagine that you are having this conversation with your Chief Financial Officer, who may or may not understand the various subtleties about the used-address marketplace.

You should base part of your choice of whether to rent or buy on the size of the block involved. Some brokers specialize in larger blocks and some won't sell or lease anything less than a /24, for example. "If you are selling a large block (say a /16 or larger) you would need to use a broker who can be an effective intermediary with the larger buyers," said Geoff Huston in an e-mail to me. Again, knowing that your broker has listed prior transactions can help you make a more informed decision. Not all brokers have pricing transparency, and many brokers are more circumspect about pricing.

IPv4.Global is one that does list their own prior auction sale data^[19], for example. Another broker, IPv4 Market Group, has assembled the overall pricing chart shown in Figure 4 from March 2019^[23]. There is no way to independently verify this information, but at least these examples show you how the market has evolved over the past decade.

Figure 4: IP Address Block Pricing Trends over Time



In early January 2020, /24s were selling at around US\$20–24 per IP address, or US\$5,000–6,000 for the entire block. Rental prices varied from 20 cents to US\$1.20 per month per address, meaning at best a 2-year payback and at worst a 10-year payback when compared to sales. I decided to sell my block: I wanted the cash, and didn't like the idea of being a landlord of my block any more than I liked being a physical landlord of an apartment that I once owned. You'll also want to ensure that the RIR that is responsible for your block recognizes the broker you eventually choose.

Second, there is no guarantee that any of these brokers is reputable and will actually deliver the goods, or even if the RIR listings and contacts for the broker are still accurate. There is no easy way to vet their operations, or even agree on overall metrics to be used as part of the vetting process. Unless you know them personally, or know someone who does, chances are the names of the brokers on the RIR lists will require additional research for you to decide whom you should use to sell your block. You can look to see their registration data with ARIN, if ARIN controls your block.

One possible vetting strategy is to inquire how the broker is involved in the various Internet governance committees in your region, or at least examine their posted attendee lists. The hypothesis for this strategy is that broker representatives who attend IETF, RIR, and network operator meetings such as The *North American Network Operators' Group* (NANOG)^[20] are more reliable than those that have never been to any of these meetings. (For example, PrefixBroker.com claims on their website that they helped author the RIPE transfer rules.)

IPv4 Market Group has a list of questions^[21] to ask a potential broker, including if they will represent only one side of the transaction (most handle both buyer and seller) and if they have appropriate legal and insurance coverage. I found that a useful starting point.

Some brokers are also involved (either as other lines of business at their own company or as a subsidiary of a larger corporation) in other network- and Internet-related businesses, such as hosting and cloud services, while others operate in real estate development and intellectual property litigation. That may be relevant, or it could cloud your evaluation if the quality of these other businesses differs from that of the brokerage.

One of the reasons I went with IPv4.Global/Heficed was because of their transparency in terms of showing me the active auctions and past sales of their blocks right on their web homepage, and they e-mailed me periodically with the active and closed auctions.

Third is how you vet the other party in your transaction. In other words, if you are a seller, what process do you use to know your buyer, and vice-versa?

You might want to consider longer-term contracts for rentals (such as 3 years) for stability and also to minimize the movement of their tenants. “I would be somewhat worried if the broker did not undertake some diligence steps directly to validate the credentials of the seller,” said Huston in an e-mail to me.

The final part of the transfer process is to understand the condition of the actual address block itself. There is no guarantee that a used block isn’t tainted with spammers or used for other less-than-legal activities. “There are no established standards of conduct, little transparency, and even less accountability,” wrote Marc Lindsey in 2018 for a blog post on *CircleID*.^[13] “Many participants in the market struggle to define, from a legal perspective, what is being bought and sold.” He also has several suggestions on vetting the other party in the transaction that are worth reviewing.

Most brokers will state that they examine prior ownership of their blocks to ensure they are spam-free and to eliminate the potential of being used for other shady dealings. The trick is understanding what tools they use to convince you of this claim. For example, some brokers require you to check the blacklists (such as those maintained at Cisco Talos, Hetrixtools.com, and IP-score.com) on your own to ensure that your block isn’t listed there. IPv4 Market Group offers a blacklist cleaning service^[22] that examines 90 blacklists. While charges vary, to give you an idea, they quoted me \$2,000 as part of their selling services for my /24 block. IPv4.Global checks 20 different blacklists as part of their services.

However, identifying whether a block is on a blacklist and removing it from a list are two different matters. If it is listed, you will have to work on removal from the blacklists before you can lease it. According to Geoff, “Once an address is blacklisted it’s exceptionally hard to get it unlisted.” None of the brokers will give you a firm price on cleansing a block, because it depends on how many blacklists it appears on.

So what happened to my sale? It took 10 days to auction off my block. I worked with ARIN to transfer my ownership to my current corporation, and paid them a second fee of \$125 for dealing with my legacy ownership. I then worked with my broker to finalize the sale. The overall elapsed time from beginning to end was 1 month, including about a week of elapsed time to conduct the initial research and select the broker.

Summary

If all that seems like a lot of work to you, then perhaps you just want to steer clear of the used marketplace for now. But if you like the challenge of doing the research, you could be a hero at your company for taking this task on. Expect the entire process to take several months from start to finish, allowing for time to get your ownership in order (if you have a legacy block), navigate the legal and other corporate approvals, research your broker, and then actually execute the transaction.

References and Further Reading

- [0] Various authors, *ConneXions—The Interoperability Report*, Volume 8, No. 5, May 1994, Special Issue: IP: The Next Generation, available from The Charles Babbage Institute:
<http://www.cbi.umn.edu/hostedpublications/Connexions/index.html>
- [1] Geoff Huston, “Opinion: The Mythology of IPv6,” *The Internet Protocol Journal*, Volume 6, No. 2, June 2003.
- [2] Geoff Huston, “IPv4 Address Depletion and Transition to IPv6,” *The Internet Protocol Journal*, Volume 10, No. 3, September 2007.
- [3] Various authors, *The Internet Protocol Journal*, Volume 14, No. 1, March 2011. This special issue is devoted entirely to the addressing and transitioning topics.
- [4] Geoff Huston, presentation slides about current issues with IP addressing, December 2019.
<https://www.potaroo.net/presentations/2019-12-11-kismet-addresses.pdf>
- [5] Daniel Karrenberg, Gerard Ross, Paul Wilson, and Leslie Nobile, “Development of the Regional Internet Registry System,” *The Internet Protocol Journal*, Volume 4, No. 4, December 2001.
- [6] Eric Bais, “Transferring IPv4 Resources in the RIPE region,” June 2016. An ebook published by Prefix Broker.
<https://www.prefixbroker.com/ebook/>
- [7] Frank Solensky, Proceedings of the 18th IETF, 1990.
<https://www.ietf.org/proceedings/18.pdf> (see p. 67 of the PDF for his original predictions of address exhaustion)
- [8] Daniel Karrenberg, Yakov Rekhter, Eliot Lear, and Geert Jan de Groot, “Address Allocation for Private Internets,” RFC 1918, February 1996.
- [9] Elise Gerich, “Guidelines for Management of IP Address Space,” RFC 1366, October 1992.
- [10] Elise Gerich, “Guidelines for Management of IP Address Space,” RFC 1466, May 1993.
- [11] Kim Hubbard, Jon Postel, Mark Kosters, Daniel Karrenberg, and David Conrad, “Internet Registry IP Allocation Guidelines,” RFC 2050, November 1996.
- [12] RIPE press release, November 2019, “The RIPE NCC Has Run Out of IPv4 Addresses,”
<https://www.ripe.net/publications/news/about-ripe-ncc-and-ripe/the-ripe-ncc-has-run-out-of-ipv4-addresses>

- [13] Marc Lindsey, *CircleID* blog post July 2018. “An Insider’s Guide to the IPv4 Market – Updated,”
http://www.circleid.com/posts/20180710_an_insiders_guide_to_the_ipv4_market_updated/
- [14] Here are the links to the RIR webpages regarding their rules for transferring resources:
- <https://www.apnic.net/manage-ip/manage-resources/transfer-resources>
- <https://www.ripe.net/manage-ips-and-asns/resource-transfers-and-mergers>
- <https://afrinic.net/resources/transfers>
- <https://www.arin.net/resources/registry/transfers>
- <https://www.lacnic.net/1019/2/lacnic/resources-transference>
- [15] Here are the links to query the WHOIS resources at each RIR:
- AFRINIC Database:
<https://www.afrinic.net/whois-web/public/query>
- APNIC Database:
<https://wq.apnic.net/apnic-bin/whois.pl>
- ARIN Database:
<https://whois.arin.net/>
- LACNIC Database:
<https://lacnic.net/cgi-bin/lacnic/whois>
- RIPE Database:
<https://www.ripe.net/manage-ips-and-asns/db>
- [16] APNIC, Registered IPv4 brokers:
<https://www.apnic.net/manage-ip/manage-resources/transfer-resources/transfer-facilitators/>
- [17] RIPE, Brokers:
<https://www.ripe.net/manage-ips-and-asns/resource-transfers-and-mergers/brokers>
- [18] ARIN Registered Transfer facilities:
https://www.arin.net/resources/registry/transfers/stls/registered_facilitators/
- [19] IPv4.Global, prior auction pricing data:
<https://auctions.ipv4.global/prior-sales>
- [20] NANOG, attendees list of meeting #77:
<https://events.nanog.org/events/nanog-77/attendees-15-13b224d66c30422494a9627a6dcb6c94.aspx>

- [21] IPv4 Market Group, “Approved IPv4 Address Facilitator for Your IPv4 Needs, a Guide to Questions You Might Want to Ask Your Broker,”
<https://ipv4marketgroup.com/ipv4-market-group/>
- [22] IPv4 Market Group, blacklist removal service.
<https://ipv4marketgroup.com/broker-services/ipv4-blacklist-removal/>
- [23] IPv4 Market Group, “IPv4 Price Trends,”
<https://ipv4marketgroup.com/ipv4-price-trends/>
- [24] Prefix Broker: <https://www.prefixbroker.com/>
- [25] Heficed: <https://www.heficed.com/>
- [26] Richard Jimmerson, “On the ‘Misuse’ of the Internet Number Resource Transfer Market,” Team ARIN Blog, August 26, 2020.
<https://teamarin.net/2020/08/26/on-the-misuse-of-the-internet-number-resource-transfer-market/>

DAVID STROM has written several articles for IPJ, most recently on fileless malware in 2018. He was the founding editor-in-chief for *Network Computing* (USA) magazine and ran overall editorial operations for Tom’s Hardware.com. He is the author of two books on computing, including one as co-author with Marshall T. Rose on Internet messaging. He lives in St. Louis and can be reached at: david@strom.com or Twitter [@dstrom](https://twitter.com/dstrom).

In Memoriam: Yngvar Lundh

by Ole Jacobsen, *The Internet Protocol Journal*

Yngvar Gundro Lundh (March 19, 1932 – August 15, 2020) was my friend, mentor, and boss at the *Norwegian Defence Research Establishment* (NDRE). I first met Yngvar around 1976 when I was still in high school working on a report about computers and society. I worked at NDRE in Yngvar's micro-computer group through my military service and later during summer vacations while at university. Yngvar was the person who introduced me to the wonders of computers, and most of all to networking. NDRE had access to the first ARPANET connection outside of the United States. It was through this link (a *TeleType* connected to the NORSAR-TIP itself connected at 9.6 kbit/s to the ARPANET via satellite) that I met many friends in the US, ultimately leading to my employment at the Network Information Center at SRI International in 1984.

Yngvar played a major role in fostering technology development in Norway through his work at NDRE, as professor of informatics at the University of Oslo, as chief engineer at Norwegian Telecom, and as consultant on a variety of projects, including the first commercial electronic mail system in Norway.



Photo: Gisle Hannemyr CC BY-SA 3.0

His group designed and built Norway's first transistor-based computer, SAM, which you can see at Norsk Teknisk Museum in Oslo.

He was perhaps best known for his early work with the ARPANET and SATNET at NDRE. Yngvar Lundh and Pål Spilling were largely responsible for getting Norway connected to the Internet in the early 1980s.^[1,2,3,4]

I fondly remember Yngvar as a patient teacher, generous with his time and always willing to help with projects large and small. He played a major role in my university and career path, and I will very much miss his guidance and inspiration.

Yngvar had many hobbies, including gardening, bee keeping, wood working, ham radio (LA72C), and above all, sailing. After retirement, he moved from Skedsmokorset near Oslo to Tolvsrød near the coastal town of Tønsberg, allowing him easy access to his sailboat.

References

- [1] Yngvar Lundh, “A Slice of Norway’s Computing History,” *IEEE Xplore*, April-June 2018.
<https://ieeexplore.ieee.org/document/8415734>
- [2] Wikipedia article on Internet Pioneers:
https://en.wikipedia.org/wiki/List_of_Internet_pioneers
- [3] Pål Spilling and Yngvar Lundh, “Features of the Internet History, The Norwegian Contribution to the Development,” *Teletronikk* 3.2004. Available from:
<https://www.usit.uio.no/om/organisasjon/sst/stab/ansatte/bness/tilkoplet/web/7/src/pal-spilling-yngvar-lundh-features-of-the-internet-history.pdf>
- [4] Wikipedia article on Pål Spilling:
https://en.wikipedia.org/wiki/P%C3%A5l_Spilling
- [5] Wikipedia article on Yngvar Lundh (in Norwegian):
https://no.wikipedia.org/wiki/Yngvar_Lundh
- [6] Norwegian Defence Research Establishment:
<https://www.ffi.no/en>
- [7] Dag Andreassen, “Internett med norske pionerer,” Norsk Teknisk Museum,
<https://www.tekniskmuseum.no/21-nyheter/354-internett-med-norske-pionerer>
- [8] Tor Sverre Lande, “Nekrolog: Yngvar Lundh,” *Aftenposten*, August 27, 2020.
<https://www.aftenposten.no/personalia/i/1ngy8e/nekrolog-yngvar-lundh>
- [9] Yngvar Lundh, *Konstruksjon av integrerte kretser*, Universitetsforlaget, 1983, ISBN13 9788200066910.

OLE J. JACOBSEN is the Editor and Publisher of *The Internet Protocol Journal*, a quarterly publication on all aspects of Internet technology. He has been active in the computer networking field since 1976, when he joined the Norwegian Defence Research Establishment, an early ARPANET site. Ole holds a B.Sc. in Electrical Engineering and Computing Science from the University of Newcastle upon Tyne, England. He serves on the board of the *Asia Pacific Network Operators Group* (APNOG), which hosts the annual *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT) conference, and has served on several ICANN and IETF nomination committees. In his spare time, Ole organizes pipe-organ concerts and demonstration events. E-mail: ole@protocoljournal.org

Book Review

Transforming Information Security

Transforming Information Security: Optimizing Five Concurrent Trends to Reduce Resource Drain, by Kathleen Moriarty, ISBN-13 978-1839099311, Emerald Publishing Limited, July 2020.

When I was asked to write a short review about Kathleen Moriarty's book, I took a copy with me on my summer holiday. I usually try to stay away from work-related literature during vacation, but in this case it was well worth it. With some 200 pages packed with facts and information, this book requires a bit of concentration and focus. But you get a lot in return for the effort.

With her extensive background and expertise in security, Kathleen analyses five trends in the current security debate: End-to-End Encryption, Strong Session Encryption, The Evolution of the Transport Protocol Stack, Data-Centric Security, and More User Control.

With these trends in mind, Kathleen comes to the conclusion that we need a fundamentally different approach to network and information security. She promotes a more manageable system with a minimised, secure operating system and layered or hosted (authorised) applications on top of it. Vendors need to take more responsibility managing vulnerability, and should enable automated updates that users can trust.

Security has become increasingly complex and requires specialised knowledge and expertise. There is already a huge shortage of security practitioners. Training more people and buying additional security tools and products is not going to scale. It is increasingly challenging to secure our networks and keep them manageable at the same time. Only wide-scale adoption of end-to-end encryption, increased capabilities of the end-points, and a change of network architecture and security practices will help in the mid and end terms.

But Kathleen doesn't leave it at these high-level statements. The book is full of practical tips and provides a wide range of *Internet Engineering Task Force* (IETF) standards, guidelines, and suggested security frameworks that IT and security staff can find useful—the list of references in itself is very useful and encourages further reading.

Kathleen walks us through various aspects of information security: from threat detection and prevention, to the use of security control frameworks, to the need for more automation and the importance of sharing information with peers in the network and security community.

She also provides an overview of many standards and protocols such as IPv6, *Quick UDP Internet Connection* (QUIC), *Manufacturer Usage Description* (MUD), routing overlay protocols, *DNS over Hypertext Transfer Protocol Secure* (DoH), and *DNS over TLS* (DoT) to name only a few, and explains their relevance in the overall security landscape.

Some sentences are packed with information, and it is worth it to read them twice.

The book provides a peek into the hopefully not-too-distant future where applying a more holistic view on security, automation, and sharing relevant information will benefit the networking and security community.

—Mirjam Kühne, mir@zu-hause.nl

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. For more information, contact us at ipj@protocoljournal.org

Check your Subscription Details!

If you have a print subscription to this journal, you will find an expiration date printed on the back cover. For several years, we have “auto-renewed” your subscription, but now we ask you to log in to our subscription system and perform this simple task yourself. Make sure that *both* your postal and e-mail addresses are up-to-date since these are the only methods by which we can contact you. If you see the words “Invalid E-mail” on your copy this means that we have been unable to contact you through the e-mail address on file. If this is the case, please contact us at ipj@protocoljournal.org with your new information. The subscription portal is located here: <https://www.ipjsubscription.org/>

Announcement Regarding IPv4 Address Block 43/8

As many people know, I have dedicated much of my career to the development of research networks and network technologies in Japan and Asia. This included the WIDE (*Widely Integrated Distributed Environment*) Project, founded in 1985 for computer networking Research and Development. In the early days of the WIDE Project, we were aware of the exciting advent of the Internet, and I was often in contact with Jon Postel and other Internet pioneers, about how it could be brought to Asia.

In the late 1980s, I recognized the Internet's importance in the world and in Asia. I requested a number of early IPv4 Class B assignments (/16s), directly from Jon Postel as NIC function delegation trial, as well as Class Cs and a Class A, for use by research networks in Japan. Since then, I have been administrating the Class A assignment, 43/8, to assist in the long-term development of the Internet in the Asia Pacific region.

In the early 1990s, I helped to establish the *Asia Pacific Network Information Centre* (APNIC) from the *Japan Network Information Center* (JPNIC), to provide continuing allocations of IPv4 address space for our region, Asia Pacific, at a time when the Internet was growing very quickly. APNIC launched in 1993 and has been very successful in managing IPv4 address space since then.

Since 1992, I continued to lead the WIDE project, which was then dedicated to the development and promotion of IPv6. Some of the 43/8 address space was used for this purpose, to assist Japanese networks with renumbering in their transition to IPv6. Some of this space, a /11 in total, was allocated by APNIC to participants in that project, and the rest retained by the WIDE project for other R&D activities.

The deployment of IPv6 has been slower than expected, but I'm very happy that finally, IPv6 is in full production around the Internet, and used by around 25% of Internet users globally. It's clear now that IPv6 will succeed and that the Internet will be greatly improved as the transition continues into the future.

IPv4 has a continuing role on the Internet, but a relatively short-term role, as IPv6 adoption increases. Therefore, IPv4 address space has a current value, but a value that will reduce and disappear over the next 10 years or so. While I have not been an active supporter of the commercialization of IP addresses, the fact is that a market for IPv4 addresses exists and the APNIC community has remained neutral by developing a proper policy framework for market transfers.

In considering the future of 43/8 I have again considered how it may be best used for its original purpose. After careful consideration, I have taken a decision to release this address block, for the purpose toward healthy development of today's Internet services and toward supporting Internet development in the AP region. This is possible by making it available on the IPv4 address market. This is an opportunity to produce a capital asset, with a significant impact on Internet development, if used well and carefully. It is an opportunity that exists today and might not be repeated at any point in the Internet's future.

As I mentioned, APNIC has now been established for 27 years, and it has performed a critical and successful role. APNIC has served very well as the *Regional Internet Registry* for our region, and it has had a great impact in the development of the Internet in our region. With the establishment of the *APNIC Foundation* in 2016, it's clear that APNIC is committed to the continuation and expansion of that good work.

Recognising APNIC's role and its successes, I have asked APNIC to receive a transfer of the unallocated portion of 43/8, on two conditions: that the block will be placed on the IPv4 address market for those who still need IPv4 addresses, and that the proceeds be used in support of Internet development in our region. I am grateful that the APNIC Executive Council has accepted this offer and is now proceeding accordingly, with the establishment of a charitable trust the *Asia Pacific Internet Development Trust*, (APIDT) to take responsibility for this asset and its disposal on the IPv4 address market.

I will remain closely involved, personally and through the WIDE Project, in the management of the Trust, and in its support of Internet development in our region, primarily through the APNIC Foundation. I am very happy to have taken this step and am looking forward to the results in the coming years and decades. I thank everyone involved in this process.

—Jun Murai, Founder, WIDE Project, March 25, 2020

For further information, contact secretariat@wide.ad.jp

WIDE Project: <http://www.wide.ad.jp/>

APNIC: <https://www.apnic.net/>

APIDT: <http://www.apidt.org/>

APNIC Foundation: <https://apnic.foundation/>

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Fabrizio Accatino	Olivier Cahagne	Y Ertur	John Hardin	Dae Young Kim
Michael Achola	Antoine Camerlo	ERNW GmbH	David Harper	William W. H.
Martin Adkins	Tracy Camp	ESdatCo	Edward Hauser	Kimandu
Christopher Affleck	Ignacio Soto Campos	Steve Esquivel	David Hauweele	John King
Scott Aitken	Fabio Caneparo	Jay Etchings	Marilyn Hay	Russell Kirk
Jacobus Akkerhuis	Roberto Canonico	Mikhail Evstiounin	Headcrafts SRLS	Gary Klesk
Antonio Cuñat Alario	David Cardwell	Bill Fenner	Hidde van der Heide	Anthony Klopp
Nicola Altan	John Cavanaugh	Paul Ferguson	Johan Helsingius	Henry Kluge
Matteo D'Ambrosio	Lj Cemerar	Ricardo Ferreira	Robert Hinden	Michael Kluk
Jens Andersson	Dave Chapman	Kent Fichtner	Asbjørn Højmark	Andrew Koch
Danish Ansari	Stefanos Charchalakakis	Michael Fiumano	Damien Holloway	Ia Kochiashvili
Finn Arildsen	Greg Chisholm	The Flirble Organisation	Alain Van Hoof	Carsten Koempe
Tim Armstrong	David Chosrova	Gary Ford	Edward Hotard	Richard Koene
Richard Artes	Marcin Cieslak	Jean-Pierre Forcioli	Bill Huber	Alexander Kogan
Michael Aschwanden	Guido Coenders	Susan Forney	Hagen Hultzs	Antonin Kral
David Atkins	Brad Clark	Christopher Forsyth	Kevin Iddles	Robert Krejčí
Jac Backus	Narelle Clark	Andrew Fox	Mika Ilvesmaki	Mathias Körber
Jaime Badua	Joseph Connolly	Craig Fox	Karsten Iwen	John Kristoff
Bent Bagger	Steve Corbató	Fausto Franceschini	David Jaffe	Terje Krogdahl
Eric Baker	Brian Courtney	Valerie Fronczak	Ashford Jaggernaut	Bobby Krupczak
Santosh Balagopalan	Dave Crocker	Tomislav Futivic	Martijn Jansen	Murray Kucherawy
Michael Bazarewsky	Kevin Croes	Edward Gallagher	Jozef Janitor	Warren Kumari
David Belson	John Curran	Andrew Gallo	John Jarvis	George Kuo
Hidde Beumer	André Danthine	Chris Gamboni	Dennis Jennings	Dirk Kurfuerst
Pier Paolo Biagi	Morgan Davis	Xosé Bravo Garcia	Edward Jennings	Darrell Lack
Tyson Blanchard	Jeff Day	Oswaldo Gazzaniga	Aart Jochem	Andrew Lamb
John Bigrow	Julien Dhallenne	Kevin Gee	Brian Johnson	Richard Lamb
Orvar Ari Bjarnason	Freek Dijkstra	Greg Giessow	Curtis Johnson	Yan Landriault
Axel Boeger	Geert Van Dijk	John Gilbert	Richard Johnson	Edwin Lang
Keith Bogart	David Dillow	Serge Van Ginderachter	Jim Johnston	Sig Lange
Mirko Bonadei	Richard Dodsworth	Greg Goddard	Jonatan Jonasson	Markus Langenmair
Roberto Bonalumi	Ernesto Doelling	Tiago Goncalves	Daniel Jones	Fred Langham
Julie Bottorff	Michael Dolan	Ron Goodheart	Gary Jones	Tracy LaQuey Parker
Photography	Eugene Doroniuk	Octavio Alfageme	Jerry Jones	Rick van Leeuwen
Gerry Boudreaux	Karlheinz Dölger	Gorostiaga	Anders Marius	Simon Leinen
L de Braal	Joshua Dreier	Barry Greene	Jørgensen	Robert Lewis
Kevin Breit	Lutz Drink	Jeffrey Greene	Amar Joshi	Christian Liberale
Thomas Bridge	Dmitriy Dudko	Richard Gregor	David Jump	Martin Lillepui
Ilia Bromberg	Andrew Dul	Martijn Groenleer	Merike Kao	Roger Lindholm
Václav Brožík	Joan Marc Riera	Geert Jan de Groot	Andrew Kaiser	Link Light Networks
Christophe Brun	Duocastella	Christopher Guemez	Christos Karayiannis	Sergio Loreti
Gareth Bryan	Pedro Duque	Gulf Coast Shots	David Kekar	Eric Louie
Stefan Buckmann	Holger Durer	Sheryll de Guzman	Stuart Kendrick	Guillermo a Loyola
Caner Budakoglu	Mark Eanes	Rex Hale	Robert Kent	Hannes Lubich
Darrell Budic	Peter Robert Egli	Jason Hall	Jithin Kesavan	Dan Lynch
Scott Burleigh	George Ehlers	James Hamilton	Jubal Kessler	Sanya Madan
Chad Burnham	Peter Eisses	Stephen Hanna	Shan Ali Khan	Miroslav Madić
Jon Harald Bøvre	Torbjörn Eklöv	Martin Hannigan	Nabeel Khatri	Alexis Madriz

Carl Malamud	Michel Nakhla	William Rawlings	Jeffrey Sicuranza	Martin Urwaleck
Jonathan Maldonado	Mazdak Rajabi Nasab	Bill Reid	Thorsten Sideboard	Betsy Vanderpool
Michael Malik	Krishna Natarajan	Petr Rejhon	Greipur Sigurdsson	Surendran
Tarmo Mamers	Naveen Nathan	Robert Remenyi	Andrew Simmons	Vangadasalam
Yogesh Mangar	Darryl Newman	Rodrigo Ribeiro	Pradeep Singh	Ramnath Vasudha
Bill Manning	Thomas Nikolajsen	Glenn Ricart	Henry Sinnreich	Philip Venables
Harold March	Paul Nikolich	Justin Richards	Geoff Sisson	Buddy Venne
Vincent Marchand	Travis Northrup	Mark Risinger	Helge Skrivervik	Alejandro Vennera
Gabriel Marroquin	Marijana Novakovic	Fernando Robayo	Darren Sleeth	Luca Ventura
David Martin	David Oates	Gregory Robinson	Richard Smit	Tom Vest
Jim Martin	Ovidiu Obersterescu	Ron Rockrohr	Bob Smith	Dario Vitali
Ruben Tripiana Martin	Tim O'Brien	Carlos Rodrigues	Courtney Smith	Michael L Wahrman
Timothy Martin	Mike O'Connor	Magnus Romedahl	Eric Smith	Laurence Walker
Carles Mateu	Mike O'Dell	Lex Van Roon	Mark Smith	Randy Watts
Juan Jose Marin	John O'Neill	Alessandra Rosi	Craig Snell	Andrew Webster
Martinez	Jim Oplotnik	David Ross	Job Snijders	Tim Weil
Ioan Maxim	Packet Consulting	William Ross	Ronald Solano	Jd Wegner
David Mazel	Limited	Boudhayan	Asit Som	Westmoreland
Miles McCredie	Carlos Astor Araujo	Roychowdhury	Ignacio Soto Campos	Engineering Inc.
Brian McCullough	Palmeira	Carlos Rubio	Evandro Sousa	Rick Wesson
Joe McEachern	Alexis Panagopoulos	Timo Ruiters	Peter Spekrijse	Peter Whimp
Alexander McKenzie	Gaurav Panwar	RustedMusic	Thayumanavan	Russ White
Jay McMaster	Manuel Uruena Pascual	Babak Saberi	Sridhar	Jurrien Wijlhuizen
Mark Mc Nicholas	Ricardo Patara	George Sadowsky	Paul Stancik	Derick Winkworth
Carsten Melberg	Dipesh Patel	Scott Sandefur	Ralf Stempfer	Pindar Wong
Kevin Menezes	Alex Parkinson	Sachin Sapkal	Matthew Stenberg	Phillip Yialeloglou
Bart Jan Menkveld	Craig Partridge	Arturas Satkovskis	Adrian Stevens	Janko Zavernik
Sean Mentzer	Dan Paynter	PS Saunders	Clinton Stevens	Muhammad Ziad
William Mills	Leif Eric Pedersen	Richard Savoy	John Streck	Ziayuddin
David Millsom	Rui Sao Pedro	John Sayer	Martin Streule	Jose Zumalave
Desiree Miloshevic	Juan Pena	Phil Scarr	David Strom	Romeo Zwart
Joost van der Minnen	Chris Perkins	Elizabeth Scheid	Viktor Sudakov	Bernd Zeimetz
Thomas Mino	Michael Petry	Jeroen Van Ingen	Edward-W. Suor	廖明沂.
Rob Minshall	Alexander Peuchert	Schenau	Vincent Surillo	
Wijnand Modderman	David Phelan	Carsten Scherb	T2Group	
Mohammad Moghaddas	Derrell Piper	Ernest Schirmer	Roman Tarasov	
Roberto Montoya	Rob Pirnie	Philip Schneck	David Theese	
Charles Monson	Marc Vives Piza	Dan Schrenk	Douglas Thompson	
Andrea Montefusco	Jorge Ivan Pincay Ponce	Richard Schultz	Lorin J Thompson	
Fernando Montenegro	Victoria Poncini	Timothy Schwab	Joseph Toste	
Joel Moore	Blahoslav Popela	Roger Schwartz	Rey Tucker	
John More	Eduard Llull Pou	SeenThere	Sandro Tumini	
Maurizio Moroni	Tim Pozar	Scott Seifel	Angelo Turetta	
Brian Mort	David Raistrick	Yury Shefer	Phil Tweedie	
Soenke Mumm	Priyan R Rajeevan	Yaron Sheffer	Steve Ulrich	
Tariq Mustafa	Balaji Rajendran	Doron Shikmoni	Unitek Engineering AG	
Stuart Nadin	Paul Rathbone	Tj Shumway	John Urbanek	



Follow us on Twitter and Facebook

@protocoljournal



<https://www.facebook.com/newipj>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Diamond Sponsors



Ruby Sponsors



Sapphire Sponsors

Your logo here!

Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

David Conrad, Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Chief Internet Technology Officer
The Internet Society

Dr. Jun Murai, Founder, WIDE Project, Dean and Professor
Faculty of Environmental and Information Studies,
Keio University, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

March 2021

Volume 24, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

FROM THE EDITOR

In This Issue

From the Editor	1
DNS Trends	2
What Have We Done?	18
Fragments	22
Thank You!	28
Call for Papers	30
Supporters and Sponsors	31

We have just completed the annual *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT). The event was to have been held in Manila in the Philippines, but because of the global pandemic it was held as a “virtual” or online event instead. This change of venue is of course not unique to APRICOT. The year 2020 saw many events cancelled, postponed, or converted to online gatherings. In most cases, the Internet remained a reliable and resilient alternative as more and more organizations and individuals took advantage of various networked conferencing systems. Already several studies have documented how the Internet performed through the pandemic. For example, RIPE Labs published “The Lockdown Effect—Implications of the COVID-19 Pandemic on Internet Traffic,” which you can find through your favorite search engine.

Many of the core protocols of the Internet have been updated or otherwise enhanced over the years, particularly protocols that originally did not have security as part of their initial design. Development continues within *The Internet Engineering Task Force* (IETF) to improve all aspects of the Internet Protocol Suite, including novel uses of existing technologies. Our first article is a look at current developments in the *Domain Name System* (DNS).

The first IETF meeting was held in January 1986 with a mere 21 attendees. Thirty-five years later the typical IETF meeting attracts about 1,000 attendees from all over the world and lasts a full week, including the pre-IETF *Hackathon* and *Code Sprint* sessions. During the pandemic, IETF meetings too have been confined to online events, this month in “Virtual Prague.”

We don’t usually publish opinion pieces in this journal, but with 35 years of IETF development and more than 50 years since the origins of the Internet, this seems like a good time to pause and examine where we are with respect to the overall state of our digital economy. Geoff Huston asks, “What have we done?” in his provocative essay that we hope will inspire you to submit your own views in the form of a Letter to the Editor or perhaps an opinion column of your own.

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

DNS Trends

by Geoff Huston, APNIC

We used to think of computer networks as being constructed using two fundamental common infrastructure components: *names* and *addresses*. Every connected device had a stable protocol address to allow all other devices to initiate a communication transaction with it by addressing a data packet to this protocol address. And every device was also associated with a name, allowing human users and human-use applications to use a more convenient alias for these protocol addresses. By mapping names to protocol addresses, the realm of human use could navigate the services of the network by using symbolic names, while at the level of packets the data flow was directed by the network based on topology information of where these device addresses were located.

But that's 1980s thinking and 1980s network architectures.

Communications architectures have evolved, and today's Internet architecture has, very surprisingly, dispensed with that view of the role of addresses. These days, in no small part because of the exhaustion of the IPv4 address pool, but equally because of an architectural evolution that had to cope with the massive explosion of numbers of devices in networks, we've shifted to a client/server network model where clients initiate connections and servers respond. So now clients don't require a permanently assigned network-wide address. Instead, they can use an address only while it is communicating with a server and pass it back to a shared address pool otherwise. Equally, on the server side we've seen the aggregation of uniquely named service points into service delivery platforms, and the multiplexing function that directs clients into the appropriate service rendezvous point is performed at an application level rather than as an infrastructure function. We're now using the address infrastructure in very different ways than the way we had envisaged in the 1980s. Addresses in today's terms look more like ephemeral session tokens on the client side, and coarse rendezvous points on the server side. It is left to the application level to define the specific client-requested service.

But the architecture of the name space and its use has not been static either. The name infrastructure of the Internet is subject to the same evolutionary pressures, and it is these pressures I'd like to look at here. How is the *Domain Name System* (DNS) responding? This survey has three parts: trust, privacy, and all the other stuff.

Trust

Can you believe what the DNS tells you? The answer is that you probably can't!

Many parties have exploited this obvious failure in the trust model in many ways. The DNS is seen as an overt control channel.

For example, you can block access to a named service if that name does not resolve in the DNS. As a consequence, we've seen the rise of deliberate lies in the DNS where content and services that are categorised as harmful are removed from access by withholding the associated name resolution in the DNS. Numerous open resolvers have turned this filtering of the DNS into a positive attribute, and there are many so-called "clean feed" resolvers that do not resolve a collection of service names where the service is deemed to be harmful or criminal in some manner.^{[1] [2]}

This selective filtering of the DNS is a distinguishing feature in the realm of competitive open resolvers. We've also seen numerous national regimes placing the onus on ISPs to block certain services, and given that addresses are no longer uniquely associated with individual services, the implementation of these national regulations is invariably performed through DNS blocking.^[3]

We have also seen exercises to attempt to monetise the DNS, where "no such domain" (**NXDOMAIN**) DNS responses are rewritten to send you to a sponsoring search site through response rewriting.

DNS lies have also been used in the IPv6 transition environment where the DNS records—the protocol addresses—are synthesised to allow you to be steered through an IPv4-IPv6 transitional environment.

The motives of all these exercises may vary, but the result is the same, in so far as the DNS answer is a lie.

Then there are the hostile efforts to replace a genuine response with a lie in order to mislead you, in addition to the technique of response guessing to try to insert a fake response before the "normal" response. You can use this technique in DNS over the *User Datagram Protocol* (UDP) transport as the first UDP response whose query section matches the original query the asker used—whether or not it is the "genuine" response. We have also seen manipulated glue records and even attacks on fragmented packets.^[4] The insidious nature of these forms of attack is that they rely on the host system to run quite normally. It's the infrastructure of the name system itself that is being perverted here, and the applications are completely unaware of this manipulation.

The response to this need to detect any form of manipulation of the DNS response that has taken place, and, even better, to withhold the lie from the user, is to add a trust mechanism to the DNS. This trust mechanism takes the form of adding digital signatures to DNS responses. The idea is that a digital signature attached to a DNS response can allow the receiver of the response to be assured that the DNS information is current, that it is authentic, that it has not been manipulated or altered, and that it cannot be repudiated. *Domain Name System Security Extensions* (DNSSEC), the framework for adding digital signatures into the DNS, was some 10 years in the making.

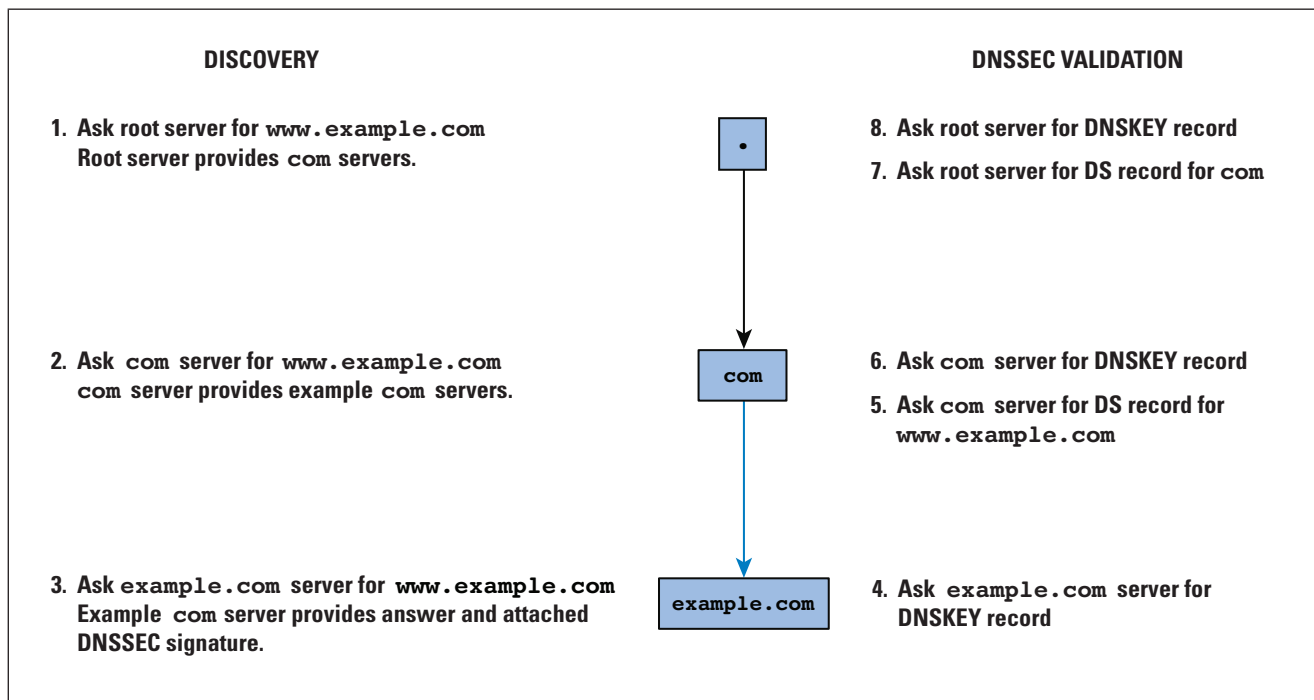
You can associate a key pair (or pairs) with a delegated zone in the DNS, and a set of five further DNS *Resource Records* (RRs) are defined in this framework to sign each entry and to aid in validation of the signature.^{[5][6][7]}

The commentary on DNSSEC deployment varies considerably. Some 25% of the world's users cannot reach a DNS-named service if they cannot validate the DNSSEC signature. That's triple the level from early 2014, so the adoption of validation is certainly gaining some momentum.^[8] At the same time, the number of DNSSEC-signed zones appears to be pitifully low. Of the hundreds of millions (perhaps billions these days) of delegated zones in the DNS, we see some 8M signed zones in one such survey.^[9] Perhaps a more relevant metric is the ranking of domain names by usage and the ratio of DNSSEC-signed zones in that set. The Alexa top 25 list is a good place to start. None of these is a DNSSEC-signed name.^[10] A scan of all **.com**, **.org**, and **.net** second-level domains found that between 0.75 and 1.0% of all domains in these three zones are signed.^[11] Zone signing is just not that common in today's DNS.^[12] It appears that turning on DNSSEC validation of DNS responses in a recursive resolver has very little downside, given that so few of the popular DNS names appear to be DNSSEC-signed in the first place!

Why is DNSSEC zone signing so uncommon? Are the content and service providers behind these DNS names unconcerned about potential misdirection of users? Of course not! They are extremely concerned because ultimately their potential revenue and reputation are at risk. Are they ignorant about DNSSEC? Again, not at all! Zone signing or not is a choice. These providers want to prevent users from being misdirected, but equally they want to reduce the dependence on intermediaries, and they want your service experience to be as efficient as possible. If DNSSEC was the only choice, then content and service providers would be using it. But it's not the only option. These days service provision uses *Transport Layer Security* (TLS). Almost every service URL out there is an HTTPS URL. Can you be misdirected with TLS? Not normally. Misdirection or deception requires leakage of the service provider's private key, or corruption of the Web *Public Key Infrastructure* (PKI) certificate system. TLS is also fast, because all the credentials needed to validate the certificate are provided in the TLS handshake.

Is DNSSEC validation as fast? Well, no. DNSSEC validation is a serial query sequence all the way back up the name delegation path (Figure 1). Is deployment of DNSSEC zone signing simple? No. There is local key management, the *Zone Signing Key* (ZSK)/*Key Signing Key* (KSK) key split, limited automation, and limited support for high-resilience hosting. And the fundamental criticism is that all this additional effort doesn't stop recursive resolvers from passing back lies in the DNS for DNSSEC-signed zones anyway, because most stub resolvers do not perform DNSSEC validation in any case.

Figure 1: DNSSEC Validation



The entire point is that lies in the DNS were just not possible with DNSSEC. But that assumes that all resolvers perform DNSSEC validation, and that's not the way we've deployed it so far. Many recursive resolvers perform DNSSEC validation. Very few stub resolvers perform DNSSEC validation. This scenario generally works in so far as the recursive resolver withholds the response if the DNSSEC validation fails. But what if the recursive resolver is the one that is telling the lie in the first place? The stub resolver is none the wiser, given that it's not validating, so the lie stands. If the ISP's recursive resolver is blocking some names, performing **NXDOMAIN** substitution, or redirecting actual names, then the stub resolver is just caught in the lie. With all that effort to sign the zone, and all that effort to validate the DNS response, there is absolutely no robust protection against being misdirected. TLS just seems to offer a solution that is faster, simpler, and more robust. No wonder few zone administrators use DNSSEC signing in the DNS service world. It's just a case of more pain, and no real gain.

Is DNSSEC good for anything else? As long as 75% of users sit behind nonvalidating DNS resolver systems—and virtually no users directly validate DNS responses in any case—we cannot place critical information in the DNS securely and expect everyone to be protected by DNSSEC. This means that the incentives for putting critical information into the DNS and protecting it with DNSSEC do not look very convincing. There is simply no natural market-based incentive for deployment of DNSSEC. This conclusion is distressing, because it would certainly be more useful for the network and its captive user population if its name system were trustworthy.

Many have said that the heart of numerous issues with the DNS lies in the choice of a transport protocol for the DNS. The use of UDP as the primary first-choice protocol and the fallback to TCP means that it's challenging to place large quantities of information in DNS answers while still operating within what we've become accustomed to in terms of parameters of speed and robustness.

Validation is a very inefficient process, and the inefficiency is increased by the DNS model where the onus is placed on the client, who is requesting the information, and not the server, who is the source of this information. End clients do not validate because every validation operation would entail further DNS queries in order to construct the validation chain, and the incremental time penalties would be unacceptable in terms of user expectation.

Frustratingly, we know how to make DNSSEC validation faster, and the approach is to pre-provision the validation answers. We can package up all the answers to the DNSSEC validation chain construction queries and include them as additional information to the original signed answer in a single chain extension in the response.^[13] However, it's unlikely that this inclusion is viable in a DNS-over-UDP framework. If we want to go down a TLS-like path and package up a validation chain into the DNSSEC-signed response, we will probably have to use DNS over TCP or *DNS over TLS* (DoT).^[14] The price of this trust solution is significant, and it creates a higher threshold for the benefits that trusted answers in the DNS can provide. If all this discussion is about protecting users from a Kaminsky-styled attack,^[15] then that's just not enough of a case. The benefit needs to be far more than helping justify the considerably higher costs in moving the DNS from UDP to a TCP-styled platform.

Privacy

Everybody looks at the DNS. Everybody. Because the Internet is funded by its users, then what users do on the Internet is of paramount interest to people who sell services to users. Because a lot of crime these days is cybercrime, the criminal and abusive behaviour on the Internet is of fundamental interest to those agencies whose role is to police such behaviours. Because the Internet is now largely about how individuals choose to live their lives and how and why they communicate with others, we've learned that what users do is of paramount interest to government.

How can you find out what users do? Easy. Look at the DNS. Every transaction on the Internet starts with a DNS query, and the DNS exposes every action. But it's worse than that. The DNS is needlessly and senselessly chatty. The DNS overexposes information. These queries and responses are collected, packaged, analysed, profiled, replayed, and traded at all points in the DNS.

How can we make the DNS not the go-to system to expose users and user behaviours to business and government alike? How can we improve its privacy?

There was little in the way of motivation to do anything about this question for years. After all, if the Internet actors are busy constructing a global economy based on surveillance capitalism, why should the parties conducting this surveillance make the task any harder than necessary? The watershed moment that changed the stance for many was the publication of material that Edward Snowden gathered. Government agencies had spent considerable sums in weaponizing the Internet and transforming it into a highly effective surveillance tool that operated at a scale of national populations. Their motivations were not overly concerned about your future purchases, but more about your personal profile. And of all the components of the Internet, the system that laid out all this information in a clear text prepackaged format was the DNS.

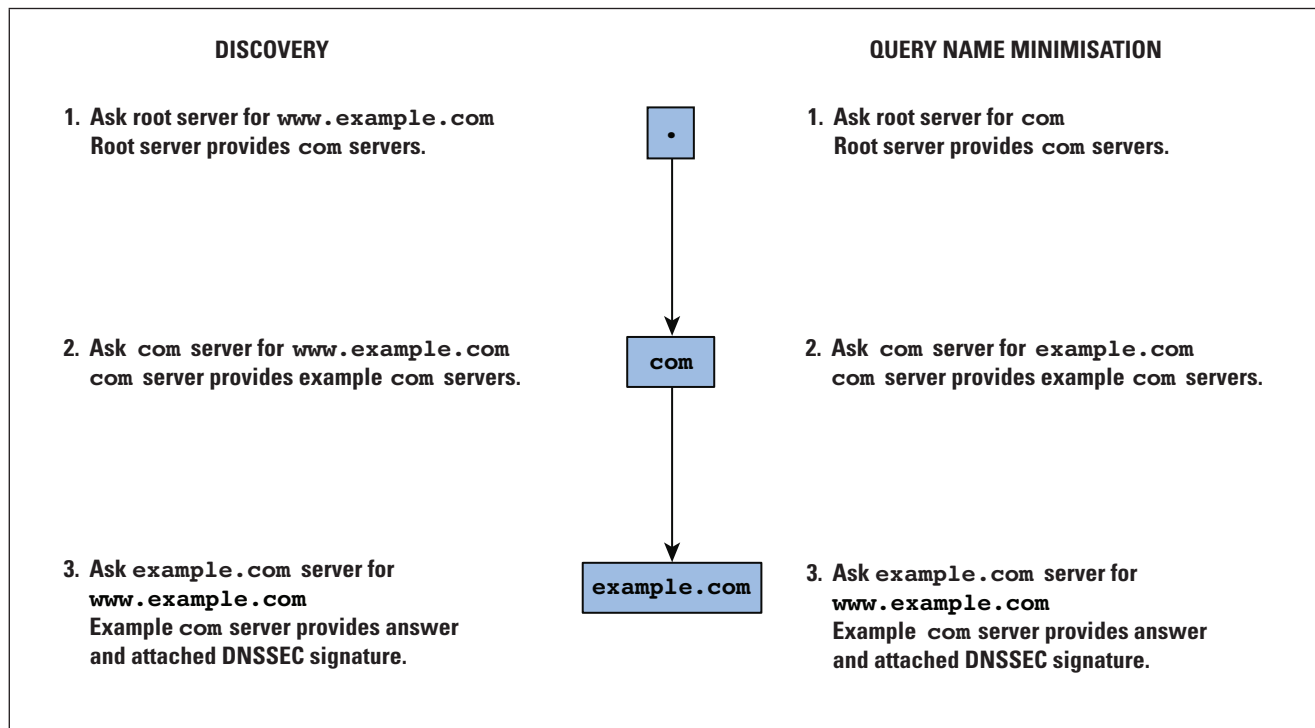
In response, we've been changing aspects of DNS behaviour to try to stop the most blatant forms of information leakage.

The first of this set of privacy-enhancing responses is called *Query Name Minimisation*.^[16] The change is to prevent the DNS name-resolution process from being an unconstrained extraneous information leak. This leakage largely relates to the interaction between recursive resolvers and the authoritative name servers. The task of the recursive resolver is to find the right name server to ask, and it starts at the root and asks the query. The response of the root-zone server will direct the queries to the name servers of the relevant delegated top-level domain name, and this process repeats as the resolver traverses down the delegation hierarchy until the resolver has an answer.

But in this process every name server in this sequence, from a root server down, is now aware of the full DNS name that is being resolved. Query Name Minimisation trims the name in these queries so that only the next label is exposed to each name server. Root servers will see only top-level domain name queries, while top-level domain name servers will see only second-level name queries, and so on (Figure 2).

There has been some further work to understand the most robust query type for this discovery process. The initial suggestion of **NS** queries has been supplanted by **A** queries in the light of experience with this approach. The issue of **CNAME** rewriting and the equally vexed question of Empty Non-Terminal domains and the variable behaviour of name servers in such situations have added some complications to this question. This technique of this approach is now widely used, although some implementations have taken some license with the specification and used their own re-interpretation of the technique. Some resolvers, apparently including Google's public resolver service, performs Query Name Minimisation to only the first three levels of the DNS name.

Figure 2: Query Name Minimisation



It appears as if the recursive resolver is deliberately withholding full query name information from the root servers and the top- and second-level domain name services, but is quite willing to disclose the full query name information to servers for zones that are deeper in the name hierarchy. Is this approach motivated by protecting user interests or by an effort to deny information to authoritative servers located at the upper levels of the name hierarchy?

Of course, if we were serious about user privacy, the *Client Subnet* extension would never have been specified.^[17] The knowledge of full query names that are emitted by a recursive resolver is to some extent mitigated by the inability to conclusively associate such queries with an end user. But if the query is also loaded with the IP address of the end client, or even the network subnet of the end client, then all pretence of privacy protection has been shredded. While Query Name Minimisation could be seen as a positive step in providing a greater level of concealing extraneous information in the DNS, the use of the Client Subnet value in queries is a gigantic leap backward!

A generic response to privacy considerations on the Internet has been channel encryption. *Telnet* was replaced by *ssh* because of the issues of running sessions over the Internet in the clear. Similarly, **HTTP** has been largely replaced by **HTTPS** for much the same reason. The DNS is increasingly an anachronism in still passing queries and responses in the clear. Not only does it permit eavesdropping, but it also enables efforts to manipulate the responses, all to the detriment of the user.

However, to repeat an earlier observation, the heart of many issues with the DNS lies in its choice of transport protocol. Encryption normally involves many steps, including the presentation and validation of credentials to confirm that clients are talking to the party they intended to talk to, and also to establish a session encryption key to allow encryption of the subsequent data exchange in a manner known to the two parties but unknown to all others. This type of encryption is challenging in UDP. The effort to implement TLS over UDP, namely *Datagram TLS* (DTLS)^[18], has the overhead of the exchange of credentials and session cipher establishment, so it's a long step away from a single packet exchange of query and response. DTLS also should avoid IP-level fragmentation, but it cannot avoid large payloads associated with this session establishment process. The result is that fragmentation is pushed up to the application layer and DTLS needs to handle payloads that extend across multiple DTLS datagrams. It appears that the additional overheads of DTLS roughly equate to the overheads of TLS over TCP, but with some added fragility relating to packet fragmentation that is not replicated in TCP. The result of this fragility of DTLS means that when we refer to DNS over TLS, we are in fact referring to *DNS over TLS over TCP* (DoT).^[14] It is this TCP-based implementation of TLS that has been implemented and deployed over the path between the stub resolver and the recursive resolver.

DoT adds encryption to the stub-to-resolver path; not only does encryption hide the query and response stream from eavesdroppers, but also DoT prevents alteration or manipulation of the response by third parties. The recursive resolver can still lie about the response, and unless the stub resolver is performing DNSSEC validation (and it's likely not) and the domain name is signed (which it most likely is not), then any DNS lie from the recursive resolver will be unnoticed, whether or not the transport channel from the recursive resolver to the stub resolver uses TLS. A lie is still a lie no matter how secure the packaging used to carry it is. DoT does not eliminate the potential for manipulation of DNS information, but limits the number of entities who are in a position to perform such manipulation and the place and method that the manipulation can be performed. It could be argued that with DoT all you really gain is being better informed as to who is lying to you!

How far should channel cloaking go? Should the identity of the other party be obscured? Should the fact that these transactions are DNS exchanges be obscured? DoT makes no effort to cloak its use. The use of TCP port 853 for DoT is a visible signal that there is an active DoT connection. The use of a novel port number is likely to cause many firewall configurations to trigger their drop filters. The IP address of the remote end is clearly visible, as is the TCP header. The TLS handshake may get around to using *Encrypted Client Hello* (ECH)^[19] and encrypt the server name at some point in the future, but in the case of DoT it probably is a minor artefact, given that name-based overloading of service IP addresses is not happening in DoT today and unlikely will in the future.

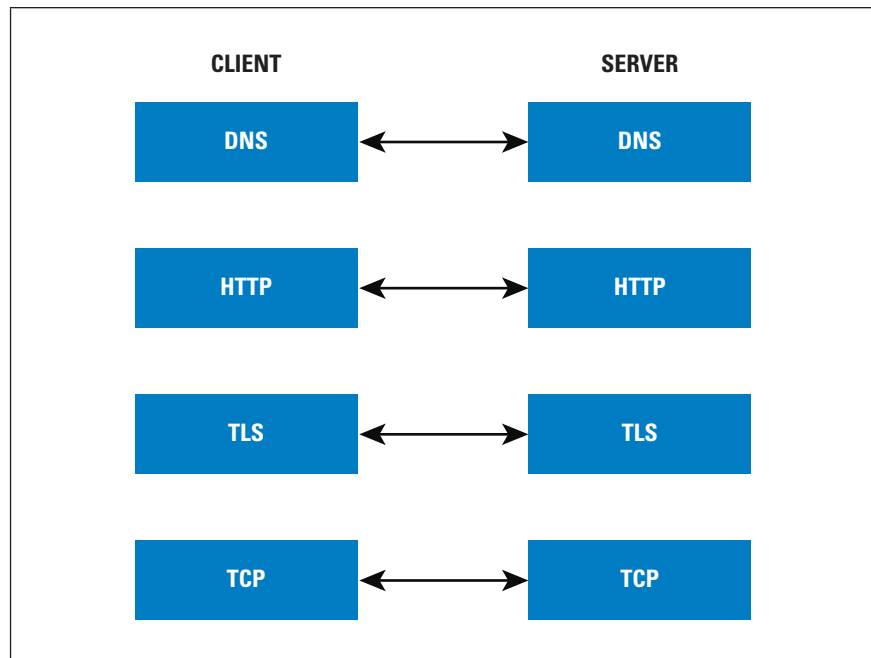
Is DoT going anywhere? It is unlikely in my view. Right now, it requires users to play with their DNS settings, and that is a massive barrier to widespread use. Some users may use it as a means to jump across one set of recursive resolver's DNS filters, such as those provided by their ISP, to hook up with another DNS resolver provider, but with the overt signalling that this is happening, an ISP can readily block this action if it wants to. In theory, the use of TCP permits larger DNS payloads, and we could possibly use *DNS Chaining*^[13] to make DNSSEC validation fast and efficient on end systems using DoT. But so many other preconditions, including server provisioning of DNS Chained responses and a reliable way for the DNS to manage large responses, mean that it is still a distant glimmer of a possibility and nothing more.

DoT is seen as a replacement for the existing DNS infrastructure service, where the DNS is a service located on the common platform and applications use the same DNS resolution calls to the platform as they always have. It's a platform approach to securing the DNS. Adoption is probably going to require some form of automated provisioning that typically involves the local access service provider.

Given that the major compromise threat actor here is the same access ISP, and given that the ISP operates the recursive resolver in any case, it's very challenging to understand the incremental benefit of DoT deployment to an ISP. Perhaps it may be that its benefit is as a barrier to other hosts in the local network. Local residential and enterprise environments are cluttered with IP stacks from many providers. A compromised stack is inside the external firewalls and is trusted merely by its physical location. DoT shifts the conversation of a DNS host into a protected channel where the protection is against other hosts on your local network!

DNS over HTTPS (DoH)^[20, 27] uses the same TCP and TLS foundations, but adds an HTTP context to the transactions. (Figure 3) A couple of changes here are interesting. The first is the switch to TCP port 443. It looks like any other HTTPS traffic and is not so readily identified in the network as being DNS traffic. Second, the DoH servers do not need to use dedicated IP addresses. Like the web itself, the HTTP protocol allows for named service points. And with TLS 1.3, with ECH you can conceal even the server name in an encrypted envelope. But there is a little more. HTTPS is an application-level protocol, and this approach allows an application to bypass the DNS services the platform provides. Therefore, no ISP-based platform-level configuration is necessarily relevant, and the application can not only conceal its DNS transactions from the local and remote networks, it can also hide these same transactions from the platform and other applications running on the same platform.

Figure 3: DNS over HTTPS (DoH)



If this development heads to DNS over HTTPS/3^[21], which uses *Quick UDP Internet Connection* (QUIC)^[22, 26], then numerous capabilities are unlocked. Not only is the transport control protocol cloaked behind an encryption envelope in QUIC, but you can make many DNS requests on a single transport channel simultaneously.

How far can we go with this effort to advance a privacy agenda in the DNS? Once we've deployed Query Name Minimisation, discarded Client Subnet, adopted DoH using HTTPS/3 as the application-to-recursive resolver protocol, and pushed DNSSEC validation to the application via attached chained DNSSEC responses, then you have realized much of the achievable trust and privacy agenda. At that point, much of the ability for a third-party onlooker to associate an end-entity identity with a DNS request is severely curtailed, and while a recursive resolver is still privy to these user transactions, the use of DNSSEC all the way to the edge makes response manipulation by any external party, even the recursive resolver itself, particularly challenging when the original DNS data is DNSSEC-signed.

The DNS privacy effort is moving on. The current question is: Should we encrypt the paths between a recursive resolver and authoritative servers? Assuming that Client Subnet has been abolished, there is little that such transactions directly reveal about the identity of the end user, and the larger the pool of clients that a recursive resolver serves, the larger the crowd each individual's queries can hide in. Irrespective of the questions of whether it is feasible (it is) and whether it is scalable (no clear answer, but it looks to have an appreciable incremental cost), the fundamental question of whether channel privacy makes any sense in a privacy context for the individual end user remains.

Other Topics

Other aspects of the technology evolution of the DNS are covered in the following sections.

Internationalized Domain Names (IDNs)

The DNS has traditionally used a 7-bit ASCII code for names. Upper and lower case are equivalent, and in addition to the Latin characters, DNS labels can use hyphens and number characters. In certain circumstances the underscore is also permitted. The expansion of the DNS into a larger character repertoire^[23] has not been a stellar success. The design decision was to preserve the capabilities of the DNS system and use encoding to map the larger character set into this restricted alphabet.

The choice of *Unicode*^[28] as the underlying character repertoire for this expanded character set was not a very good choice. Unicode involves a contract between an application and a printer. It does not matter that Unicode has multiple ways to print the same glyph on a printer. The printer does not care. But the DNS cares. The DNS has no concept of “what it means,” and alternate Unicode strings that are presented in an identical way on a screen actually map to distinct DNS names. So, the effort in the use of Unicode in the DNS has been one of trying to push the Unicode glyph set back into the box and try to specify canonical subsets of Unicode that minimise display similarity. This request is tough, and made even harder by the increasing variance in display glyphs used to display the same Unicode code point. This challenge is most evident in emoji characters.

Why is it a problem? Because the Internet still works on a rather crude model of “what you see is what you get.” If alternate ways of coding the same visual outcome are possible, then they are distinct labels in the DNS and can be associated with distinct service points. The possibilities to dupe unsuspecting users is of course a natural and inevitable outcome of this process.

DNS Abuse

These days “DNS Abuse” is a current topic, particularly in the *Internet Corporation for Assigned Names and Numbers* (ICANN) world. The phrase describes an effort to engender a level of self-regulation in the DNS supply industry, where behaviours are governed largely by contractual provisions between the registrant and the registrar, and between the registrar and a common registry. It allows various forms of abusive behaviours, including criminal activities, that use the DNS to be sanctioned by contractual enforcement including takedown of the DNS names. It’s a lot like the self-regulatory measures that are common in the finance industry, but without the reporting framework, without any common legal framework for enforcement, and without any penalties for breaches.

My suspicion is that it will turn out to be no more effective than the similar measures to undertake self-regulation in the finance sector, and probably even more ineffectual than the rather unimpressive results that the banking sector has posted. It's unlikely to be successful in reducing the levels of abusive and criminal behaviour that use the DNS and the Internet.

DNS Fragmentation

DNS fragmentation is also a perennial topic in the evolution of the name space. The pressures for a single, consistent name space are embedded in the concept of a single network. Communications systems rely on assumptions of referential integrity, and referential integrity typically implies that the same DNS name refers to the same resource.^[24]

We've seen this concept tested many times, from alternate root systems of a couple of decades ago to private name spaces today in the enterprise environment. A good case in point about referential fragmentation is the use of search terms as a replacement for the DNS. The objective of a search engine is to try to customise the responses to best match the known preferences of the querier, and when *I* attempt to pass a pointer to you about a digital resource, the search term that I use that will expose this resource may not be exposed when *you* enter the same search term in your context. This possible difference is not only an attribute of search engines, but a feature.

However, DoH enables other forms of DNS fragmentation. It enables you to lift a name space out of common network infrastructure and place it into the context of an attribute of an application. The application can direct DoH queries to server infrastructures of its own choosing and provide responses that pertain to the application as distinct from a lookup in a common distributed database. The ability in HTTPS to push objects to the application client also allows you to use so-called *Resolverless-DNS*^[25], where an application can improve the performance of name resolution functions by performing them in advance of the time they are needed.

Name Flattening

DNS *name flattening* has been a constant pressure in the DNS. Nobody wants to have their critical service names **buried.deep.down.in.the.dns.under.a.bunch.of.other.names**. Not only do such names take longer to resolve, they increase the set of dependencies in the same way because presumably a greater number of service providers all the way down in the name hierarchy exist. DNS users want shorter names. The shorter the better. The result is that the name space is under constant erosive pressure to flatten down. The ultimate place to land is in the top level of the DNS, in the root zone, and as the price premium for top-level domain comes down, the pressure to inflate this zone with significantly larger numbers of entries is an inevitable consequence.

The Future of Names

But perhaps the forces of evolutionary pressure are more fundamental and parallel the evolutionary forces of the Internet itself.

The silicon industry is indeed prodigious, and there are many more processors in this world than people. While we have constructed the DNS name space using an analogue of natural language terms as a means of facilitation of human use, this use pattern is not necessarily the dominant use pattern of the DNS any longer. One view of the DNS today is a universal signalling and tunnelling protocol, and the use of the DNS as a command-and-control channel for malware bot armies testifies to the efficacy of such use of the DNS!

It's likely that as the number of such devices increases, the use of the DNS as an orchestration mechanism increases in importance and the human use of the DNS becomes increasingly marginalised. Human-use DNS may well become an esoteric luxury business. The high-touch activity of DNS name management is unsustainable in a shift from human to largely automated use, and the business models and institutions that populate this space will need to adjust to a names business that provides names not as a branding attribute using natural language tokens, but as an undistinguished commodity activity. In the same way that we have transformed IP addresses from end-point identifiers to ephemeral session tokens, we may well see the DNS as a code base for command and control of highly distributed automated systems, and that is very different from the distributed database lookup that we originally constructed for the human-use model of the DNS.

When we think of a DNS query as a set of instructions to a DNS server, and the DNS server as a distributed processing environment, the DNS changes from a distributed database to a distributed computation and signalling environment. The composition of labels in such a DNS is no longer roughly derived from dictionaries of known words from human languages, but instead is encoded instructions where the labels are in effect a coded program for a name resolver to execute. It is certainly a different future for the DNS as we know it, but its probable commoditisation in the future is in line with the plight of carriage, switching, and content in the Internet!

From this perspective, the evolution of the DNS parallels the larger evolution of the Internet itself, where the infrastructure is not about a human-usable framework any longer, but instead is focussed on providing a highly automated environment where the elements are themselves programs and automata.

That does not mean that the human-use DNS will disappear. But the DNS as we know it today may end up as a small set of high-end luxury boutique activities that make a feature of the luxury of custom procedures to manage persistent names.

In the meantime, the rest of the DNS heads deeper into a commodity utility world of large sets of algorithmically generated transient names that are managed entirely automatically and tailored for one-off use by other processes. It may be that the overwhelming use of tomorrow's DNS has nothing much to do with human names any longer and will be concentrated on serving a largely automated framework that uses the DNS to support a general command-and-control signalling framework. Ephemeral names are as good as, if not better than, persistent names. Registration and attribution processes are largely irrelevant.

The DNS may still be valuable, but individual names will be completely worthless!

References and Further Reading

- [1] Quad 9 Open DNS Resolver
<https://www.quad9.net/about/>
- [2] Cloudflare 1.1.1.1 for Families
<https://blog.cloudflare.com/introducing-1-1-1-1-for-families/>
- [3] Young Xu, "Deconstructing the Great Firewall of China," *Thousand Eyes*, March 8, 2016.
<https://blog.thousandeyes.com/deconstructing-great-firewall-china/>
- [4] Kazunori Fujiwara and Paul Vixie, "Fragmentation Avoidance in DNS," November 2020. Internet-Draft, work in progress.
<https://tools.ietf.org/html/draft-ietf-dnsop-avoid-fragmentation-03>
- [5] Scott Rose, Matt Larson, Dan Massey, Rob Austein, and Roy Arends, "DNS Security Introduction and Requirements," RFC 4033, March 2005.
- [6] Scott Rose, Matt Larson, Dan Massey, Rob Austein, and Roy Arends, "Resource Records for the DNS Security Extensions," RFC 4034, March 2005.
- [7] Scott Rose, Matt Larson, Dan Massey, Rob Austein, and Roy Arends, "Protocol Modifications for the DNS Security Extensions," RFC 4035, March 2005.
- [8] APNIC Labs, DNSSEC Validation Report
<https://stats.labs.apnic.net/dnssec>
- [9] DNSSEC Signed Zone Survey
<https://www.secspider.net>
- [10] DNSSEC Name and Shame
<https://dnssec-name-and-shame.com>

- [11] StatDNS website
<https://www.statdns.com>
- [12] Matthäus Wander, “Measurement survey of server-side DNSSEC adoption,” 2017 Network Traffic Measurement and Analysis Conference (TMA), Dublin, Ireland, 2017, pp. 1–9, DOI: 10.23919/TMA.2017.8002913.
- [13] Paul Wouters, “Chain Query Requests in DNS,” RFC 7901, June 2016.
- [14] John Heidemann, Duane Wessels, Allison Mankin, Paul Hoffman, and Liang Zhu, “Specification for DNS over Transport Layer Security (TLS),” RFC 7858, May 2016.
- [15] Dan Kaminsky, “The Great DNS Vulnerability of 2008,” <https://duo.com/blog/the-great-dns-vulnerability-of-2008-by-dan-kaminsky>
- [16] Stephane Bortzmeyer, “DNS Query Name Minimisation to Improve Privacy,” RFC 7816, March 2016.
- [17] Wilmer van der Gaast, Carlo Contavalli, and Warren Kumari, “Client Subnet in DNS Queries,” RFC 7871, May 2016.
- [18] Eric Rescorla and Nagendra Modadugu, “Datagram Transport Layer Security Version 1.2,” RFC 6347, January 2012.
- [19] Christopher Wood, Kazuho Oku, Eric Rescorla, and Nick Sullivan, “TLS Encrypted Client Hello,” Internet-Draft, work in progress, March 2021.
<https://datatracker.ietf.org/doc/draft-ietf-tlsesni/>
- [20] Paul Hoffman and Patrick McManus, “DNS Queries over HTTPS (DoH),” RFC 8484, October 2018.
- [21] Mike Bishop, Ed., “Hypertext Transport Protocol Version 3 (HTTP/3),” Internet-Draft, work in progress, October 2020.
<https://datatracker.ietf.org/doc/draft-ietf-quic-http/>
- [22] Christian Huitema, Melinda Shore, Allison Mankin, Sara Dickinson, and Jana Iyengar, “Specification of DNS over Dedicated QUIC Connections,” Internet-Draft, work in progress, September 2019.
<https://tools.ietf.org/html/draft-huitema-quic-dnsquic-07>
- [23] John Klensin, “Internationalized Domain Names for Applications (IDNA): Definitions and Document Framework,” RFC 5890, August 2010.

- [24] Internet Architecture Board, “IAB Technical Comment in the Unique DNS Root,” RFC 2826, May 2000.
- [25] Erik Sy, “Enhanced Performance and Privacy via Resolver-less DNS,” August 2019.
https://svs.informatik.uni-hamburg.de/publications/2019/2019-08-13-Sy-preprint-Enhanced_Performance_and_Privacy_via_Resolver-Less_DNS.pdf
- [26] Geoff Huston, “A Quick Look at QUIC,” *The Internet Protocol Journal*, Volume 22, No. 1, March 2019.
- [27] Geoff Huston, “DNS Privacy and the IETF,” *The Internet Protocol Journal*, Volume 22, No. 2, July 2019.
- [28] Unicode: <https://home.unicode.org/>
- [29] George Michaelson, “DoH the right thing,” APNIC Blog, February 10, 2021.
<https://blog.apnic.net/2021/02/10/doh-the-right-thing/>
See also page 24.
- [30] Geoff Huston, “DNS Oblivion,” APNIC Labs, December 15, 2020.
<https://labs.apnic.net/?p=1392>

GEOFF HUSTON, B.Sc., M.Sc. A.M., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

Opinion: What Have We Done?

by Geoff Huston

One of the roles of an opinion piece is to challenge your assumptions and present alternative perspectives, and that is certainly what I plan to do here. You may not agree with my views. I'm not even sure that I agree with them all of the time, because some of these opinions are pretty bleak. But if this article provokes you to make your own assessment of the Internet in a broader context of the evolving relationship between society and this technology, then that is perhaps as much as I could ever hope to achieve here. I should also say that I'm writing this opinion piece as an individual and nothing more. I am not pretending to speak on behalf of my employer in any way. These are *my* words and thoughts.

I was asked to speak at an *Internet Governance Forum* during the COVID-cursed year of 2020. I was briefed that “the most useful thing would be to hear your thoughts on what are the big issues at the moment. Where you see things heading. It would mean that you’d be speaking on your areas of interest from your perspective, not necessarily trying to channel some sort of universal Internet zeitgeist.”

I have found this brief a challenging one. In some decades of working in this space I’ve heard many boom-and-bust talks. In addition, I have seen techno-exuberance reach dizzying heights—and then expositions of sobering realities bring it all back to Earth. But behind this phenomenon I have not seen many perspectives that challenge the very fundamentals of the Internet. We appear to assume that the technology is either beneficial, or at worst neutral, and it’s the humans in the loop that overreact. Perhaps, even more dangerously, we assume that the technology is competently implemented. This assumption is perhaps the most dangerous one!

My personal view is that we are heading to a Bad Place. A *very* Bad Place.

A Revolution

Compared to our somewhat naive expectations about the role of computers and networking in the 1980’s, we have come a long way down a path that now seems to have taken a turn into some dark—and possibly malign—spaces. How else could we have ended up in accusations of rigged elections, “fake” news, and truly bizarre paranoid notions of some form of “deep state” that seems to sit within the collective social psyche these days. But it’s not all just a parade of some ridiculous memes that appear to be rooted in human credulity, because we also have to acknowledge the wholesale destruction of livelihoods and the creation of a new technology economy that is based largely on surveillance capitalism. The digital automation of our society has a highly disruptive aspect, and I think we can confidently assert that we are in the middle of a social revolution as fundamental as the industrial revolution. However, in this case we seem to have backed into this one with our eyes closed.

How could we have missed all the signals? Why are we still thinking that the old social contracts are still valid when they are clearly broken? What went wrong? Well, I'm sure exploring that subject would make a great thesis, but we have two problems. Firstly, I have only a few minutes of your time with this article, and secondly, I really don't know why it all went so wrong anyway! So, without truly knowing how it happened, we find ourselves trapped in another massive revolution.

What advice can I offer? Well, if we are talking about social revolutions, then I should say, "Don't trust that Robespierre guy. He's going to kill us all!" Or perhaps, "Napoleon is a genocidal maniac! He hasn't come back from Elbe to make it all better!" But such dire warnings are ineffective because no one listens.

What should I say here?

Perhaps I should simply apologise for my small part in this mess we find ourselves in.

Because it has all turned out so horrendously bad, I think we should have been more aware of the risks, even if at the time they may have sounded totally far-fetched. We said of the Internet: "This is so good everyone should be able to play." And we said: "The Internet is for everyone." But we never really thought about what we really meant when we proclaimed the universality of the Internet. "Everyone should be able to do this?" has turned into "What have we done to ourselves?"

Code

I am probably not a brilliant programmer. In fact, I should admit that I'm a shocking programmer—and I know I'm not the only one. In fact, I'm probably pretty average as a coder. And if that's the average in our profession, then all I can say is that we are all shockingly bad programmers.

We are building these massive edifices of mind-boggling complexity and then replicating all this rather shoddy software in billions of devices. We were told to "move fast and break things," and we did exactly that. We learned to use the end user as the test case. But the consequences are ugly. Your average car has at least 300 processors and huge amounts of code. It mostly works, but just remember that the network that contains the drive control systems also probably contains the entertainment system. *And all this complexity is probably provided by the lowest bidder!* It's cheaper that way. And much riskier. Modern machinery is now at a level of complexity that visibly defies human understanding or control. God alone knows exactly what is in the software-controlled systems on a Boeing 737 Max 8. Boeing apparently does not. Or even the firmware in your fancy digital front doorbell. Bitter experience has taught us that we can turn a few hundred million baby web cams into a massively destructive attack force within seconds.

The *Key Performance Indicators* in our industry are best described by the currently oh-so-fashionable Agile process: “Let’s write even crappier code even faster, and let’s break more of it!”

Nobody knows how these systems work anymore. Nobody truly understands the dependencies anymore, if they ever did, and the continual stream of software upgrades should give you ample evidence that we are only just bailing out the bilge as fast as we can to stop the entire ship from sinking!

We spend hundreds of millions of dollars on staffing shiny cyber defence bodies to try to show what a great job we are doing to defend ourselves when, in fact, the problem is not the folks who are driving the hostile trucks through the wide-open doors. The real problem is that it’s the people just like me who produced the insanely poor code in the first place who left all these gaping holes behind them. Because none of us really is up to the task. And I don’t know about everybody else, but I am still on the keyboard. Still writing code. Collectively we have done an amazing job. The Internet is now busted! And it’s not clear that we can fix the problem. We can’t make it better. Sorry.

Security

It is evident that we have no desire to build truly secure systems. In the rush to digitise our world of services we are taking extraordinary risks. The term “web security” is the punchline to some demented sick joke because the online world is held together by a level of naive trust that makes all other forms of human credulity look restrained and cautious! Even when we thought about what better security might look like, the response was that we have neither the time nor the money to do a better job. We believe that the consumer is so impatient that milliseconds matter far more than security. We continue to cut corners and build fast, faulty code. Maybe we should have said “no” and walked away from the keyboard. But we didn’t. Sorry.

We thought we were helping people communicate, because after all, communication is what drives the human experience. We knew that if you change how we communicate you change the nature of human society. We knew that. But we didn’t consider that message seriously. None of us envisioned the perversion of that nobly motivated ambition into the incessant deluge of waste products from the social media factory. We only appreciated the role of content mediators when we eliminated them from the planet. This situation is not pretty. We choose to listen only to what we agree with. The Internet has become a vanity-reinforcing gigantic distorted selfie. Sorry.

No Rules

We built this new world so quickly that we outpaced everything else. This new technology has no controls, no regulation, no competition. In the rush to be the first to unleash the ruthless forms of surveillance capitalism on an unsuspecting populace, we have bypassed all the conventional forms of care and restraint. Just seven digital giants dominate our world.

Their unstinting efforts to lobby politicians has turned the political process into a fatally corrupted empty shell. We moved too quickly and no one else kept up.

We wrote our own rules, and Rule Number One was: “Just do it.” From Uber to Google the word was “disruption.” But the wholesale destruction of the old-world business environment wasn’t the worst thing we did. Destruction of retail shopping wasn’t the worst thing we did. Far worse is that we privatised the public communications space. We turned our culture and our public discourse into private property. We privatised our intellectual achievements. Who owns antibiotics? Who owns my genetic code? Who owns my personal profile? We turned everything into a transaction. We destroyed our libraries and replaced them with search engines. We replaced journalism with tweets. Our world is no longer a collection of public spaces, but a collection of private enclosures. In some small way, I helped build that reality. Sorry.

No Way Back

Can I provide some helpful suggestions, offer some motivation, or provide some palliative comfort by asserting that our voices matter, and we can change our world for the better? No. I think that we already betrayed you 30 years ago. The glittering prizes that this new technology promised us turned out to be tawdry, corrupted, and debased. We thought technology would be a compelling force for good. We were wrong. I am truly sorry.

The task before us right now is not to make it better. That is way too ambitious. We just can’t make it better. There is no way to back out now. Having unleashed these digital monstrosities, we cannot just tie them up again and put everything back into a box. That we cannot do. The best we can do is to somehow accept the terrible situation and the betrayal of trust that got us here and try to deal with it without making it even worse.

Sorry.

—*Geoff Huston*, gih@apnic.net

Upcoming Articles in IPJ

“Automatic Disaggregation in the Routing in Fat Trees (RIFT) Protocol,” by Bruno Rijsman. RIFT is a new routing protocol being defined in the IETF. This article focuses on one particular feature of RIFT, namely automatic aggregation and disaggregation.

“Network Functions Virtualization (NFV),” by William Stallings. NFV provides a powerful, vendor-independent approach to implementing complex networks with dynamic demands. NFV builds on well-established technologies, including virtual machines, containers, and virtual networks. With the demand from 5G and cloud service providers, as well as enterprises with large internal networks, NFV is becoming an increasingly widespread technology.

Postal Service Award Presented to Onno W. Purbo

The Internet Society, a global nonprofit organization that promotes the development and use of an open, globally connected, and secure Internet, recently presented the prestigious *Jonathan B. Postel Service Award* to Onno W. Purbo for his sustained and substantial technical contributions, leadership, and service to the global Internet community.

Named in honor of the technical community legend Jonathan Postel, this award recognizes extraordinary people like Mr. Purbo who have committed themselves to the technological development, growth, and strength of the Internet. Known as “Indonesia’s Internet



Liberator,” Mr. Purbo is a prolific and well-published Internet advocate who has played a key role in democratizing Internet access, making it more affordable especially in Indonesia’s rural areas.

“Mr. Purbo’s contribution to the digital sector is invaluable and this award marks what he has achieved and inspired others to achieve. His initiative of meaningful Internet access and Community Networks have instilled the growth of not only affordable but accessible Internet in various areas across Indonesia. I am confident this award will embolden others to innovate and follow his steps and overcome the challenges in their communities especially in improving digitalization,” said Johnny Plate, Minister of Communication and Information Technology for Indonesia.

Of his many achievements, Mr. Purbo is best known for pioneering the Internet in Indonesia through sophisticated use of wireless and *Voice over Internet Protocol* technologies. He led the first Internet connection at the Institute of Technology in Bandung and used it to build the first Indonesian educational network. He also championed the deregulation of WiFi frequencies and introduced cyber cafes, neighborhood networks, and community cellular networks to Indonesia. Mr. Purbo organized the first community telephony network over Internet and led the re-introduction of ICT into the Indonesian high school curriculum.

Currently, he is involved in the largest Indonesian FREE e-Learning service, which has brought more than 700 courses to nearly 40,000 participants and trained more than 8,000 teachers on e-learning operations.

“It is an honor to receive the highest and priceless acknowledgment given to Indonesia from the Internet communities,” said Mr. Purbo.

“With modified simple off-the-shelf gadgets and equipment, one may fulfill the right to access information and knowledge, which is the necessary foundation for any nation to move forward. The Internet Society has acknowledged the approach is one of the right routes towards the Internet for all. The job is indeed not finished. The Postel Service Award sheds light on the way to go for all of us and inspires extraordinary enthusiasm for moving towards a knowledge-based society.”

Mr. Purbo was selected by a distinguished international committee comprised of former Postel Award winners which includes Internet visionaries and luminaries. Now in its 21st year, the Postel Award was established in 1999 by the Internet Society to honor individuals and organizations that, through their work, embody the spirit of Jonathan Postel, whose technical influence can be seen at the very heart of many of the protocols which make the Internet work. Andrew Sullivan, President and CEO of the Internet Society, presented the award, which includes a US\$20,000 honorarium and a crystal engraved globe, during a virtual ceremony as part of the 109th *Internet Engineering Task Force* (IETF) meeting which took place November 16–20, 2020.

For more information, please visit:
<https://www.internetsociety.org>.

History of Networking Recordings

Russ White writes: “In 2017, I realized a lot of the people I’ve worked with over the years were retiring. When these people leave the networking community, they take a wealth of knowledge about the intent, challenges, and inventions of the early Internet. I decided to capture as much of this history in oral format as possible—hence the history of networking recordings were started. I thought, at first, this would be a small, short-lived series, but I have been amazed by the reaction of the community, and the number of technologies and organizations involved in the design and operation of computer networks.

If you know of someone who should be here, please contact me, as I would like to collect as much oral history in this area as I can for this and future generations. These recordings are released under Creative Commons License (CC BY-NC-ND 4.0). This means recordings can be distributed for any noncommercial purposes by anyone, so long as they are released in full (with no modifications).”

The recordings can be found here:
<https://rule11.tech/history-of-networking/>

NSA Recommends How Enterprises Can Securely Adopt Encrypted DNS

The *National Security Agency* (NSA) recently released a cybersecurity document, “Adopting Encrypted DNS in Enterprise Environments,”^[1] explaining the benefits and risks of adopting the encrypted *Domain Name System* (DNS) protocol, *DNS over HTTPs* (DoH), in enterprise environments. The document provides solutions for secure implementation based on enterprise network needs.

DNS translates domain names in URLs into IP addresses, making the Internet easier to navigate. However, it has become a popular attack vector for malicious cyber actors. DNS shares its requests and responses in plaintext, which can be easily viewed by unauthorized third parties. Encrypted DNS is increasingly being used to prevent eavesdropping and manipulation of DNS traffic. As encrypted DNS becomes more popular, enterprise network owners and administrators should fully understand how to properly adopt it on their own systems. Even if not formally adopted by the enterprise, newer browsers and other software may try to use encrypted DNS anyway and bypass the enterprise’s traditional DNS-based defenses.

DoH encrypts DNS requests, preventing eavesdropping and manipulation of DNS traffic. While good for ensuring privacy in home networks, DoH can present risks to enterprise networks if it isn’t appropriately implemented. The recommendations detailed will assist enterprise network owners and administrators in balancing DNS privacy and governance for their networks. It outlines the importance of configuring enterprise networks appropriately to add benefits to, and not hinder, their DNS security controls. These enterprise DNS controls can prevent numerous threat techniques used by cyber threat actors for initial access, command and control, and exfiltration.

NSA recommends that an enterprise network’s DNS traffic, encrypted or not, be sent only to the designated enterprise DNS resolver. This ensures proper use of essential enterprise security controls, facilitates access to local network resources, and protects internal network information. All other DNS resolvers should be disabled and blocked.

NSA seeks to regularly release unique, actionable, and timely cybersecurity guidance to secure the Department of Defense, National Security Systems, and the Defense Industrial Base. For more information or other cybersecurity products, visit:

<https://www.NSA.gov/cybersecurity-guidance>.

[1] https://media.defense.gov/2021/Jan/14/2002564889/-1/-1/0/CSI_ADOPTING_ENCRYPTED_DNS_U_OO_102904_21.PDF

WebRTC Becomes a Standard

The *World Wide Web Consortium* (W3C) and the *Internet Engineering Task Force* (IETF) recently announced that *Web Real-Time Communications* (WebRTC), which powers myriad services, is now an official standard, bringing audio and video communications anywhere on the Web.

WebRTC, comprised of a *JavaScript* API for Web Real-Time Communications and a suite of communications protocols, allows any connected device, on any network, to be a potential communication endpoint, on the Web. WebRTC already serves as a cornerstone of online communication and collaboration services. The WebRTC framework provides the building blocks from which web and app developers can seamlessly add video chat to a range of applications, including tele-education and tele-health, entertainment and gaming, professional and workforce collaboration.

With the foundations standardized and deployed as a royalty-free feature in Web browsers and other devices and platforms, setting up a secure audio-video communication system with WebRTC has become a built-in capability, eliminating the need to install plugins or download separate applications.

WebRTC is massively deployed as a communications platform and powers video conferences and collaboration systems across all major browsers, both on desktop and mobile. Billions of users can interact now that WebRTC makes live video chat easier than ever on the Web. In commercial products and open source projects, WebRTC has vastly expanded the ability to deploy real-time interaction solutions to customers and users.

The year 2020 has shown both how critical WebRTC already is in a world where travel and physical contacts need to be limited, as well as the many improvements that can be brought to the technology to address new usages that have emerged. Organizations are leveraging WebRTC to conduct training, interviews, strategic planning or as a substitute for in-person meetings. Schools and universities have shifted to virtual learning platforms. Families and friends make daily use of products that are built with WebRTC or parts of it.

With the use of WebRTC expanding beyond the initial core design to power video conferences and collaboration systems in web browsers and other ecosystems, more features and more optimizations are now needed. The IETF *WebTransport* work is aiming to build out additional web support for a variety of transport properties. The *WebRTC Ingest Signaling over HTTPS* work is focusing on the development of a protocol to support one-way WebRTC-based audiovisual sessions between broadcasting tools and real-time media broadcast networks. Similar work to expand the use cases of WebRTC is ongoing in the W3C. For more information visit:

<https://www.w3.org/TR/webrtc/>

<https://www.ietf.org/blog/webrtc-milestone/>

Nominations Open for Prestigious Internet Hall of Fame

The Internet Society recently announced that nominations are now open for the next *Internet Hall of Fame* class of inductees. The nomination period will close April 23, 2021 and inductees will be announced at an awards ceremony to be held later this year. The Internet Hall of Fame, now in its tenth year, recognizes a select group of visionaries, leaders and luminaries who have made significant contributions to the development and advancement of the open, global Internet.

Through the work of these individuals, including Vint Cerf, Robert Kahn, Leonard Kleinrock, Tim Berners-Lee, and Elizabeth Feinler, among many others, the Internet Hall of Fame reflects the history of the Internet's development and evolution.

"At no point in time has the importance of the Internet and its chief characteristic—to connect—been felt so broadly, and so acutely," said Andrew Sullivan, President and CEO of the Internet Society.

"The critical role the Internet has played throughout the pandemic reinforces now, more than ever, the significance of the people who originally conceived, built, guided and promoted this global network. It is our privilege to highlight their work and contributions."

Individuals worldwide who have played an extraordinary role in the conceptualization, building, and development of the Internet globally will be considered for induction. In addition to those who have been more visible, the Internet Hall of Fame also seeks nominees who have made crucial, behind-the-scenes contributions. Criteria for evaluation include:

Impact: The contribution has made an extraordinary impact on the development or growth of the Internet, and was and may still be directly relevant to the Internet's ongoing advancement and evolution.

Influence: The contribution, relative to the Internet, has significantly influenced: 1) the work of others in the field; 2) society at large; or 3) another more defined but critical audience or region.

Innovation: The contribution has broken new ground with original thinking/creativity that has established new paradigms, eliminated significant obstacles, or accelerated Internet advancements.

Reach: The contribution has significantly impacted the Internet's reach among society at large, within key audiences or specific geographies, with global impact.

Founded in 2012, the Internet Hall of Fame is an ongoing awards program established by the Internet Society to recognize a distinguished and select group of leaders and luminaries who have made significant contributions to the development and advancement of the global open Internet. More information on the program can be found at <http://www.internethalloffame.org/>.

Domain Abuse Activity Reporting

ICANN's *Domain Abuse Activity Reporting* (DAAR) project is a system for studying and reporting on domain name registration and security threat (domain abuse) behavior across *top-level domain* (TLD) registries. The overarching purpose of DAAR is to develop a robust, reliable, reproducible, and replicable methodology for analyzing security threat activity that can then be later used by the ICANN community to facilitate informed policy decisions.

The system collects TLD zone data and complements these data sets with a large set of high-confidence reputation (security threat) data feeds. The aggregated and anonymized data collected by the DAAR system can serve as a platform for studying or reporting daily or historical registration or abuse activity by each registry. The data is currently being pushed to registries using the ICANN *Service Level Agreement Monitoring* (SLAM) system.

The data collected out of the DAAR system is being used to generate the DAAR monthly reports. The reports are point-in-time analysis of all TLDs for which data was available. The report provides aggregated statistics and time-series analysis about security threats of interest to DAAR namely phishing, malware, spam, and botnet command-and-control. For more information visit:

<https://www.icann.org/octo-ssr/daar>

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. For more information, contact us at ipj@protocoljournal.org

Check your Subscription Details!

If you have a print subscription to this journal, you will find an expiration date printed on the back cover. For several years, we have “auto-renewed” your subscription, but now we ask you to log in to our subscription system and perform this simple task yourself. Make sure that *both* your postal and e-mail addresses are up-to-date since these are the only methods by which we can contact you. If you see the words “Invalid E-mail” on your copy this means that we have been unable to contact you through the e-mail address on file. If this is the case, please contact us at ipj@protocoljournal.org with your new information. The subscription portal is located here: <https://www.ipjsubscription.org/>

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Kjetil Aas	Darrell Budic	Holger Durer	Christopher Guemez	Christos Karayiannis
Fabrizio Accatino	Scott Burleigh	Mark Eanes	Gulf Coast Shots	David Kekar
Michael Achola	Chad Burnham	Andrew Edwards	Sheryll de Guzman	Stuart Kendrick
Martin Adkins	Jon Harald Bøvre	Peter Robert Egli	Rex Hale	Robert Kent
Melchior Aelmans	Olivier Cahagne	George Ehlers	Jason Hall	Jithin Kesavan
Christopher Affleck	Antoine Camerlo	Peter Eisses	James Hamilton	Jubal Kessler
Scott Aitken	Tracy Camp	Torbjörn Eklöv	Stephen Hanna	Shan Ali Khan
Jacobus Akkerhuis	Ignacio Soto Campos	Y Ertur	Martin Hannigan	Nabeel Khatri
Antonio Cuñat Alario	Fabio Caneparo	ERNW GmbH	John Hardin	Dae Young Kim
Nicola Altan	Roberto Canonico	ESdatCo	David Harper	William W. H. Kimandu
Matteo D'Ambrosio	David Cardwell	Steve Esquivel	Edward Hauser	John King
Selva Anandavel	John Cavanaugh	Jay Etchings	David Hauweel	Russell Kirk
Jens Andersson	Lj Cemerar	Mikhail Evstiounin	Marilyn Hay	Gary Klesk
Danish Ansari	Dave Chapman	Bill Fenner	Headcrafts SRLS	Anthony Klopp
Finn Arildsen	Stefanos Charchalakis	Paul Ferguson	Hidde van der Heide	Henry Kluge
Tim Armstrong	Greg Chisholm	Ricardo Ferreira	Johan Helsingius	Michael Kluk
Richard Artes	David Chosrova	Kent Fichtner	Robert Hinden	Andrew Koch
Michael Aschwanden	Marcin Cieslak	Armin Fisslthaler	Asbjørn Højmark	Ia Kochiashvili
David Atkins	Lauris Cikovskis	Michael Fiumano	Damien Holloway	Carsten Koempfe
Jac Backus	Guido Coenders	The Flirble Organisation	Alain Van Hoof	Richard Koene
Jaime Badua	Brad Clark	Gary Ford	Edward Hotard	Alexander Kogan
Bent Bagger	Narelle Clark	Jean-Pierre Forcioli	Bill Huber	Antonin Kral
Eric Baker	Horst Clausen	Susan Forney	Hagen Hultzs	Robert Krejčí
Santosh Balagopalan	Joseph Connolly	Christopher Forsyth	Kevin Iddles	Mathias Körber
Benjamin Barkin-Wilkins	Steve Corbató	Andrew Fox	Mika Ilvesmaki	John Kristoff
Michael Bazarewsky	Brian Courtney	Craig Fox	Karsten Iwen	Terje Krogdahl
David Belson	Beth and Steve Crocker	Fausto Franceschini	David Jaffe	Bobby Krupczak
Hidde Beumer	Dave Crocker	Valerie Fronczak	Ashford Jaggernaut	Murray Kucherawy
Pier Paolo Biagi	Kevin Croes	Tomislav Futiv	Martijn Jansen	Warren Kumari
Tyson Blanchard	John Curran	Laurence Gagliani	Jozef Janitor	George Kuo
John Bigrow	André Danthine	Edward Gallagher	John Jarvis	Dirk Kurfuerst
Orvar Ari Bjarnason	Morgan Davis	Andrew Gallo	Dennis Jennings	Darrell Lack
Axel Boeger	Jeff Day	Chris Gamboni	Edward Jennings	Andrew Lamb
Keith Bogart	Julien Dhallenne	Xosé Bravo Garcia	Aart Jochem	Richard Lamb
Mirko Bonadei	Freek Dijkstra	Oswaldo Gazzaniga	Brian Johnson	Yan Landriault
Roberto Bonalumi	Geert Van Dijk	Kevin Gee	Curtis Johnson	Edwin Lang
Julie Bottorff Photography	David Dillow	Greg Giessow	Richard Johnson	Sig Lange
Gerry Boudreaux	Richard Dodsworth	John Gilbert	Jim Johnston	Markus Langenmair
L de Braal	Ernesto Doelling	Serge Van Ginderachter	Jonatan Jonasson	Fred Langham
Kevin Breit	Michael Dolan	Greg Goddard	Daniel Jones	Tracy LaQuey Parker
Thomas Bridge	Eugene Doroniuk	Tiago Goncalves	Gary Jones	Rick van Leeuwen
Ilia Bromberg	Karlheinz Dölger	Ron Goodheart	Jerry Jones	Simon Leinen
Václav Brožík	Joshua Dreier	Octavio Alfageme	Anders Marius Jørgensen	Robert Lewis
Christophe Brun	Lutz Drink	Gorostiaga	Amar Joshi	Christian Liberale
Gareth Bryan	Dmitriy Dudko	Barry Greene	Javier Juan	Martin Lillepuu
Stefan Buckmann	Andrew Dul	Jeffrey Greene	David Jump	Roger Lindholm
Caner Budakoglu	Joan Marc Riera	Richard Gregor	Merike Kao	Link Light Networks
	Duocastella	Martijn Groenleer	Andrew Kaiser	Sergio Loreti
	Pedro Duque	Geert Jan de Groot		

Eric Louie	Joel Moore	Eduard Llull Pou	Dan Schrenk	Douglas Thompson
Adam Loveless	John More	Tim Pozar	Richard Schultz	Kerry Thompson
Guillermo a Loyola	Maurizio Moroni	David Raistrick	Timothy Schwab	Lorin J Thompson
Hannes Lubich	Brian Mort	Priyan R Rajeevan	Roger Schwartz	Fabrizio Tivano
Dan Lynch	Soenke Mumm	Balaji Rajendran	SeenThere	Joseph Toste
Sanya Madan	Tariq Mustafa	Paul Rathbone	Scott Seifel	Rey Tucker
Miroslav Madić	Stuart Nadin	William Rawlings	Yury Shefer	Sandro Tumini
Alexis Madriz	Michel Nakhla	Mujtiba Raza Rizvi	Yaron Sheffer	Angelo Turetta
Carl Malamud	Mazdak Rajabi Nasab	Bill Reid	Doron Shikmoni	Phil Tweedie
Jonathan Maldonado	Krishna Natarajan	Petr Rejhon	Tj Shumway	Steve Ulrich
Michael Malik	Naveen Nathan	Robert Remenyi	Jeffrey Sicuranza	Unitek Engineering AG
Tarmo Mamers	Darryl Newman	Rodrigo Ribeiro	Thorsten Sideboard	John Urbanek
Yogesh Mangar	Thomas Nikolajsen	Glenn Ricart	Greipur Sigurdsson	Martin Urwaleck
Bill Manning	Paul Nikolich	Justin Richards	Andrew Simmons	Betsy Vanderpool
Harold March	Travis Northrup	Rafael Riera	Pradeep Singh	Surendran
Vincent Marchand	Marijana Novakovic	Mark Risinger	Henry Sinnreich	Vangadasalam
Gabriel Marroquin	David Oates	Fernando Robayo	Geoff Sisson	Ramnath Vasudha
David Martin	Ovidiu Obersterescu	Gregory Robinson	Helge Skrivervik	Philip Venable
Jim Martin	Tim O'Brien	Ron Rockrohr	Darren Sleeth	Buddy Venne
Ruben Tripiana Martin	Mike O'Connor	Carlos Rodrigues	Richard Smit	Alejandro Vennera
Timothy Martin	Mike O'Dell	Magnus Romedahl	Bob Smith	Luca Ventura
Carles Mateu	John O'Neill	Lex Van Roon	Courtney Smith	Tom Vest
Juan Jose Marin	Jim Oplotnik	Alessandra Rosi	Eric Smith	Dario Vitali
Martinez	Packet Consulting	David Ross	Mark Smith	Jeffrey Wagner
Ioan Maxim	Limited	William Ross	Craig Snell	Don Wahl
David Mazel	Carlos Astor Araujo	Boudhayan	Job Snijders	Michael L Wahrman
Miles McCredie	Palmeira	Roychowdhury	Ronald Solano	Laurence Walker
Brian McCullough	Alexis Panagopoulos	Carlos Rubio	Asit Som	Randy Watts
Joe McEachern	Gaurav Panwar	Rainer Rudigier	Ignacio Soto Campos	Andrew Webster
Alexander McKenzie	Manuel Uruena Pascual	Timo Ruiters	Evandro Sousa	Tim Weil
Jay McMaster	Ricardo Patara	RustedMusic	Peter Spekrijse	Jd Wegner
Mark Mc Nicholas	Dipesh Patel	Babak Saberi	Thayumanavan Sridhar	Westmoreland
Carsten Melberg	Alex Parkinson	George Sadowsky	Paul Stancik	Engineering Inc.
Kevin Menezes	Craig Partridge	Scott Sandefur	Ralf Stempfer	Rick Wesson
Bart Jan Menkveld	Dan Paynter	Sachin Sapkal	Matthew Stenberg	Peter Whimp
Sean Mentzer	Leif Eric Pedersen	Arturas Satkovskis	Adrian Stevens	Russ White
William Mills	Rui Sao Pedro	PS Saunders	Clinton Stevens	Jurrien Wijlhuizen
David Millsom	Juan Pena	Richard Savoy	John Streck	Derick Winkworth
Desiree Miloshevic	Chris Perkins	John Sayer	Martin Streule	Pindar Wong
Joost van der Minnen	Michael Petry	Phil Scarr	David Strom	Phillip Yialeloglou
Thomas Mino	Alexander Peuchert	Gianpaolo	Viktor Sudakov	Janko Zavernik
Rob Minshall	David Phelan	Scassellati	Edward-W. Suor	Muhammad Ziad
Wijnand Modderman	Derrell Piper	Elizabeth Scheid	Vincent Surillo	Ziayuddin
Mohammad Moghaddas	Rob Pirnie	Jeroen Van Ingen	Terence Charles	Jose Zumalave
Roberto Montoya	Marc Vives Piza	Schenau	Sweetser	Romeo Zwart
Charles Monson	Jorge Ivan Pincay Ponce	Carsten Scherb	T2Group	Bernd Zeimet
Andrea Montefusco	Victoria Poncini	Ernest Schirmer	Roman Tarasov	廖明沂
Fernando Montenegro	Blahoslav Popela	Philip Schneck	David Theese	



Follow us on Twitter and Facebook

@protocoljournal



<https://www.facebook.com/newipj>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Diamond Sponsors



Ruby Sponsors



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

David Conrad, Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Principal – Internet Technology, Policy, and Advocacy
The Internet Society

Dr. Jun Murai, Founder, WIDE Project
Distinguished Professor, Keio University
Co-Director, Keio University Cyber Civilization Research Center, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

July 2021

Volume 24, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

FROM THE EDITOR

In This Issue

From the Editor	1
Disaggregation in RIFT	2
Network Functions Virtualization.....	15
Fragments	28
Thank You!	32
Call for Papers	34
Supporters and Sponsors	35

Since the launch of *The Internet Protocol Journal* in 1998, we have covered several aspects of the core technologies used in the global Internet and in enterprise networks. Routing protocols such as the *Border Gateway Protocol* (BGP) continue to play a critical role in the operation of these networks, but other routing protocols are being developed by the *Internet Engineering Task Force* (IETF) for use in data center environments. One such protocol is the *Routing in Fat Trees Protocol* (RIFT).

Route *aggregation* is used to summarize a set of specific routing-table entries into a single, less-specific route, in order to reduce the size of routing tables. Disaggregation is the opposite of aggregation whereby an aggregate route is divided into several more-specific routes. In our first article, Bruno Rijsman explains how automatic disaggregation is accomplished in RIFT.

Security has also been a recurring theme in this journal. Most of the protocols used in today's Internet were originally designed without comprehensive security in mind, but the IETF has produced security enhancements for many of the core protocols. Securing the routing system itself has proven challenging because it requires widespread deployment in order to be effective. Starting with our next issue, Geoff Huston presents a two-part article entitled "A Survey on Securing Inter-Domain Routing." Make sure your subscription details are up to date!

The IETF is, of course, not the only organization that produces standards for computer networks. Our second article, by William Stallings, is an overview of *Network Functions Virtualization*, an emerging set of standards being developed by the *European Telecommunications Standards Institute* (ETSI).

The generous support of individuals and organizations makes publication of this journal possible. We are pleased to welcome our latest sponsor, *The APNIC Foundation*. More information about the foundation is available on page 30 of this issue.

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

Automatic Disaggregation in the Routing in Fat Trees Protocol

by Bruno Rijsman

Routing in Fat Trees (RIFT) is a new routing protocol being defined in the *Internet Engineering Task Force* (IETF).^[1] RIFT is optimized for large networks that have a highly structured topology such as fat tree, *Clos*, or similar topologies. It is typically used as a scalable and fast-converging *Interior Gateway Protocol* (IGP) for the underlay in data centers, but it has other use cases as well.^[2]

RIFT brings several innovations to the table without requiring any changes to existing networking hardware (“silicon”), including:

- *Zero Touch Provisioning* (ZTP) virtually eliminates the need for configuration and auto-detects miscabling.
- RIFT is *anisotropic*: it is a link-state protocol north-bound and a distance-vector protocol south-bound, combining the advantages of link-state with the advantages of distance-vector/path-vector.
- RIFT is inherently loop-free, allowing it to distribute traffic across all available paths (not just the shortest path).
- A built-in flooding-reduction mechanism greatly reduces flooding traffic in densely connected topologies such as fat trees.
- With the automatic aggregation feature, in the absence of failures, each node needs only a single multi-path default route pointing north. This feature reduces the size of the routing tables at or close to the leaf nodes, and hence the cost of top-of-rack switches.
- With the automatic disaggregation feature, if a failure occurs the north-bound default route is automatically disaggregated into more specific routes, but only to the extent needed to route around the failure.
- RIFT supports large data-center networks without the need for splitting the network into multiple areas.
- RIFT offers very fast convergence—even in very large networks.
- Because of the simplicity of RIFT functionality on leaf switches, you can easily run it on servers; this feature is also known as *Routing on The Host* (RoTH). It enables support for multi-homed servers with automatic recovery from link and node failures.
- Model-based (Thrift) specification of the routing protocol messages accelerates development, enhances interoperability, and most importantly, improves security by removing most message-parsing vulnerabilities.

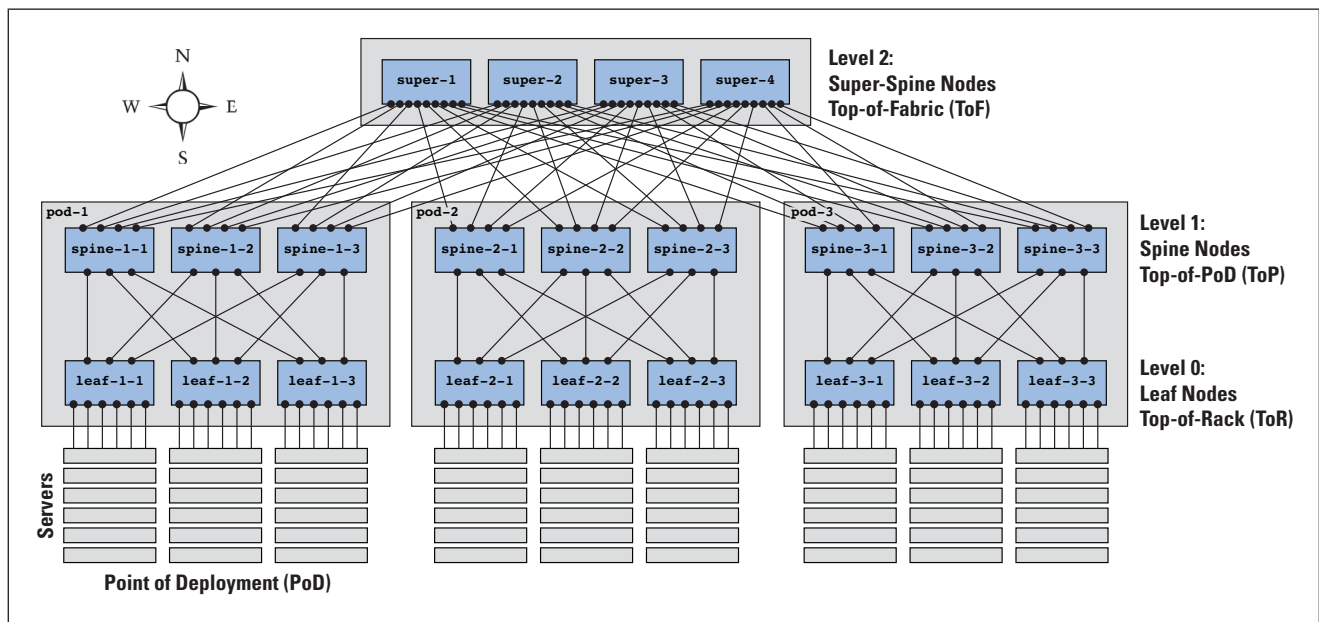
RIFT offers an open-source implementation^[3] and at least one commercial implementation.^[4]

In this article, we focus on one feature of RIFT: automatic aggregation and disaggregation, which is one of the most novel and most interesting innovations in the RIFT protocol. For a more general overview of RIFT, see the presentation at APNIC^[5] or the recently released (and free) *Day One book on RIFT*.^[6] For a discussion of link-state routing in data centers (including RIFT), see the article “Recent Developments in Link State on Data-Center Fabrics,” also published in this journal.^[7]

Introduction to RIFT

You can use RIFT in topologies where it makes sense to speak of north and south directions, which allows you to divide the nodes into levels, including fat-tree data center topologies such as the one shown in Figure 1.

Figure 1: Fat-Tree Data Center Topology



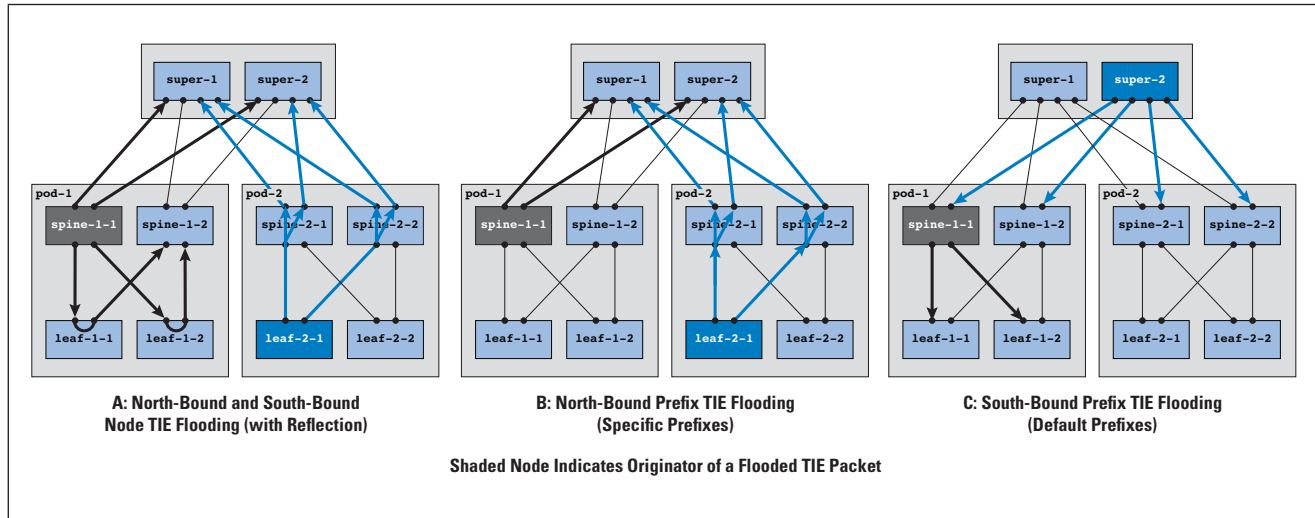
In many respects, RIFT is a link-state protocol similar to *Intermediate System-to-Intermediate System* (IS-IS):

- RIFT routers exchange hello packets, called *Link Information Element* (LIE) packets, to establish adjacencies with neighbor routers.
- RIFT routers originate link-state packets, called *Topology Information Element* (TIE) packets, to describe the state, adjacencies, originated prefixes, disaggregated prefixes, and other information about the router.
- RIFT reliably floods the link-state packets across the network. It uses *Topology Information Description Element* (TIDE) packets to summarize the contents of the link-state database and *Topology Information Request Element* (TIRE) packets to acknowledge and request TIE packets. Together, TIDEs and TIREs are used to make the flooding reliable.

- RIFT stores link-state packets in its *Link State Database* (LSDB).
- RIFT runs the *Shortest Path First* (SPF) algorithm on the topology stored in the link-state database to compute the shortest path to each destination.

RIFT is unique in that it has different rules for flooding TIE packets in both the north-bound south-bound directions, as shown in Figure 2:

Figure 2: TIE Flooding Rules



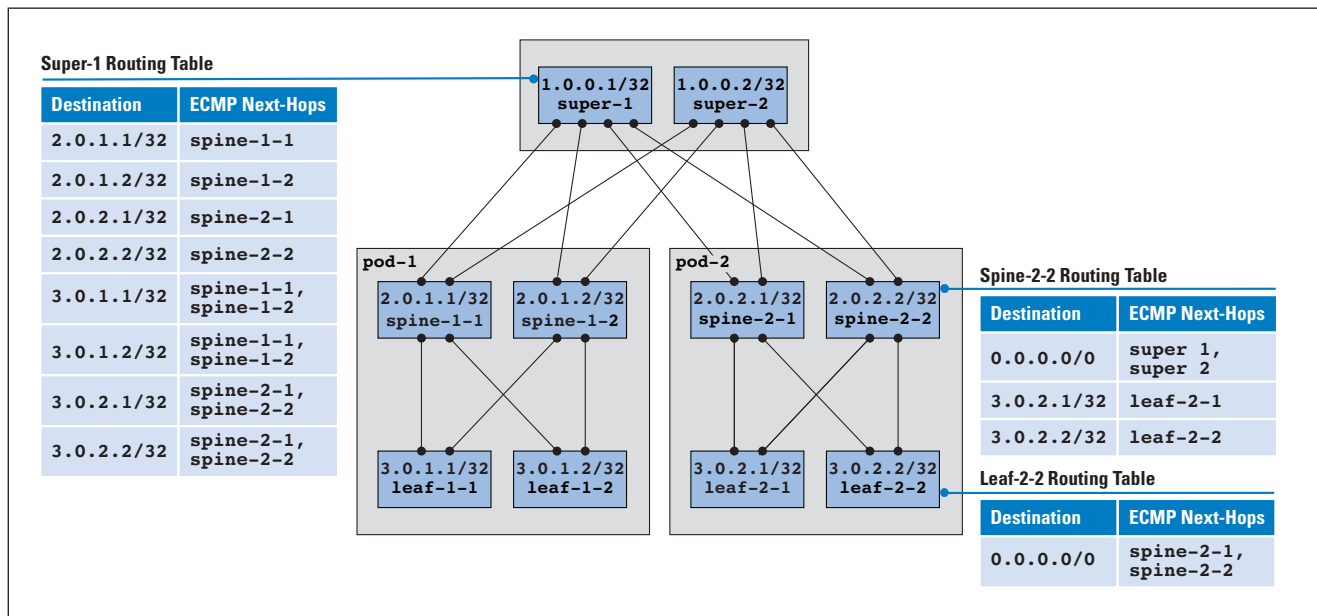
- Each node advertises its adjacencies in node TIEs that are flooded in both the north-bound and south-bound directions. Furthermore, the level just below the originator of the node TIE “reflects” the node TIEs back-up. This reflection allows nodes to discover the adjacencies of other nodes at the same level (refer to diagram A in Figure 2).
- Each node advertises its local prefixes in prefix TIEs, which are flooded only in the north-bound direction. (Shown in diagram B in Figure 2).
- Each node advertises a fabric default route (typically `0.0.0.0/0` and `::/0`) in prefix TIEs that are flooded exactly one hop (but no further) in the south-bound direction (see diagram C in Figure 2). The top-of-fabric nodes always originate a default, and the lower nodes originate a default only if they have received at least one default from a parent. This model makes the south-bound flooding similar to distance-vector routing and it is the reason that RIFT is colloquially described as link-state towards the spine and distance-vector towards the leaves.
- RIFT also allows for east-west “short-cut” links and has flooding rules for those links (not shown in the figure).

- RIFT also includes a flooding-reduction mechanism that avoids multiple copies of the same TIE being sent to the same node (that mechanism is not shown here). For example, in diagrams A and B in Figure 2 node super-1 receives two identical copies of the TIE from leaf-2-1.

After the TIEs are flooded across the network in the manner described previously, the RIFT nodes run the SPF algorithm to compute the routing tables. Actually, RIFT does at least two SPF runs: one for the north-bound and one for the south-bound direction.

Figure 3 shows an example of typical RIFT routing tables (in the absence of failures):

Figure 3: Typical RIFT Routing Tables in the Absence of Failures



We can see the following routes:

- *Specific routes for all south-bound traffic:* These routes are typically host /32 (for IPv4) or /128 (for IPv6) routes.
- *Fabric default routes for all north-bound traffic:* These routes are typically 0.0.0.0/0 or ::/0 routes.

Both the north-bound default routes and the south-bound, specific routes, are multi-path routes, distributing the traffic across all available paths. The next-hops can be weighted according to the bandwidth available on each path.

RIFT Automatic Aggregation and Disaggregation

The *aggregation*^[8] concept has existed in routing protocols since the beginning. Aggregation allows you to summarize a set of specific routes by a single, less-specific route, called the *aggregate route*. The most common use case for aggregation is to reduce the size of the routing table by summarizing unneeded details.

The concept of *disaggregation*^[9] has also been used for a long time. Disaggregation is the opposite of aggregation: it takes a single less-specific route (the aggregate route) and divides it into several more-specific routes. The most common use case for disaggregation is traffic engineering.

In existing protocols such as the *Border Gateway Protocol* (BGP), *Open Shortest Path First* (OSPF), and IS-IS, aggregation and disaggregation typically are manually configured for optimization purposes. In RIFT, on the other hand, aggregation and disaggregation are automatic and always enabled. RIFT automatically aggregates routes (typically to the default route) wherever possible. And RIFT automatically disaggregates routes wherever needed, for example, because of link or node failures.

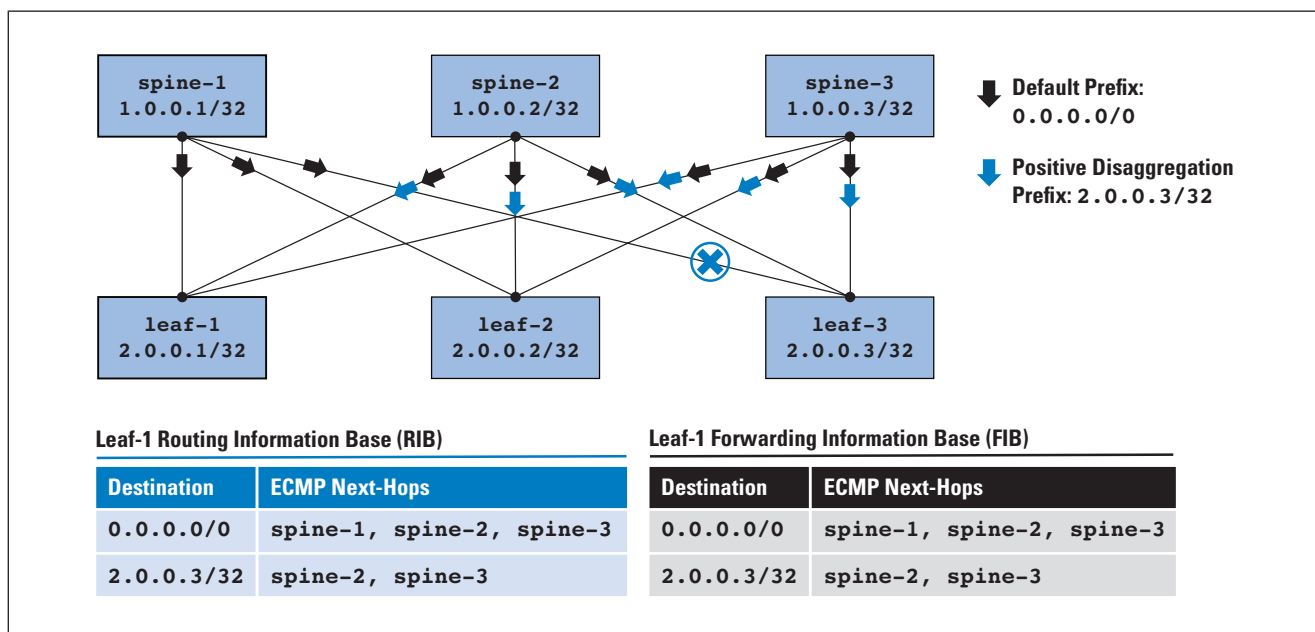
Disaggregation actually has two modes in RIFT:

- *Positive Disaggregation* works by advertising a more specific route to “attract” traffic to a repair path away from a failed path. Advertising more-specific prefixes is exactly how disaggregation works in existing protocols.
- *Negative Disaggregation* works by advertising a so-called *negative prefix* to “repel” traffic away from a failed path towards a repair path. This new mechanism does not have an equivalent in existing widely deployed protocols. Negative disaggregation is needed only in certain large topologies, namely multi-plane topologies.

RIFT Positive Disaggregation

Earlier we saw that RIFT normally uses default routes for north-bound traffic, which reduces the size of the forwarding tables, but it may cause traffic to be black-holed when a link failure occurs (refer to Figure 4):

Figure 4: Positive Disaggregation in a Two-Level Fabric



Consider a link failure between spine-1 and leaf-3, as shown in Figure 4. When leaf-1 wants to send traffic to leaf-3 and follows its north-bound *Equal-Cost Multi-path Routing* (ECMP) default route, it might select spine-1 as the next hop, which black-holes the traffic.

To avoid such black-holing of traffic, spine-2 and spine-3 each automatically triggers positive disaggregation. Following is the sequence of events from the perspective of spine-2, but the same sequence of events happens at spine-3:

1. Spine-2 has a south-bound route **2.0.0.3/32**, whose next-hop is leaf-3.
2. Spine-2 knows the adjacencies of spine-1 because the leaves reflected the spine-1 node TIE of spine-1 to spine-2. Hence, spine-2 knows that spine-1 does not have an adjacency with leaf-3 and that spine-1 cannot reach **2.0.0.3/32**.
3. Spine-2 automatically initiates positive disaggregation by flooding a positive disaggregation prefix TIE containing prefix **2.0.0.3/32** in the south-bound direction (the blue arrows in Figure 4).
4. Leaf-1 and leaf-2 install the more-specific route prefix **2.0.0.3/32** in their forwarding table. In the end, this route ends up being a two-way ECMP route across spine-2 and spine-3 (but not spine-1). Note that it takes a finite amount of time for the route to reach its full ECMP next-hop set, which may cause transitory in-cast issues (these issues can be addressed with implementation-specific mechanisms).
5. Leaf-1 and leaf-2 still rely on the default route for all other destinations. This route is a three-way ECMP route across all three spines.

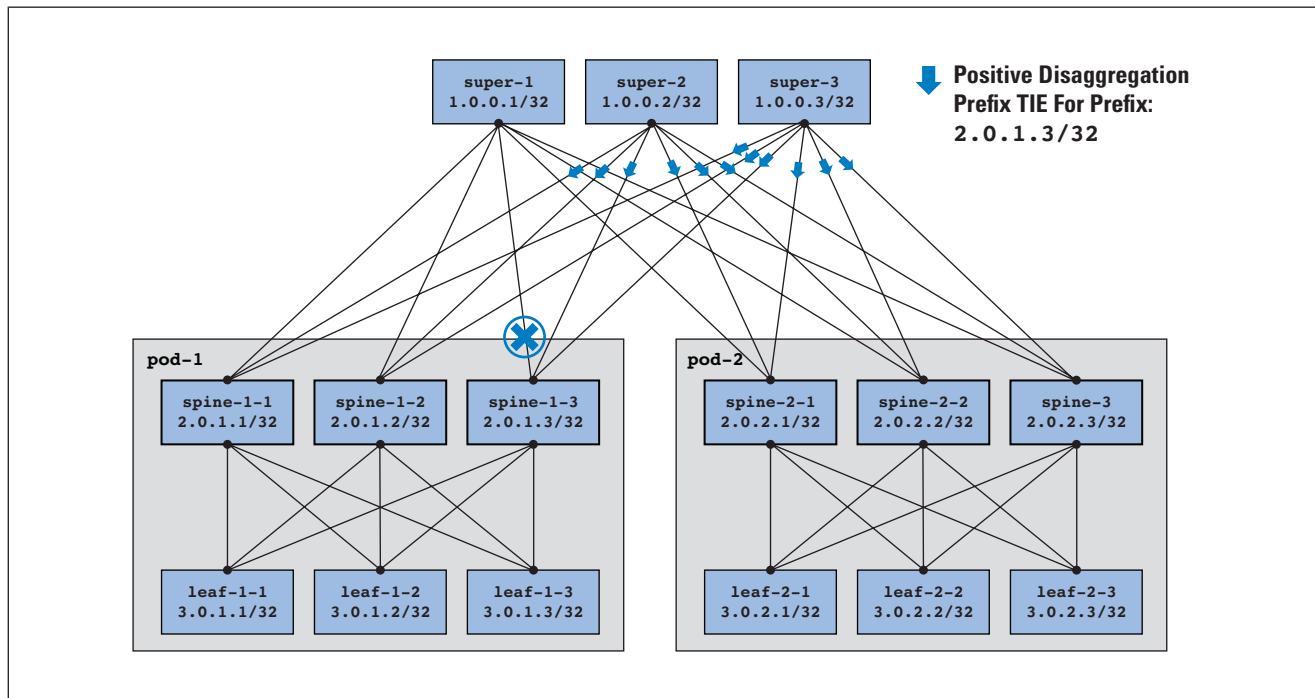
In summary, spine-2 and spine-3 detected that a link was broken in the topology, and they automatically initiated positive disaggregation to “attract” the traffic away from the failed path (spine-1) towards a repair path (themselves).

We now consider positive disaggregation in a more-complex scenario, namely a three-level fabric. In Figure 5, the link from super-1 to spine-1-3 has failed.

Super-2 and super-3 will automatically initiate positive disaggregation for prefix **2.0.1.3/32** because:

1. Super-2 and super-3 have a south-bound route for prefix **2.0.1.3/32** with only one next-hop, namely spine-1-3.
2. Super-2 and super-3 know that super-1 does not have an adjacency with spine-1-3.
3. Super-2 and super-3 conclude that super-1 can no longer reach **2.0.1.3/32**. Hence, they initiate positive disaggregation for that prefix.

Figure 5: Positive Disaggregation Repairs a Single Failure in a Three-level Fabric



However, at this point super-2 and super-3 will not yet initiate positive disaggregation for prefixes 3.0.1.1/32, 3.0.1.2/32, and 3.0.1.3/32 because:

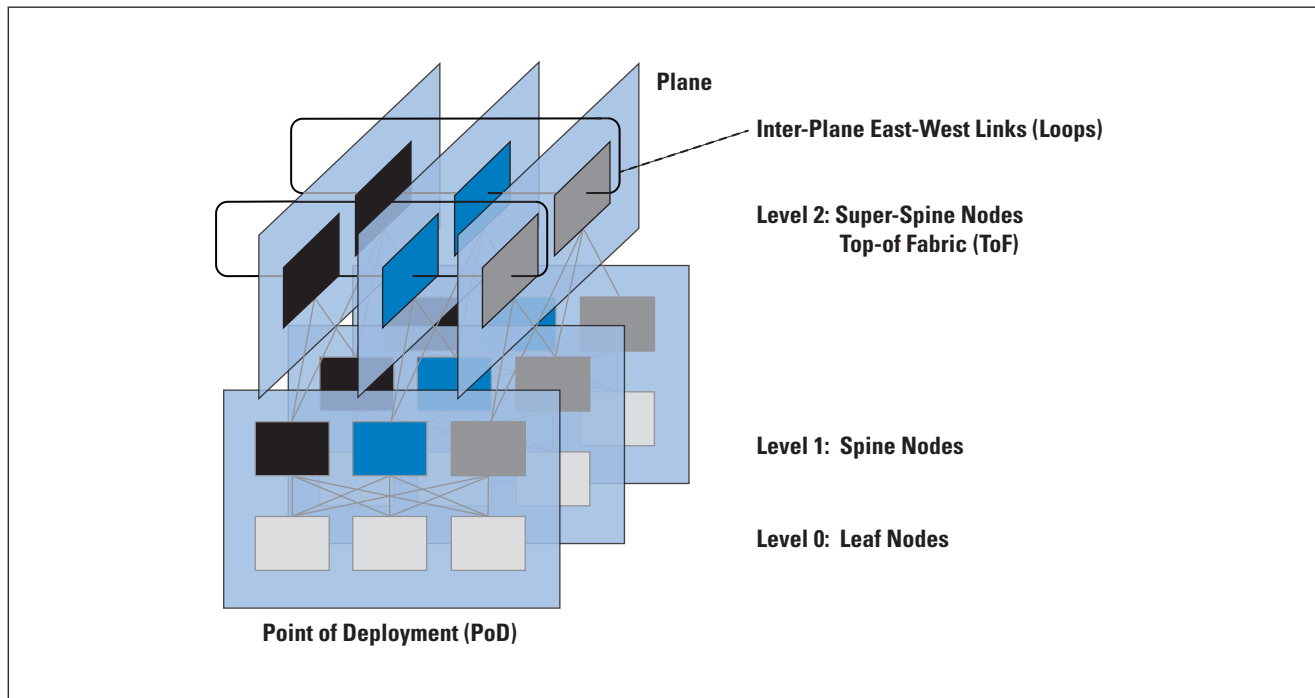
1. Super-2 and super-3 have routes for prefixes 3.0.1.1/32, 3.0.1.2/32, and 3.0.1.3/32, each with three ECMP next-hops, namely spine-1-1, spine-1-2, and spine-1-3.
2. Super-2 and super-3 know that super-1 does not have an adjacency with spine-1-3, but it does still have an adjacency with spine-1-1 and spine-1-2.
3. Super-2 and super-3 conclude that although super-1 can still reach 3.0.1.1/32, 3.0.1.2/32, and 3.0.1.3/32: not through spine-1-3 but still through spine-1-1 and spine-1-2. Hence, they do not initiate positive disaggregation for those prefixes.

We leave it as an exercise for the reader to verify that super-2 and super-3 will initiate disaggregation for 3.0.1.1/32, 3.0.1.2/32, and 3.0.1.3/32 when all links from super-1 to pod-1 are broken.

Multi-Plane Topologies

In the fat-tree topologies that we have considered thus far, every spine is connected to every super-spine. When the network becomes large, you reach a point where the super-spines don't have enough ports to connect to every spine. Such networks often use a multi-plane topology such as the one shown in Figure 6 on the following page.

Figure 6: A Multi-Plane Topology (with East-West Inter-Plane Links)



For now, ignore the loops that connect the super-spines together; they are explained later. In a multi-plane topology, the spines and super-spines are partitioned into planes. In Figure 6, we have a blue plane, a black plane, and a dark grey plane. The super-spines in a plane are connected only to the spines in that same plane.

In such a multi-plane topology, the RIFT positive disaggregation mechanism does not always work because node TIEs are reflected only between super-spines in the same plane. The dark grey super-spines, for example, do not know the adjacencies of the blue or black super-spines. Hence, a super-spine in one plane cannot detect that a super-spine in a different plane has lost connectivity to some set of prefixes.

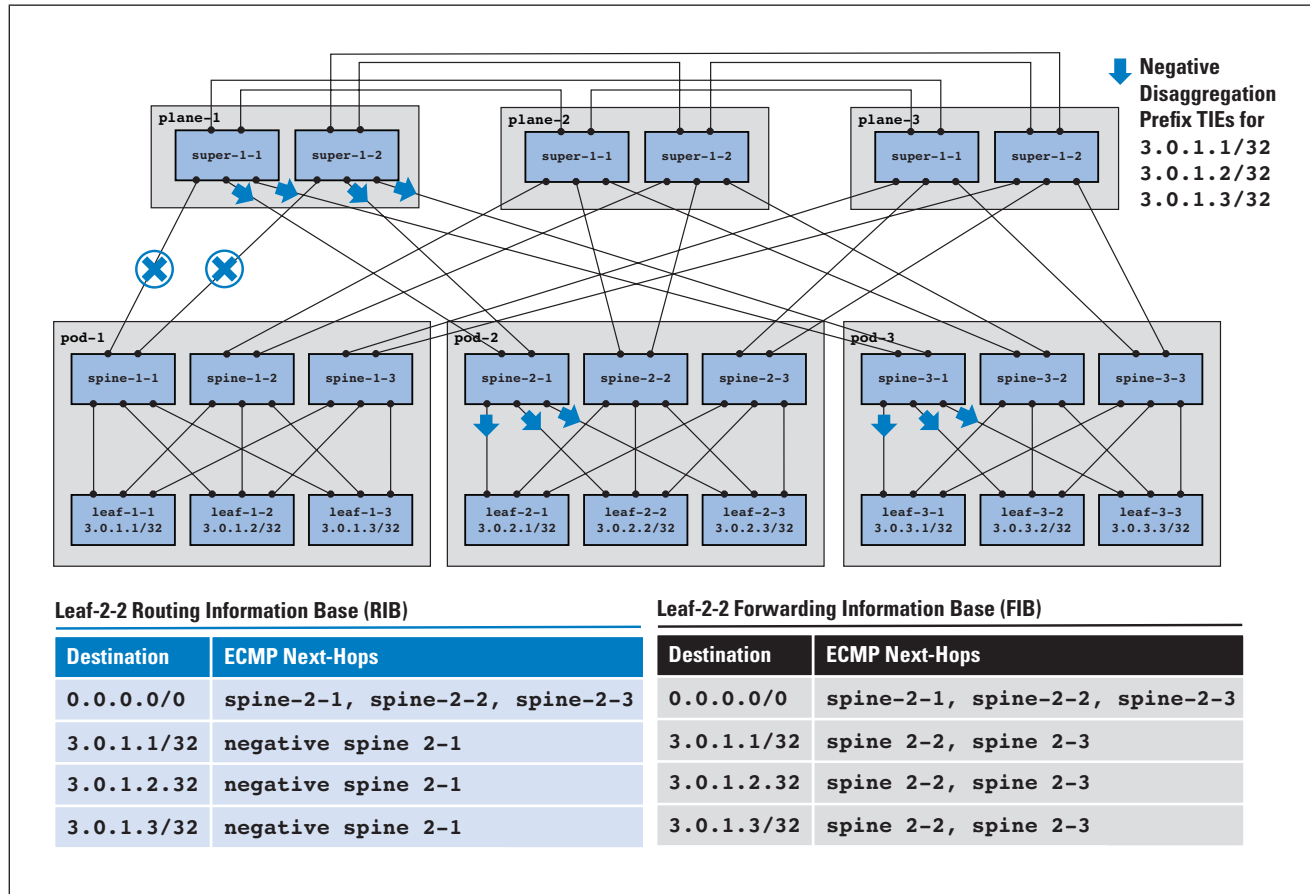
Negative Disaggregation

RIFT uses a different disaggregation mechanism, called *negative disaggregation*, to recover from failures in multi-plane topologies.

To make negative disaggregation work, the super-spines in different planes need to be interconnected using east-west inter-plane links as shown in Figures 6 and 7. These east-west inter-plane links are used only for control-plane traffic, and they do not carry user traffic (so they can be low-bandwidth links).

To understand negative disaggregation, consider the following multi-plane topology:

Figure 7: Negative Disaggregation



Triggering Negative Disaggregation

The super-spines run the south-bound *Shortest Path First* (SPF) algorithm twice:

1. The *normal SPF run* excludes the east-west inter-plane links. The resulting routes are installed in the *Routing Information Base* (RIB) and the *Forwarding Information Base* (FIB).
2. The *special SPF run* includes the east-west inter-plane links. The resulting routes are not installed in the RIB or FIB. If the special SPF run finds any extra reachable prefixes that were not reachable in the normal SPF run, then those extra prefixes are declared to be “fallen leaves,” and they trigger negative disaggregation.

In Figure 7, from the perspective of super-1-1 and super-1-2, the prefixes in pod-1 are fallen leaves because they can be reached only through other planes (that is, using east-west inter-plane links). The special SPF run will find the prefixes in pod-1, but the normal SPF run will not find the prefixes in pod-1.

After a super-spine detects fallen-leaf prefixes, it advertises those prefixes in a negative disaggregation prefix TIE, which is flooded south in the topology.

The super-spine is telling the rest of the network “don’t send any traffic destined to the fallen-leaf prefix to me because I cannot reach it.”

In a sense, negative disaggregation is the opposite of positive disaggregation. In positive disaggregation, the repair path advertises a positive disaggregation route to *attract* the traffic away from the broken path. In negative disaggregation, the broken path advertises a negative disaggregation prefix to *repel* traffic away towards the repair path. The mechanism for choosing the repair path is described in the sections that explain negative next-hop-to-positive next-hop translation.

Propagation of Negative Disaggregation

Unlike positive disaggregation (which is never propagated), negative disaggregation can be recursively propagated southwards. RIFT uses special rules for south-bound flooding of negative disaggregation prefix TIEs: a node propagates a negative disaggregation prefix *only* if it was received from *all* of the parent nodes, meaning that this node does not have any path left to the fallen leaf.

In Figure 7, spine-2-1 has received a negative disaggregation prefix TIE for the prefixes in pod-1 from both of its parent nodes, namely super-1-1 and super-1-2. Hence, spine-2-1 propagates the negative disaggregation prefix TIE further south-bound. The same happens at spine-3-1.

Negative Disaggregation in the RIB

When a node receives a negative disaggregation prefix TIE, it is stored in the LSDB and it takes part in the SPF calculation, just like a normal prefix TIE. However, the resulting route is installed in the RIB using a negative next-hop instead of a positive next-hop.

In Figure 7, leaf-2-2 has a north-bound default route `0.0.0.0/0` with three ECMP next-hops: spine-2-1, spine-2-2, and spine-2-3. These next-hops are normal (that is, positive) next-hops; the traffic will be distributed across all three spines in the pod.

Leaf-2-2 also has north-bound more-specific routes `3.0.1.x/32` (the prefixes in pod-1) with a negative next-hop spine-2-1. A negative next-hop in the RIB is a control-plane construct, meaning “don’t send the traffic to this next-hop.” The intent of this negative next-hop is to avoid sending traffic for `3.0.1.x/32` into plane-1 because plane-1 is disconnected from pod-1.

Note that a negative next-hop is something different from a discard next-hop. A discard next-hop causes traffic to be dropped. A negative next-hop causes traffic to be sent somewhere else using a less-specific route. We will now explain how it works.

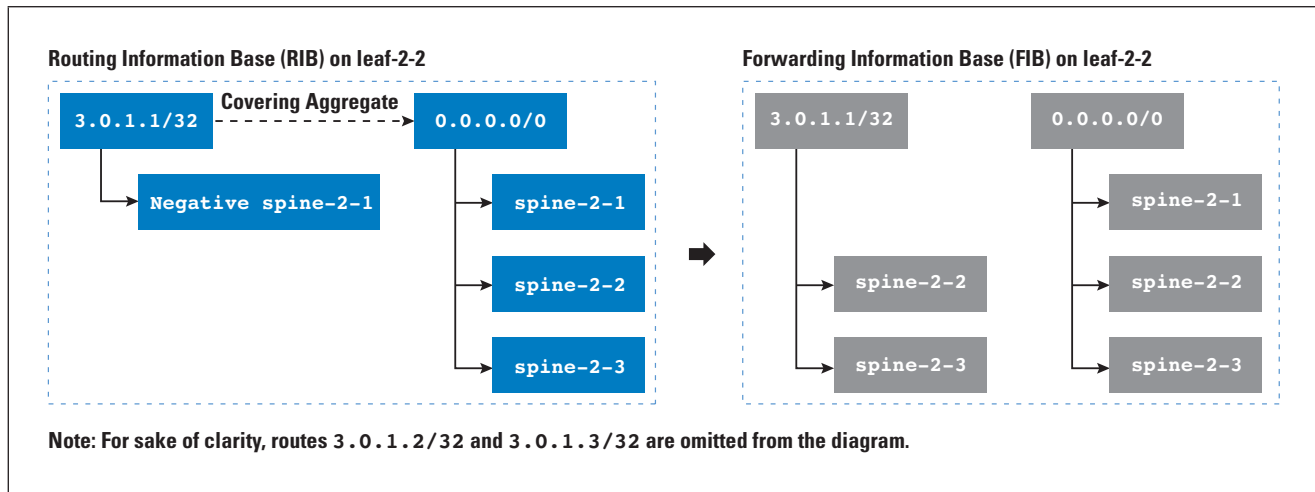
Negative Disaggregation in the FIB

Negative next-hops do not exist in current-generation forwarding hardware; they are a RIFT abstraction that exists only in the control plane and not in the forwarding plane.

When a RIFT route is installed from the RIB into the FIB, the RIB negative next-hop (where not to send the traffic) is translated into positive next-hops (where to send the traffic to instead).

Figure 8 illustrates how this translation works:

Figure 8: Translating Negative Next-Hops in the RIB into Positive Next-Hops in the FIB



What is happening in this simple example is the following:

- We have a route to `3.0.1.1/32`, which has a negative next-hop.
- We find the most specific aggregate route that covers this route, which is the default route, `0.0.0.0/0` in this case.
- We add up the next-hops of routes `3.0.1.1/32` and `0.0.0.0/0`, keeping in mind that a negative and positive next-hop cancel each other out.

We started with a route for `3.0.1.1/32` with negative next-hop `spine-2-1`. We translated the negative next-hop `spine-2-1` into the complementary positive ECMP next-hops `spine-2-2` and `spine-2-3`. These translated next-hops are stored in the FIB.

Further Reading

For more-detailed information about RIFT disaggregation, see Pascal Thubert's slides on negative disaggregation presented at the IETF [10], the RIFT-Python open-source documentation,^[11,12] or my blog post on the topic,^[13] which goes into more detail.

Conclusion

In this article, we have introduced the RIFT protocol and described how RIFT uses automatic aggregation (north-bound default routes) to reduce the size of the routing table.

We have explained how RIFT uses automatic disaggregation to reroute traffic around failed links and nodes. We have further described the two types of disaggregation in RIFT, namely *Positive Disaggregation* and *Negative Disaggregation*.

Positive Disaggregation *attracts* traffic to the repair path by advertising more-specific routes, except that RIFT advertises these routes automatically instead of as a result of manual configuration. Negative disaggregation *repels* traffic away from the broken path. This approach is novel in that it relies on the new concept of a “negative next-hop.” These negative next-hops are translated into normal positive next-hops in the data-plane hardware.

Acknowledgements

I would like to thank Tony Przygienda (the father of RIFT), Pascal Thubert (the father of negative disaggregation), and Melchior Aelmans (one of my co-authors on the RIFT book) for our frequent and deep conversations on the RIFT protocol, and for their review of this article. I would like to thank Mariano Scazzariello and Tommaso Caiazzì (both PhD students at Roma University in Italy) for implementing negative disaggregation in RIFT-Python and their extensive testing of RIFT at scale.^[16]

References

- [1] Tony Przygienda, Alankar Sharma, Pascal Thubert, Bruno Rijsman, and Dmitry Afanasiev, “RIFT: Routing in Fat Trees,” Internet Draft, Work in Progress, **draft-ietf-rift-rift-12**, May 2020
- [2] Yuehua Wei, Zheng Zhang, Dmitry Afanasiev, Tom Verhaeg, Jaroslaw Kowalczyk, and Pascal Thubert, “RIFT Applicability,” Internet Draft, Work in Progress, **draft-ietf-rift-applicability-03**, October 2020.
- [3] RIFT-Python, an open-source implementation of RIFT in Python, Bruno Rijsman and other contributors:
<https://github.com/brunorijsman/rift-python>
- [4] Juniper Networks, “RIFT User Guide for Junos OS,”
https://www.juniper.net/documentation/en_US/junos/information-products/pathway-pages/config-guide-routing/config-guide-routing-rift.html
- [5] Antoni Przygienda and Zhaohui (Jeffrey) Zhang, “Routing in Fat Trees; A New DC Routing Protocol,”
<https://www.slideshare.net/apnic/routing-in-fat-trees>
- [6] Melchior Aelmans, Olivier Vandezande, Bruno Rijsman, Jordan Head, Christan Graf, Leonardo Alberro, Hitesh Mali, and Oliver Steudler, *Day One: Routing In Fat Trees (RIFT), A complete look at the cutting edge protocol*, Juniper Networks Books.
https://www.juniper.net/documentation/en_US/day-one-books/DO_RIFT.pdf

- [7] Russ White and Melchior Aelmans, “Recent Developments in Link State on Data-Center Fabrics,” *The Internet Protocol Journal*, Volume 22, Number 2, September 2020.
- [8] Juniper Networks, “Understanding Route Aggregation,” https://www.juniper.net/documentation/en_US/junos/topics/concept/policy-aggregate-routes.html
- [9] Geoff Huston, “BGP More Specifics: Routing Vandalism or Useful?” Published on the RIPE NCC website: <https://labs.ripe.net/Members/gih/bgp-more-specifics-routing-vandalism-or-useful>
- [10] Pascal Thubert, “Negative Disaggregation,” <https://datatracker.ietf.org/doc/slides-103-rift-negative-disaggregation/>
- [11] Bruno Rijsman, “RIFT-Python Positive Disaggregation Feature Guide,” <https://github.com/brunorijsman/rift-python/blob/master/doc/positive-disaggregation-feature-guide.md>
- [12] Bruno Rijsman, “RIFT-Python Negative Disaggregation Feature Guide,” <https://github.com/brunorijsman/rift-python/blob/master/doc/negative-disaggregation-feature-guide.md>
- [13] Bruno Rijsman Blog, “Automatic Disaggregation in the Routing in Fat Trees (RIFT) Protocol,” <https://hikingandcoding.wordpress.com/2020/07/22/rift-disaggregation/>
- [14] Bruno Rijsman GitHub Page: <https://github.com/brunorijsman>
- [15] Wojciech Kozlowski, Stephanie Wehner, Rodney Van Meter, Bruno Rijsman, Angela Cacciapuoti, Marcello Caleffi, and Shota Nagayama, “Architectural Principles for a Quantum Internet,” Internet Draft, Work in Progress, September 2020, **draft-irtf-qirg-principles-05**
- [16] Tommaso Caiazzì, Mariano Scazzariello, Lorenzo Ariemma, “VFTGen: a Tool to Perform Experiments in Virtual Fat Tree Topologies,” <http://dl.ifip.org/db/conf/im/im2021demo/213179.pdf>

BRUNO RIJSMAN is a software engineer and architect working mainly on networking protocols. Over the past 25 years, he has held technical and leadership roles at network equipment vendors, including Juniper Networks, Verivue, and Lucent Technologies. He currently spends most of his time on open-source projects ^[14], which include RIFT and quantum networking^[15].
E-mail: brunorijsman@gmail.com

Network Functions Virtualization

by William Stallings

Network Functions Virtualization (NFV) originated from discussions among major network operators and carriers about how to improve network operations in the high-volume multimedia era. These discussions resulted in the publication of the original 2012 NFV White Paper by an NFV Industry Specification Group within the *European Telecommunications Standards Institute* (ETSI).^[1] In the white paper, the group listed as the overall objective of NFV the leveraging of standard IT virtualization technology to consolidate many network equipment types onto industry-standard high-volume servers, switches, and storage, which could be located in data centers, network nodes, and at the end-user premises.

The white paper highlights that the source of the need for this new approach is that networks include a large and growing variety of proprietary hardware appliances, leading to the following negative consequences:

- New network services may require additional different types of hardware appliances and finding the space and power to accommodate these boxes is becoming increasingly difficult.
- New hardware means additional capital expenditures.
- After new types of hardware appliances are acquired, operators are faced with the rarity of skills necessary to design, integrate, and operate increasingly complex hardware-based appliances.
- Hardware-based appliances rapidly reach end of life, requiring much of the procure-design-integrate-deploy cycle to be repeated with little or no revenue benefit.
- As technology and services innovation accelerates to meet the demands of an increasingly network-centric IT environment, the need for an increasing variety of hardware platforms inhibits the introduction of new revenue-earning network services.

The NFV approach moves away from the dependence on a variety of hardware platforms to the use of a small number of standardized platform types, with virtualization techniques used to provide the needed network functions. In the white paper, the group expresses the belief that the NFV approach is applicable to any data-plane packet-processing and control-plane function in fixed and mobile network infrastructures.

NFV deployment has become increasingly widespread, being used by telecommunications providers, cloud service providers, and large enterprises, such as in the banking and financial services industry.^[2] Perhaps the main driver for NFV is 5G wireless networks.^[3] NFV is an integral part of 5G and is indeed required by 5G standards.^[4]

Concepts

NFV builds on standard *Virtual Machine* (VM) technologies, extending their use into the networking domain. This departure from traditional approaches to the design, deployment, and management of networking services is significant. NFV decouples network functions, such as *Network Address Translation* (NAT), firewalling, intrusion detection, *Domain Name System* (DNS), and caching, from proprietary hardware appliances so they can run as software on VMs.

Virtual-machine technology enables migration of dedicated application and database servers to *Commercial Off-The-Shelf* (COTS) x86 servers. You can apply the same technology to network-based devices, including:

- *Network Function Devices*: Such as switches, routers, network access points, and deep packet inspectors
- *Network-related Compute Devices*: Such as firewalls, intrusion detection systems, and network management systems
- *Network-attached Storage*: File and database servers attached to the network

In traditional networks, all network elements are enclosed boxes, and hardware cannot be shared. Each device requires additional hardware for increased capacity, but this hardware is idle when the system is running below capacity. With NFV, however, network elements are independent applications that are flexibly deployed on a unified platform comprising standard servers, storage devices, and switches. In this way, software and hardware are decoupled, and capacity for each application is increased or decreased by adding or reducing virtual resources.

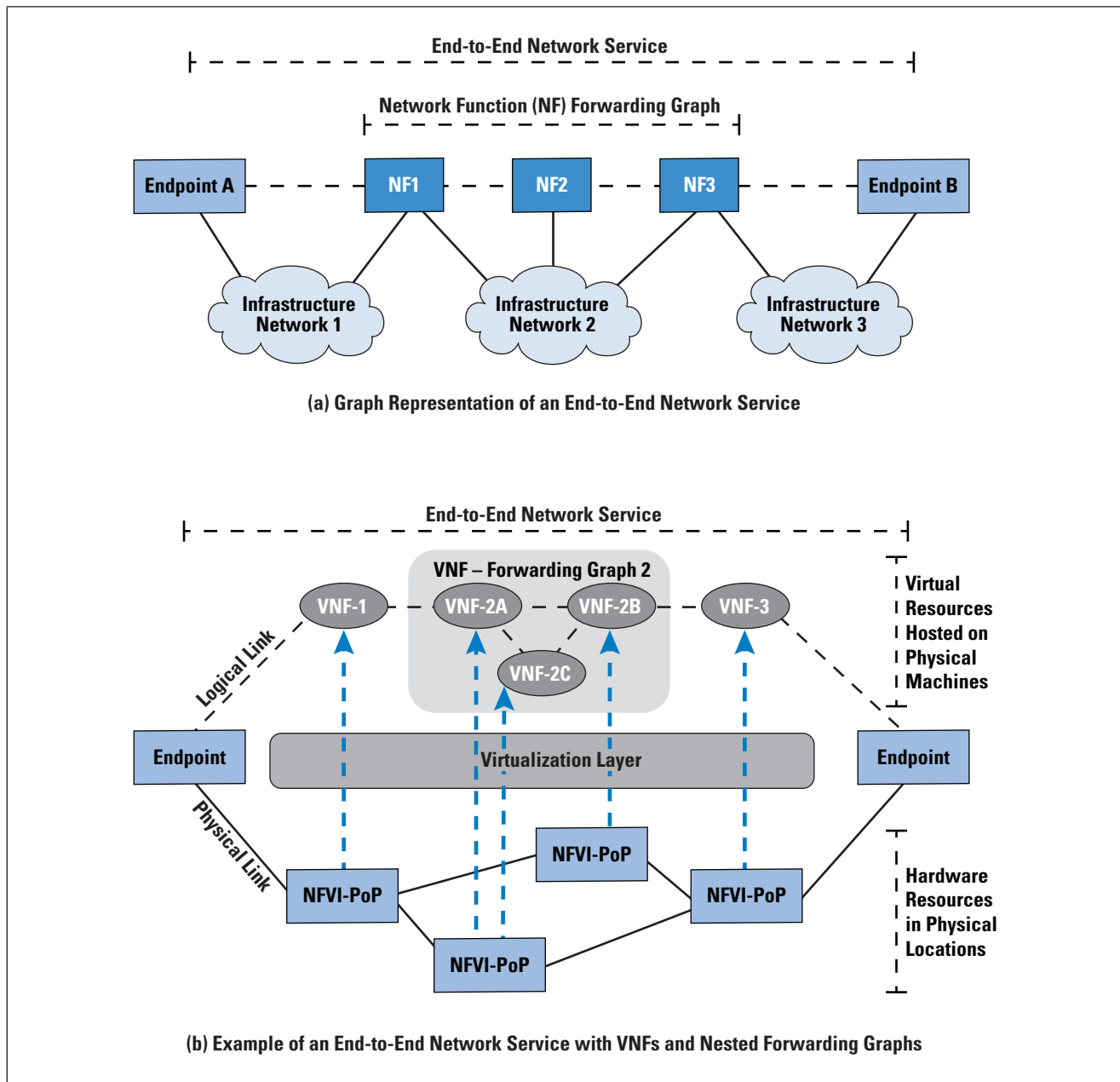
Consider a simple example from the *NFV Architectural Framework* document. Figure 1a shows a physical realization of a network service. At a top level, the network service consists of *endpoints* connected by a forwarding graph of network functional blocks, called *Network Functions* (NFs). Examples of NFs are firewalls, load balancers, and wireless network access points. In the Architectural Framework, NFs are viewed as distinct physical nodes. The endpoints are outside the scope of the NFV specifications and include all customer-owned devices. So, in the figure, endpoint A could be a smartphone and endpoint B a *Content Delivery Network* (CDN) server.

Figure 1a highlights the network functions that are relevant to the service provider and customer. The interconnections among the NFs and endpoints are depicted by dashed lines, representing logical links. These logical links are supported by physical paths through infrastructure networks (wired or wireless).

Figure 1b shows a virtualized network service configuration that could be implemented on the physical configuration of Figure 1a. *Virtual Network Function* (VNF) 1 provides network access for endpoint A, and VNF 2 provides network access for B.

The figure also depicts the case of a nested VNF forwarding graph (VNF-FG-2) constructed from other VNFs (that is, VNF-2A, VNF-2B, and VNF-2C). All of these VNFs run as virtual machines on physical machines, called *Points of Presence* (PoPs). This configuration illustrates several important points. First, VNF-FG-2 consists of three VNFs even though ultimately all of the traffic transiting VNF-FG-2 is between VNF-1 and VNF-3. The reason for this situation is that three separate and distinct network functions are being performed. For example, it may be that some traffic flows need to be subjected to a traffic policing or shaping function, which could be performed by VNF-2C. So, some flows would be routed through VNF-2C while others would bypass this network function.

Figure 1: A Simple NFV Configuration Example



A second observation is that two of the VMs in VNF-FG-2 are hosted on the same physical machine. Because the two VMs perform different functions, they need to be distinct at the virtual resource level but can be supported by the same physical machine. But this setup is not required, and a network management function may at some point decide to migrate one of the VMs to another physical machine, for reasons of performance. This movement is transparent at the virtual resource level.

Principles

As Figure 1 suggests, the VNFs are the building blocks used to create end-to-end network services. Three key NFV principles are involved in creating practical network services:

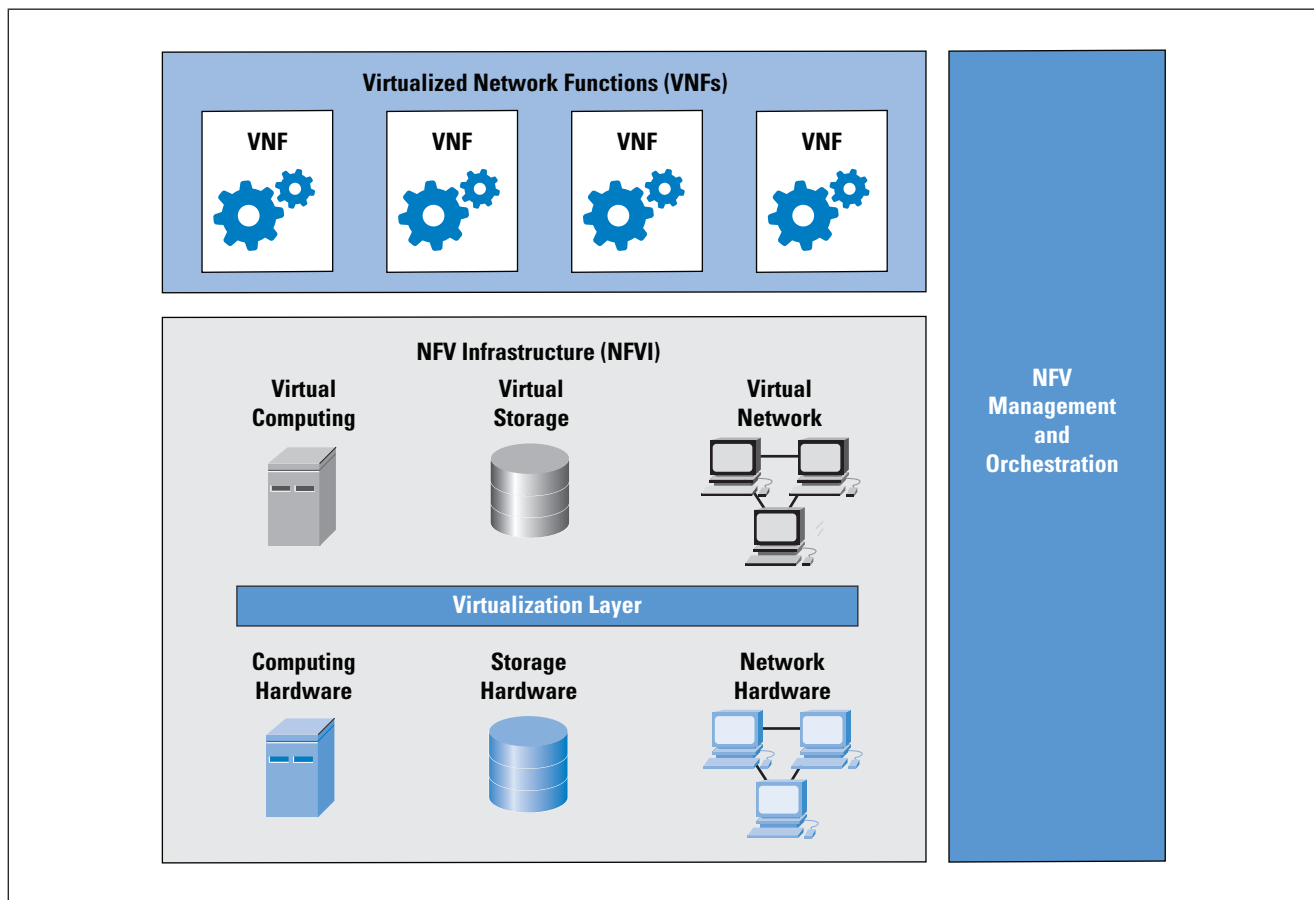
- *Service Chaining*: VNFs are modular and each VNF provides limited functionality on its own. For a given traffic flow within a given application, the service provider steers the flow through multiple VNFs to achieve the desired network functions. This practice is referred to as *service chaining*.
- *Management and Orchestration* (MANO): This feature involves deploying and managing the lifecycle of VNF instances. Examples of functions are VNF instance creation, VNF service chaining, monitoring, relocation, shutdown, and billing. MANO also manages the NFV infrastructure elements.
- *Distributed Architecture*: A VNF may be made up of one or more *VNF Components* (VNFC), each of which implements a subset of the VNF functions. Each VNFC may be deployed in one or multiple instances. These instances may be deployed on separate, distributed hosts in order to provide scalability and redundancy.

Figure 2 shows a high-level view of the NFV framework defined by ISG NFV. This framework supports the implementation of network functions as software-only VNFs. Figure 2 provides an overview of the NFV architecture, which is examined in more detail subsequently.

The NFV framework consists of three domains of operation:

- *Virtualized Network Functions*: These functions are a collection of VNFs, implemented in software, that run over the NFVI.
- *NFV Infrastructure* (NFVI): The NFVI performs a virtualization function on the three main categories of devices in the network service environment: computer devices, storage devices, and network devices.
- *MANO*: This function encompasses the orchestration and lifecycle management of physical and/or software resources that support the infrastructure virtualization and lifecycle management of VNFs. NFV management and orchestration focuses on all virtualization-specific management tasks necessary in the NFV framework.

Figure 2: High-Level NFV Framework



The ISG NFV Architectural Framework document specifies that in the deployment, operation, management, and orchestration of VNFs two types of relations between VNFs are supported:

- *VNF Forwarding Graph (VNF-FG):* Covers the case where network connectivity between VNFs is specified, such as a chain of VNFs on the path to a web server tier (for example, firewall, network address translator, or load balancer).
- *VNF Set:* Covers the case where the connectivity between VNFs is not specified, such as a Web server pool.

NFV Reference Architecture

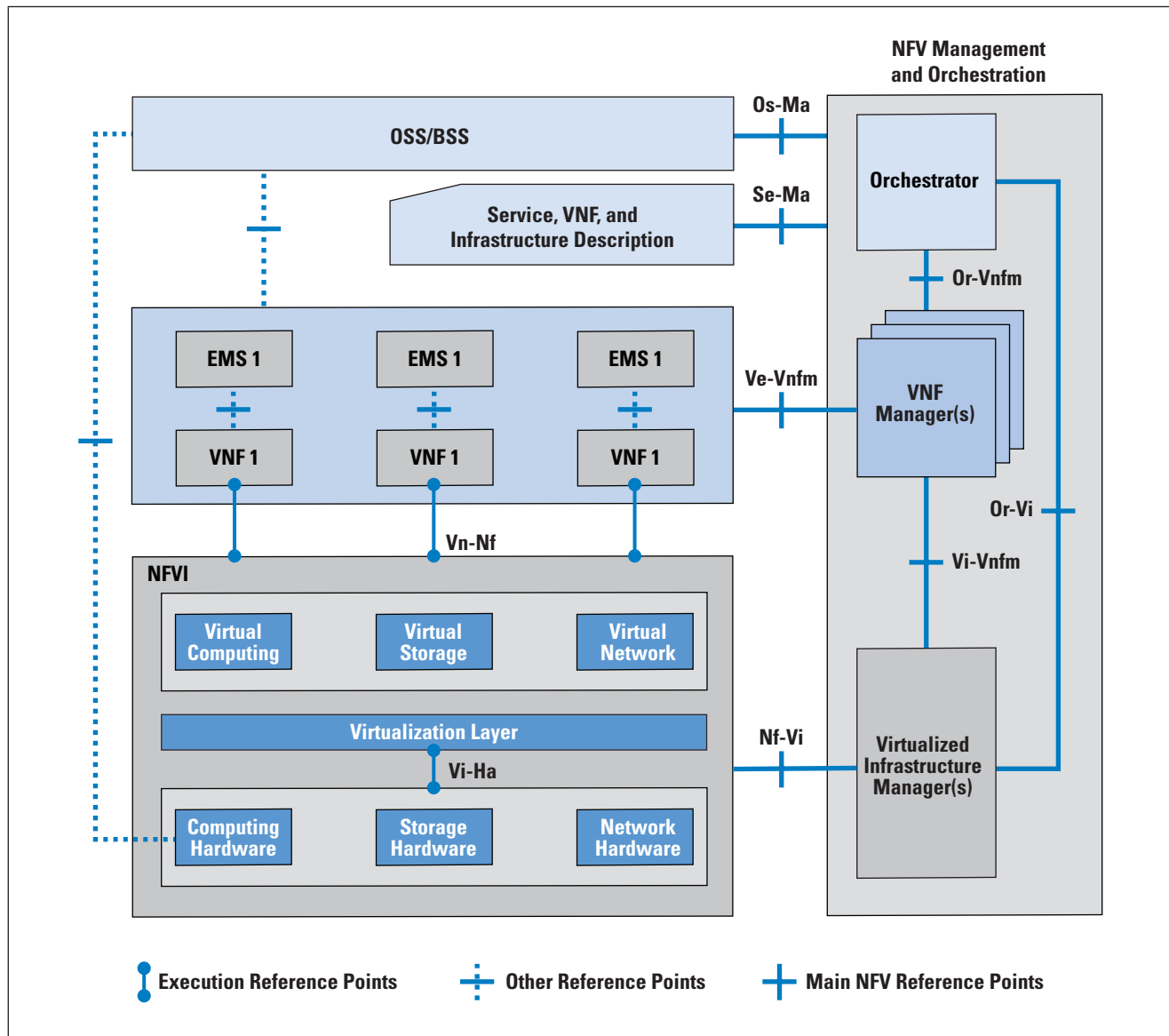
Figure 3 shows a more detailed look at the ISG NFV reference architectural framework.

The architecture consists of four major blocks:

- *NFV Infrastructure (NFVI):* This block comprises the hardware and software resources that create the environment in which VNFs are deployed. NFVI virtualizes physical computing, storage, and networking and places them into resource pools.

- *VNF/EMS*: This collection of VNFs is implemented in software to run on virtual computing, storage, and networking resources, together with a collection of element management systems that manage the VNFs.
- *NFV Management and Orchestration* (NFV-MANO): This framework manages and orchestrates all resources in the NFV environment, including computing, networking, storage, and VM resources
- *Operational and Business Support Systems* (OSS/BSS): The NFV service provider implements this system.

Figure 3: NFV Reference Architectural Framework



It also is useful to view the architecture as consisting of three layers. The NFVI together with the virtualized infrastructure manager provides and manages the virtual resource environment and its underlying physical resources.

The VNF layer provides the software implementation of network functions, together with element management systems and one or more VNF managers. Finally, there is a management, orchestration, and control layer consisting of OSS/BSS and the NFV orchestrator.

NFV Management and Orchestration

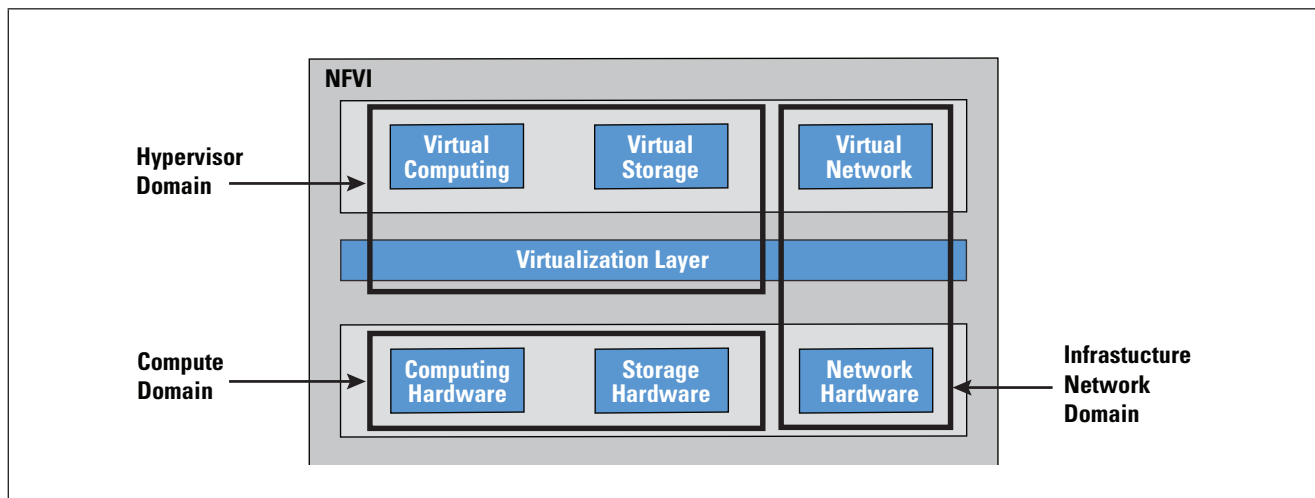
The NFV management and orchestration facility includes the following functional blocks:

- *NFV Orchestrator*: Responsible for installing and configuring new *Network Services* (NS) and VNF packages; NS lifecycle management; global resource management; and validation and authorization of NFVI resource requests.
- *VNF Manager*: Oversees lifecycle management of VNF instances.
- *Virtualized Infrastructure Manager*: Controls and manages the interaction of a VNF with computing, storage, and network resources under its authority, as well as their virtualization.

NFV Infrastructure

The heart of the NFV architecture is a collection of resources and functions known as the *NFV Infrastructure* (NFVI). The NFVI encompasses three domains (Figure 4):

Figure 4: NFV Domains

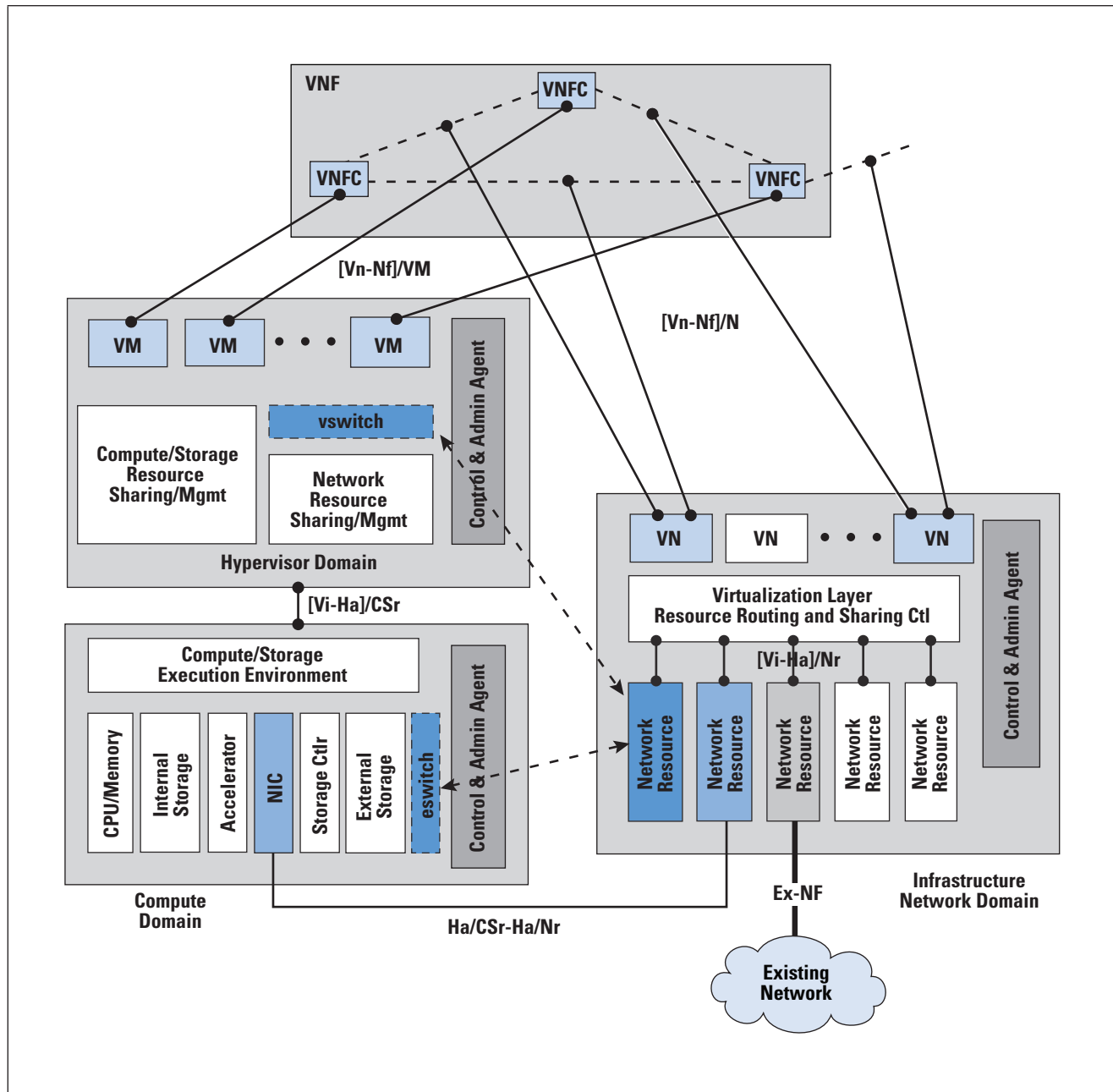


- *Compute Domain*: Provides COTS high-volume servers and storage
- *Hypervisor Domain*: Mediates the resources of the compute domain to the VMs of the software appliances, providing an abstraction of the hardware
- *Infrastructure Network Domain*: Comprises all the generic high-volume switches interconnected into a network that you can configure to supply infrastructure network services.

Logical Structure of NFVI Domains

The ISG NFV standards documents lay out the logical structure of the NFVI domains and their interconnections. The specifics of the actual implementation of the elements of this architecture will evolve in both open-source and proprietary implementation efforts. The NFVI domain logical structure provides a framework for such development and identifies the interfaces between the main components, as shown in Figure 5.

Figure 5: Logical Structure of NFVI Domains



Compute Domain

The principal elements in a typical compute domain may include the following:

- *CPU/Memory*: A COTS processor, with main memory, that executes the VNFC code
- *Internal Storage*: Non-volatile storage housed in the same physical structure as the processor, such as flash memory
- *Accelerator*: Accelerator functions for security, networking, and packet processing
- *External Storage with Storage Controller*: Access to secondary memory devices
- *Network Interface Card (NIC)*: An adapter circuit board installed in a computer to provide a physical connection to a network; it provides the physical interconnection with the infrastructure network domain
- *Control & Admin Agent*: Connects to the *Virtualized Infrastructure Manager (VIM)*; see Figure 2
- *eswitch*: Server-embedded switch; the eswitch function, described in the following paragraph, is implemented in the compute domain, but functionally it forms an integral part of the infrastructure network domain
- *Compute/Storage Execution Environment*: The execution environment that the server or storage device presents to the hypervisor software

To understand the functions of the eswitch, first note that broadly speaking, VNFs deal with two different kinds of workloads: control plane and data plane. Control-plane workloads are concerned with signaling and control-plane protocols such as the *Border Gateway Protocol (BGP)*. Typically, these workloads are more processor- than I/O-intensive, and they do not place a significant burden on the I/O system. Data-plane workloads are concerned with the routing, switching, relaying, or processing of network traffic payloads. Such workloads can require high I/O throughput.

In a virtualized environment such as NFV, all VNF network traffic would go through a virtual switch in the hypervisor domain, which invokes a layer of software between virtualized VNF software and host networking hardware. This situation can create a significant performance penalty. The purpose of the eswitch is to bypass the virtualization software and provide the VNF with a *Direct Memory Access (DMA)* path to the NIC. The eswitch approach accelerates packet processing without any processor overhead.

Hypervisor Domain

The hypervisor domain is a software environment that abstracts hardware and implements services, such as starting a VM, terminating a VM, acting on policies, scaling, live migration, and high availability. The principal elements in the hypervisor domain follow:

- *Compute/storage Resource Sharing/Management*: This service manages these resources and provides virtualized resource access for VMs.
- *Network Resource Sharing/Management*: This service manages these resources and provides virtualized resource access for VMs.
- *Virtual Machine Management and Application Programming Interface (API)*: This service provides the execution environment of a single VNFC instance.
- *Control & Admin Agent*: This agent connects to the *Virtualized Infrastructure Manager (VIM)*; see Figure 3.
- *vswitch*: The vswitch function, described in the following paragraph, is implemented in the hypervisor domain. However, functionally it forms an integral part of the infrastructure network domain.

The vswitch is an Ethernet switch implemented by the hypervisor that interconnects virtual NICs of VMs with each other and with the NIC of the compute node. If two VNFs are on the same physical server, they are connected through the same vswitch. If two VNFs are on different servers, the connection passes through the first vswitch to the NIC and then to an external switch. This switch forwards the connection to the NIC of the desired server. Finally, this NIC forwards it to its internal vswitch and then to the destination VNF.

Infrastructure Network Domain

The *Infrastructure Network Domain (IND)* performs numerous roles. It provides:

- The communication channel between the VNFCs of a distributed VNF
- The communications channel between different VNFs
- The communication channel between VNFs and their orchestration and management
- The communication channel between components of the NFVI and their orchestration and management
- The means of remote deployment of VNFCs
- The means of interconnection with the existing carrier network

An important distinction is to be made between the virtualization function provided by the hypervisor domain and that provided by the infrastructure network domain. Virtualization in the hypervisor domain uses VM technology to create an execution environment for individual VNFCs.

Virtualization in IND creates virtual networks for interconnection of VNFCs with each other and with network nodes outside the NFV ecosystem. These latter types of nodes are called *Physical Network Functions* (PNFs).

Virtualized Network Functions

A VNF is a virtualized implementation of a traditional network function. Table 1 contains examples of functions that could be virtualized.

Table 1: Potential Network Functions to Be Virtualized

Network Element	Function
Switching elements	Broadband network gateways, carrier-grade Network Address Translation (NAT), and routers
Mobile network nodes	Home Location Register/Home Subscriber Server, gateway, GPRS support node, radio network controller, and various node B functions
Customer premises equipment	Home routers and set-top boxes
Tunneling gateway elements	<i>IP Security</i> (IPSec)/SSL virtual private network gateways
Traffic analysis	<i>Deep packet inspection</i> (DPI) and <i>quality of experience</i> (QoE) measurement
Assurance	Service assurance, <i>service-level agreement</i> (SLA) monitoring, and testing and diagnostics
Signaling	Session border controllers and IP Multimedia Subsystem components
Control plane/access functions	<i>Authentication, Authorization, and Accounting</i> (AAA) servers, policy control and charging platforms, and <i>Dynamic Host Configuration Protocol</i> (DHCP) servers
Application optimization	Content-delivery networks, cache servers, load balancers, and accelerators
Security	Firewalls, virus scanners, intrusion detection systems, and spam protection
Support for General Topologies (not just DC fabrics)	No

As discussed earlier, a VNF comprises one or more *VNF Components* (VNFCs). The VNFCs of a single VNF are connected internal to the VNF. This internal structure is not visible to other VNFs or to the VNF user. An important property of VNFs is elasticity, which means being able to perform one or more of the following:

- *Scale up*: Expand capability by adding resources to a single physical machine or virtual machine.
- *Scale down*: Reduce capability by removing resources from a single physical machine or virtual machine.
- *Scale out*: Expand capability by adding additional physical or virtual machines.
- *Scale in*: Reduce capability by removing physical or virtual machines.

Every VNF has an associated elasticity parameter of no elasticity, scale up/down only, scale out/in only, or both scale up/down and scale out/in.

A VNF is scaled by scaling one or more of its constituent VNFCs. Scale out/in is implemented by adding/removing VNFC instance(s) that belong to the VNF being scaled. Scale up/down is implemented by adding/removing resources from existing VNFC instance(s) that belong to the VNF being scaled.

Summary

NFV provides a powerful, vendor-independent approach to implementing complex networks with dynamic demands. NFV builds on well-established technologies, including virtual machines, containers, and virtual networks. With the demand from 5G and cloud service providers, as well as enterprises with large internal networks, NFV is becoming an increasingly widespread technology.

Further Reading

Greater technical detail is available in many survey papers on NFV.^[5, 6, 7, 8] ETSI maintains an NFV web site that includes the ETSI NFV specifications, white papers, tutorials, and a variety of other documents and links (<https://www.etsi.org/technologies/nfv/>). A detailed discussion of the role of NFV in 5G is in [9].

References

- [1] ISG NFV, “Network Functions Virtualization: An Introduction, Benefits, Enablers, Challenges & Call for Action,” ISG NFV White Paper, October 2012.
- [2] Bloomberg L.P., “Network Function Virtualization (NFV) Market Worth \$36.3 Billion by 2024,” January 15, 2020.
<https://www.bloomberg.com/press-releases/2020-01-15/network-function-virtualization-nfv-market-worth-36-3-billion-by-2024-exclusive-report-by-marketsand-markets>
- [3] ISG NFV, “Network Operator Perspectives on NFV Priorities for 5G,” ISG NFV White Paper, February 2017.
- [4] ITU-T, “Requirements of the IMT-2020 Network,” ITU-T Recommendation Y.3101, April 2018.
- [5] Mijumbi, R., et al. “Network Function Virtualization: State-of-the-Art and Research Challenges,” *IEEE Communications Surveys & Tutorials*, First Quarter, 2016.
- [6] Li, Y., and Chen, M. “Software-Defined Network Function Virtualization: A Survey,” *IEEE Access*, December 16, 2016.

- [7] Li, X., and Qian, C. “A Survey of Network Function Placement.” *13th IEEE Annual Consumer Communications & Networking Conference (CCNC)*, 2016.
- [8] Veeraraghavan, M., et al. “Network Function Virtualization: A Survey,” *IEEE Transactions on Communications*, November 2017.
- [9] Stallings, W., *5G Wireless: A Comprehensive Introduction*, ISBN-13: 9780136767145, Pearson Education, Inc., 2021.

WILLIAM STALLINGS is an independent consultant and author of numerous books on security, computer networking, and computer architecture. His latest book is *5G Wireless: A Comprehensive Introduction*, (Pearson, 2021). He maintains a computer science resource site for computer science students and professionals at ComputerScienceStudent.com and is on the editorial board of *Cryptologia*. He has a Ph.D. in computer science from M.I.T. He can be reached at: wllmst@icloud.com

Our Privacy Policy

The *General Data Protection Regulation* (GDPR) is a regulation for data protection and privacy for all individual citizens of the *European Union* (EU) and the *European Economic Area* (EEA). Its implementation in May 2018 led many organizations worldwide to post or update privacy statements regarding how they handle information collected in the course of business. Such statements tend to be long and include carefully crafted legal language. We realize that we may need to provide similar language on our website and in the printed edition, but until such a statement has been developed here is an explanation of how we use any information you have supplied relating to your subscription:

- The mailing list for *The Internet Protocol Journal* (IPJ) is entirely “opt in.” We never have and never will use mailing lists from other organizations for any purpose.
- You may unsubscribe at any time using our online subscription system or by contacting us via e-mail. We will honor any request to remove your name and contact information from our database.
- We will use your contact information only to communicate with you about your subscription; for example, to inform you that a new issue is available, that your subscription needs to be renewed, or that your printed copy has been returned to us as undeliverable by the postal authorities.
- We will never use your contact information for any other purpose or provide the subscription list to any third party other than for the purpose of distributing IPJ by post or by electronic means.
- If you make a donation in support of the journal, your name will be listed on our website and in print unless you tell us otherwise.

Workshop: Measuring Network Quality for End-Users

The *Internet Architecture Board* (IAB) is organizing a virtual workshop, September 14–16, 2021. The Internet in 2021 is quite different from what it was 10 years ago. Today, it is a crucial part of everyone’s daily life. People use the Internet for their social life, for their daily jobs, for routine shopping, and for keeping up with major events. An increasing number of people can access a Gigabit connection, which would be hard to imagine a decade ago. And, thanks to improvements in security, people trust the Internet for both planning their finances and for everyday payments.

At the same time, some aspects of end-user experience have not improved as much. Many users have typical connection latency that remains at decade-old levels. Despite significant reliability improvements in data center environments, end users often see interruptions in service. Transport refinements, such as QUIC, Multipath TCP, and TCP Fast Open are still not fully supported in some networks. Likewise, various advances in the security and privacy of user data are not widely supported, such as encrypted DNS to the local resolver. We believe that one of the major factors behind this lack of progress is the popular perception that throughput is often the sole measure of the quality of Internet connectivity. With such narrow focus, people don’t consider questions such as:

- What is the latency under typical working conditions?
- How reliable is the connectivity across longer time periods?
- Does the network allow the use of a broad range of protocols?
- What services can be run by clients of the network?
- What kind of IPv4, NAT or IPv6 connectivity is offered, and are there firewalls?
- What security mechanisms are available for local services, such as DNS?
- To what degree are the privacy, confidentiality, integrity and authenticity of user communications guarded?

Improving these aspects of network quality will likely depend on measurement and exposing metrics to all involved parties, including to end users in a meaningful way. Such measurements and exposure of the right metrics will allow service providers and network operators to focus on the aspects that impacts the users’ experience most and at the same time empowers users to choose the Internet service that will give them the best experience. The IAB is holding this workshop to convene interested researchers, network operators, and Internet technologists to share their experiences and to collaborate on the steps needed to define properties and metrics with the goal of improving Internet access for all users. The workshop will discuss the following questions:

- What are the fundamental properties of a network that contribute to good user experience?
- What metrics quantify these properties, and how to collect such metrics in a practical way?
- What are the best practices for interpreting those metrics, and incorporating those in a decision-making process?
- What are the best ways to communicate these properties to service providers and network operators?
- How can these metrics be displayed to users in a meaningful way?

We realize that the answers to these questions will vary depending on the different experiences of the participants. For example, a commercial video-streaming platform may prioritize higher throughput and to rely on latency-hiding techniques, while a massive multiplayer online game may prioritize lower jitter, and invest into techniques for graceful degradation of the user experience in case of reduced network capacity. At the same time, researchers from the academia may be looking at properties and metrics that haven't been adopted by the industry at all. Likewise, participants may endorse different methodologies for interpreting the metrics and for making decisions. We are actively looking for identifying such methodologies and for capturing the respective best practices.

While this workshop isn't focusing on the solution space, we are welcoming submissions that dive into particular technologies, to the extent of helping to set the context for the discussion. Comparing the merits of specific solutions, however, is outside of the workshop's scope. Interested participants are invited to submit position papers on the workshop questions. Paper size is not limited, but brevity is encouraged. Interested participants who have published relevant academic papers may submit these as a position paper, optionally with a short abstract. The workshop itself will be a virtual meeting over several sessions, with focused discussion based on the position paper topics received. The logistics for the workshop is as follows:

- Submissions Due: August 2, 2021, midnight AOE (Anywhere On Earth)
- Invitations Issued by: August 16, 2021
- Workshop Dates: September 14–16, 2021 (1400–1800 UTC each day)
- Send Submissions to: **network-quality-workshop-pc@iab.org**

The Program Committee members are Jari Arkko, Olivier Bonaventure, Vint Cerf, Stuart Cheshire, Sam Crawford, Nick Feamster, Jim Gettys, Toke Høiland-Jørgensen, Geoff Huston, Cullen Jennings, Mirja Kuehlewind, Jason Livingood, Matt Mathias, Randall Meyer, Kathleen Nichols, Christoph Paasch, Tommy Pauly, Greg White, and Keith Winstein. The workshop co-chairs are Wes Hardaker, Eugeny Khorov, and Omer Shapira.

Position papers from academia, industry, the open source community and others that focus on measurements, experiences, observations and advice for the future are welcome. Papers that reflect experience based on deployed services are especially welcome. The organizers understand that specific actions taken by operators are unlikely to be discussed in detail, so papers discussing general categories of actions and issues without naming specific technologies, products, or other players in the ecosystem are expected. Papers should not focus on specific protocol solutions. The workshop will be by invitation only. Those wishing to attend should submit a position paper to the address above; it may take the form of an Internet-Draft.

All inputs submitted and considered relevant will be published on the workshop website. The organizers will decide whom to invite based on the submissions received. Sessions will be organized according to content, and not every accepted submission or invited attendee will have an opportunity to present as the intent is to foster discussion and not simply to have a sequence of presentations. Position papers from those not planning to attend the virtual sessions themselves are also encouraged. A workshop report will be published afterwards.

For more information, see:

<https://www.iab.org/activities/workshops/network-quality/>

The APNIC Foundation

The *Asia Pacific Network Information Centre* (APNIC) and the *APNIC Foundation* share a common vision of “a global, open, stable, and secure Internet that serves the entire Asia Pacific community.” Under its charter, the Foundation seeks to “advance education, on a non-profit making basis, in technical, operational and policy matters relating to Internet infrastructure, through undertaking or funding activities in Hong Kong and elsewhere in the Asia and the Pacific region.”

Incorporated in September 2016 and operational in early 2017, the Foundation was first discussed by the APNIC *Executive Council* (EC) in 2014, when it set out to explore a mechanism to support and expand the APNIC Development Program. The EC wanted to do this by raising funds, independent from APNIC membership contributions, to support regional Internet development efforts in the future

Projects and activities funded by the Foundation are designed and managed by APNIC, in collaboration with funding partners interested in Internet development. These activities are implemented by APNIC and our partners, which include a growing group of community trainers and technical advisors, and other like-minded organizations.

The Foundation is guided by an independent Board of Directors—selected by the APNIC EC—that includes recognized and respected experts from the Asia Pacific Internet community.

The Foundation's staff are based in the APNIC office in Brisbane, Australia. The Foundation welcomes support from, and collaboration with, other foundations, agencies and organizations working to develop the Internet in the Asia Pacific.

With more than 13,000 direct and indirect Members in almost every economy of the Asia Pacific, APNIC has spent over 20 years supporting the Internet to serve the region's 3 billion citizens. Many of its 80-plus staff travel regularly in the region to support Members, provide training and technical assistance, or share expertise and information. APNIC also partners with many organizations through MoUs, sponsorships and informally to support the continuing development of the Internet. APNIC's success in partnering and seeking financial support for its activities is founded on five important assets:

- A strong technical focus and regional recognition as a source of best practice and expertise.
- Neutrality and independence from any particular vendors, services, or technologies.
- A non-profit organization with financial strength and transparency.
- Robust regional networks and relationships.
- Long track record of successful management and implementation.

The APNIC Foundation builds on and supports these strengths and APNIC's strong history of success in training and community development.

APNIC development partners have included the Australian *Department of Foreign Affairs and Trade* (DFAT); Canada's *International Development Research Centre* (IDRC); the *Swedish International Development Cooperation Agency* (Sida); the *Japan International Cooperation Agency* (JICA); the World Bank; the United Nations' *International Telecommunications Union* (ITU), the *Internet Corporation for Assigned Names and Numbers* (ICANN), the DotAsia Organization, and the Internet Society. For more information, visit: <https://apnic.foundation/>

EU Launches COVID Certificate

In June 2021, the *European Union* (EU) announced the *Digital Green Certificate*, also known as "Corona Pass" or *Digital COVID Certificate* (DCC), to certify that a European resident has been vaccinated, has recently received a COVID test, or has recovered from the COVID-19 virus. The certificate is used to facilitate travel within EU, and in some cases to allow entrance to some large indoor events. The certificate itself is a QR code, and the majority of its components rely on standards developed by the *Internet Engineering Task Force* (IETF). Éric Vynce explains the details in a blog post here:

<http://evyncke.blogspot.com/2021/06/open-source-standards-at-rescue-to.html>

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Kjetil Aas	Darrell Budic	Holger Durer	Gulf Coast Shots	David Kekar
Fabrizio Accatino	BugWorks	Mark Eanes	Sheryll de Guzman	Stuart Kendrick
Michael Achola	Scott Burleigh	Andrew Edwards	Rex Hale	Robert Kent
Martin Adkins	Chad Burnham	Peter Robert Egli	Jason Hall	Jithin Kesavan
Melchior Aelmans	Jon Harald Bøvre	George Ehlers	James Hamilton	Jubal Kessler
Christopher Affleck	Olivier Cahagne	Peter Eisses	Stephen Hanna	Shan Ali Khan
Scott Aitken	Antoine Camerlo	Torbjörn Eklöv	Martin Hannigan	Nabeel Khatri
Jacobus Akkerhuis	Tracy Camp	Y Ertur	John Hardin	Dae Young Kim
Antonio Cuiat Alario	Ignacio Soto Campos	ERNW GmbH	David Harper	William W. H. Kimandu
Nicola Altan	Fabio Caneparo	ESdatCo	Edward Hauser	John King
Marcelo do Amaral	Roberto Canonico	Steve Esquivel	David Hauweele	Russell Kirk
Matteo D'Ambrosio	David Cardwell	Jay Etchings	Marilyn Hay	Gary Klesk
Selva Anandavel	John Cavanaugh	Mikhail Evstiounin	Headcrafts SRLS	Anthony Klopp
Jens Andersson	Lj Cemerar	Bill Fenner	Hidde van der Heide	Henry Kluge
Danish Ansari	Dave Chapman	Paul Ferguson	Johan Helsingius	Michael Kluk
Finn Arildsen	Stefanos Charchalakakis	Ricardo Ferreira	Robert Hinden	Andrew Koch
Tim Armstrong	Greg Chisholm	Kent Fichtner	Asbjørn Højmark	Ia Kochiashvili
Richard Artes	David Chosrova	Armin Fisslthaler	Damien Holloway	Carsten Koempe
Michael Aschwanden	Marcin Cieslak	Michael Fiumano	Alain Van Hoof	Richard Koene
David Atkins	Lauris Cikovskis	The Flirble Organisation	Edward Hotard	Alexader Kogan
Jac Backus	Guido Coenders	Gary Ford	Bill Huber	Antonin Kral
Jaime Badua	Brad Clark	Jean-Pierre Forcioli	Hagen Hultzsich	Robert Krejčí
Bent Bagger	Narelle Clark	Susan Forney	Kevin Iddles	Mathias Körber
Eric Baker	Horst Clausen	Christopher Forsyth	Mika Ilvesmaki	John Kristoff
Santosh Balagopalan	Joseph Connolly	Andrew Fox	Karsten Iwen	Terje Krogdahl
Benjamin Barkin	Steve Corbató	Craig Fox	David Jaffe	Bobby Krupczak
Wilkins	Brian Courtney	Fausto Franceschini	Ashford Jaggernauth	Murray Kucherauw
Michael Bazarewsky	Beth and Steve Crocker	Valerie Fronczak	Martijn Jansen	Warren Kumari
David Belson	Dave Crocker	Tomislav Futivic	Jozef Janitor	George Kuo
Hidde Beumer	Kevin Croes	Laurence Gagliani	John Jarvis	Dirk Kurfuerst
Pier Paolo Biagi	John Curran	Edward Gallagher	Dennis Jennings	Darrell Lack
Tyson Blanchard	André Danthine	Andrew Gallo	Edward Jennings	Andrew Lamb
John Bigrow	Morgan Davis	Chris Gamboni	Aart Jochem	Richard Lamb
Orvar Ari Bjarnason	Jeff Day	Xosé Bravo Garcia	Brian Johnson	Yan Landriault
Axel Boeger	Julien Dhallenne	Oswaldo Gazzaniga	Curtis Johnson	Edwin Lang
Keith Bogart	Freek Dijkstra	Kevin Gee	Richard Johnson	Sig Lange
Mirko Bonadei	Geert Van Dijk	Greg Giessow	Jim Johnston	Markus Langenmair
Roberto Bonalumi	David Dillow	John Gilbert	Jonatan Jonasson	Fred Langham
Julie Bottorff	Richard Dodsworth	Serge Van Ginderachter	Daniel Jones	Tracy LaQuey Parker
Photography	Ernesto Doelling	Greg Goddard	Gary Jones	Jose Antonio Lazaro
Gerry Boudreaux	Michael Dolan	Tiago Goncalves	Jerry Jones	Lazaro
L de Braal	Eugene Doroniuk	Ron Goodheart	Anders Marius	Rick van Leeuwen
Kevin Breit	Karlheinz Dölger	Octavio Alfageme	Jørgensen	Simon Leinen
Thomas Bridge	Joshua Dreier	Gorostiaga	Amar Joshi	Robert Lewis
Ilia Bromberg	Lutz Drink	Barry Greene	Javier Juan	Christian Liberale
Václav Brožík	Dmitriy Dudko	Jeffrey Greene	David Jump	Martin Lillepuu
Christophe Brun	Andrew Dul	Richard Gregor	Merike Kaeo	Roger Lindholm
Gareth Bryan	Joan Marc Riera	Martijn Groenleer	Andrew Kaiser	Link Light Networks
Stefan Buckmann	Duocastella	Geert Jan de Groot	Christos Karayiannis	Sergio Loreti
Caner Budakoglu	Pedro Duque	Christopher Guemez	Daniel Karrenberg	Eric Louie

Adam Loveless	Maurizio Moroni	David Raistrick	Timothy Schwab	Kerry Thompson
Guillermo a Loyola	Brian Mort	Priyan R Rajeevan	Roger Schwartz	Lorin J Thompson
Hannes Lubich	Soenke Mumm	Balaji Rajendran	SeenThere	Fabrizio Tivano
Dan Lynch	Tariq Mustafa	Paul Rathbone	Scott Seifel	Joseph Toste
Sanya Madan	Stuart Nadin	William Rawlings	Yury Shefer	Rey Tucker
Miroslav Madić	Michel Nakhla	Mujtiba Raza Rizvi	Yaron Sheffer	Sandro Tumini
Alexis Madriz	Mazdak Rajabi Nasab	Bill Reid	Doron Shikmoni	Angelo Turetta
Carl Malamud	Krishna Natarajan	Petr Rejhon	Tj Shumway	Phil Tweedie
Jonathan Maldonado	Naveen Nathan	Robert Remenyi	Jeffrey Sicuranza	Steve Ulrich
Michael Malik	Darryl Newman	Rodrigo Ribeiro	Thorsten Sideboard	Unitek Engineering AG
Tarmo Mamers	Thomas Nikolajsen	Glenn Ricart	Greipur Sigurdsson	John Urbanek
Yogesh Mangar	Paul Nikolich	Justin Richards	Andrew Simmons	Martin Urwaleck
Bill Manning	Travis Northrup	Rafael Riera	Pradeep Singh	Betsy Vanderpool
Harold March	Marijana Novakovic	Mark Risinger	Henry Sinnreich	Surendran Vangadasalam
Vincent Marchand	David Oates	Fernando Robayo	Geoff Sisson	Ramnath Vasudha
Gabriel Marroquin	Ovidiu Obersterescu	Gregory Robinson	Helge Skrivervik	Philip Venables
David Martin	Tim O'Brien	Ron Rockrohr	Terry Slattery	Buddy Venne
Jim Martin	Mike O'Connor	Carlos Rodrigues	Darren Sleeth	Alejandro Vennera
Ruben Tripiana Martin	Mike O'Dell	Magnus Romedahl	Richard Smit	Luca Ventura
Timothy Martin	John O'Neill	Lex Van Roon	Bob Smith	Scott Vermillion
Carles Mateu	Jim Oplotnik	Alessandra Rosi	Courtney Smith	Tom Vest
Juan Jose Marin	Packet Consulting	David Ross	Eric Smith	Dario Vitali
Martinez	Limited	William Ross	Mark Smith	Jeffrey Wagner
Ioan Maxim	Carlos Astor Araujo	Boudhayan	Tim Sneddon	Don Wahl
David Mazel	Palmeira	Roychowdhury	Craig Snell	Michael L Wahrman
Miles McCredie	Alexis Panagopoulos	Carlos Rubio	Job Snijders	Laurence Walker
Brian McCullough	Gaurav Panwar	Rainer Rudigier	Ronald Solano	Randy Watts
Joe McEachern	Manuel Uruena Pascual	Timo Ruiters	Asit Som	Andrew Webster
Alexander McKenzie	Ricardo Patara	RustedMusic	Ignacio Soto Campos	Tim Weil
Jay McMaster	Dipesh Patel	Babak Saberi	Evandro Sousa	Jd Wegner
Mark Mc Nicholas	Alex Parkinson	George Sadowsky	Peter Spekrijse	Westmoreland
Carsten Melberg	Craig Partridge	Scott Sandefur	Thayumanavan Sridhar	Engineering Inc.
Kevin Menezes	Dan Paynter	Sachin Sapkal	Paul Stancik	Rick Wesson
Bart Jan Menkveld	Leif Eric Pedersen	Arturas Satkovskis	Ralf Stempfer	Peter Whimp
Sean Mentzer	Rui Sao Pedro	PS Saunders	Matthew Stenberg	Russ White
William Mills	Juan Pena	Richard Savoy	Adrian Stevens	Jurrien Wijlhuizen
David Millsom	Chris Perkins	John Sayer	Clinton Stevens	Derick Winkworth
Desiree Miloshevic	Michael Petry	Phil Scarr	John Streck	Pindar Wong
Joost van der Minnen	Alexander Peuchert	Gianpaolo	Martin Streule	Phillip Yialeloglou
Thomas Mino	David Phelan	Scassellati	David Strom	Janko Zavernik
Rob Minshall	Derrell Piper	Elizabeth Scheid	Viktor Sudakov	Bernd Zeimetz
Wijnand Modderman	Rob Pirnie	Jeroen Van Ingen	Edward-W. Suor	Muhammad Ziad
Mohammad Moghaddas	Marc Vives Piza	Schenau	Vincent Surillo	Ziayuddin
Roberto Montoya	Jorge Ivan Pincay Ponce	Carsten Scherb	Terence Charles	Tom Zingale
Charles Monson	Victoria Poncini	Ernest Schirmer	Sweetser	Jose Zumalave
Andrea Montefusco	Blahoslav Popela	Philip Schneck	T2Group	Romeo Zwart
Fernando Montenegro	Andrew Potter	Peter Schoo	Roman Tarasov	廖明沂.
Joel Moore	Eduard Llull Pou	Dan Schrenk	David Theese	
John More	Tim Pozar	Richard Schultz	Douglas Thompson	



Follow us on Twitter and Facebook

@protocoljournal



<https://www.facebook.com/newipj>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Internet
Society



Diamond Sponsors

Your logo here!

Ruby Sponsors



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

David Conrad, Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Principal – Internet Technology, Policy, and Advocacy
The Internet Society

Dr. Jun Murai, Founder, WIDE Project
Distinguished Professor, Keio University
Co-Director, Keio University Cyber Civilization Research Center, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

October 2021

Volume 24, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

FROM THE EDITOR

In This Issue

From the Editor	1
Autonomic Networking.....	2
Securing Inter-Domain Routing.....	19
Thank You!	44
Call for Papers.....	46
Supporters and Sponsors	47

According to Wikipedia, *Autonomic Computing* refers to the self-managing characteristics of distributed computing resources, adapting to unpredictable changes while hiding intrinsic complexity to operators and users. The concept has been expanded to computer networks by the *Autonomic Networking Integrated Model and Approach* (ANIMA) working group of the *Internet Engineering Task Force* (IETF). They recently published six *Request For Comments* (RFCs) about autonomic networking. Our first article provides an overview of the ANIMA model and describes these specifications and several usage scenarios in detail.

During the last two weeks of September 2021, several *Voice over IP* (VoIP) providers became the target of a *Distributed Denial of Service* (DDoS) attack. The victims included the provider that I use for my office telephone service, as well as its upstream provider. According to some reports, the attack left several critical institutions without telephone service, including some 911 emergency call centers. As I write this, my service appears to have been restored, but only after a large-scale re-engineering of the network that my provider uses. DDoS mitigation is not an easy task, especially for services that are real-time in nature such as telephone calls. Although I don't expect to learn all of the details of this incident, I do hope that we can cover the topic in more general terms in future articles. If you know any experts on DDoS mitigation, please ask them to get in touch! In the meantime, check out the article entitled "May I ask who's calling, please? A recent rise in VoIP DDoS attacks," which you can find on *The Cloudflare Blog*.

Security has been a recurring theme in this journal. Most of the protocols used in today's Internet were originally designed without comprehensive security in mind, but the IETF has produced security enhancements for many of the core protocols. Securing the routing system itself has proven challenging because it requires wide-spread deployment in order to be effective. In this issue, Geoff Huston presents Part One of a two-part article entitled "A Survey on Securing Inter-Domain Routing."

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

Autonomic Networking Gets Serious

by Michael Behringer, independent;
Carsten Bormann, Universität Bremen TZI;
Brian E. Carpenter, The University of Auckland;
Toerless Eckert, Futurewei;
Jéferson Campos Nobre, UFRGS;
Sheng Jiang, Huawei Technologies;
Yizhou Li, Huawei Technologies; and
Michael C. Richardson, Sandelman Software Works Inc.

In May 2021, six *Request For Comments* (RFCs) about autonomic networking were published [5–10] as a result of the work of the *Autonomic Networking Integrated Model and Approach* (ANIMA) working group of the *Internet Engineering Task Force* (IETF). These RFCs complete the initial charter of that working group, which was started in late 2014 (see [11] for a summary of its inception); however, the first documents to be discussed in the IETF and *Internet Research Task Force* (IRTF) were posted in 2012^[13]. This foundation now allows the industry to build IETF-standardized network solutions for an *Autonomic Networking Infrastructure* (ANI) into every network device.

This article starts with an overview of the reasoning behind autonomic networking and a description of an early usage scenario. It then gives an overview of the newly published specifications and how they will interwork with existing network management, before concluding with several specific use cases.

One way to summarize autonomic networking is “plug and play” for professional networks. It can mean plug and play “for the ISP” or “for the enterprise” or “for industrial networks.” This step is a significant one forward from the well-known idea of plug and play for home networks, which the IETF addresses in the HOMENET working group.

IBM coined the term “autonomic computing” in 2001. The autonomic nervous system acts largely unconsciously and regulates bodily functions such as heart rate. IBM defined autonomic computing as “self-managing distributed computing resources, adapting to unpredictable changes while hiding intrinsic complexity from operators and users.” This definition led naturally to the idea of autonomic networking, which became a topic of discussion and work in the IRTF *Network Management Research Group*. The result was RFCs [1,2], which describe the outline of an envisioned autonomic networking infrastructure and ultimately resulted in the creation of the ANIMA working group. Since then, various aspects of the problem space were addressed in research, and in proprietary implementations by some vendors. But as always, the need is for interoperability, so proprietary methods have to give way to industry standards. This interoperability task is the job of the ANIMA working group.

The goal is self-management of networks, including self-configuration, self-optimization, self-healing, and self-protection (sometimes collectively called *self-X*). Autonomic networking puts operational intelligence into algorithms at the node level, to minimize dependency on human administrators and central management. Autonomic nodes will discover information about the surrounding network and negotiate parameter settings with their neighbors and other nodes. Later, nodes may also have learning and cognitive capability, that is, the ability to self-adapt their decision-making process based on information and knowledge sensed from their environment.

Science fiction? Not really. Distributed routing protocols as introduced with the ARPANET in the 1970s and later in the Internet are at their core autonomic: self-configuring, self-optimizing, and self-healing. Examples include *Open Shortest Path First* (OSPF) and *Intermediate System-to-Intermediate System* (IS-IS). But over the decades, even those protocols have evolved to become provisioning monsters requiring the human configuration of obscure parameters and policies. A whole industry and research discipline for network *Operations, Administration, and Management* (OAM) evolved to define architectures consisting of ever-more-complex layers between the human intent for the service-level objectives of the network (and by implication its protocols) and all the detailed parameters that need to be provisioned consistently and dynamically into each network device whenever there is any change. (As evidence, consider that the IETF alone has published more than 120 *Yet Another Next Generation* (YANG) modules and sub-modules, each of which contains many individual parameters.)

In today's networks, routing and traffic-engineering parameters are almost exclusively implemented through a centralized set of *Software-Defined Networking* (SDN) controller and orchestrator tools configured by human operators. Although a great improvement on older methods, these solutions are still difficult and expensive to build, maintain, validate, predict, secure, and above all to make reliable and resilient. These problems are rarely seen from the outside, except when network services are under oversight of regulatory entities that publish reports of those problems, such as [12].

SDN architectures are also highly proprietary—very often from a single vendor—and they typically require significant customization through programming for any multi-vendor network deployment. They therefore require network owners to not only hire network operators, but also have them become SDN developers. And sometimes, expensive experts have to travel unexpectedly at any hour of the day to fix or update systems. These issues largely arise because of the lack of the automation inside switches and routers that autonomic networking aims to enable.

Nevertheless, these SDN methods are the best option for existing large networks. They are marketed with terms that evolved in the last few years, such as *Zero-Touch Networks*, *Intent-Based Networking*, or *Self-Driving Networks*. In the metaphor of a network being a car, today's networks are like children's pedal cars guided from behind by an attentive parent, whereas ANIMA wants them to be like a self-driving car.

The long-term vision for autonomic networking is broader than the newly published standards. The autonomic networking infrastructure defined in the recent ANIMA RFCs is intended to provide foundational building blocks. These building blocks are meant to fit seamlessly with existing network and SDN/OAM designs and to improve their metrics such as simplicity, reliability, and security. Likewise, the ANI allows designers to more easily embed automation into network devices whenever there is a need. It is worth noting that today, unlike in the past, it is economic to provide enough computing power in network elements to support autonomy.

What Can the Autonomic Networking Infrastructure Do for You?

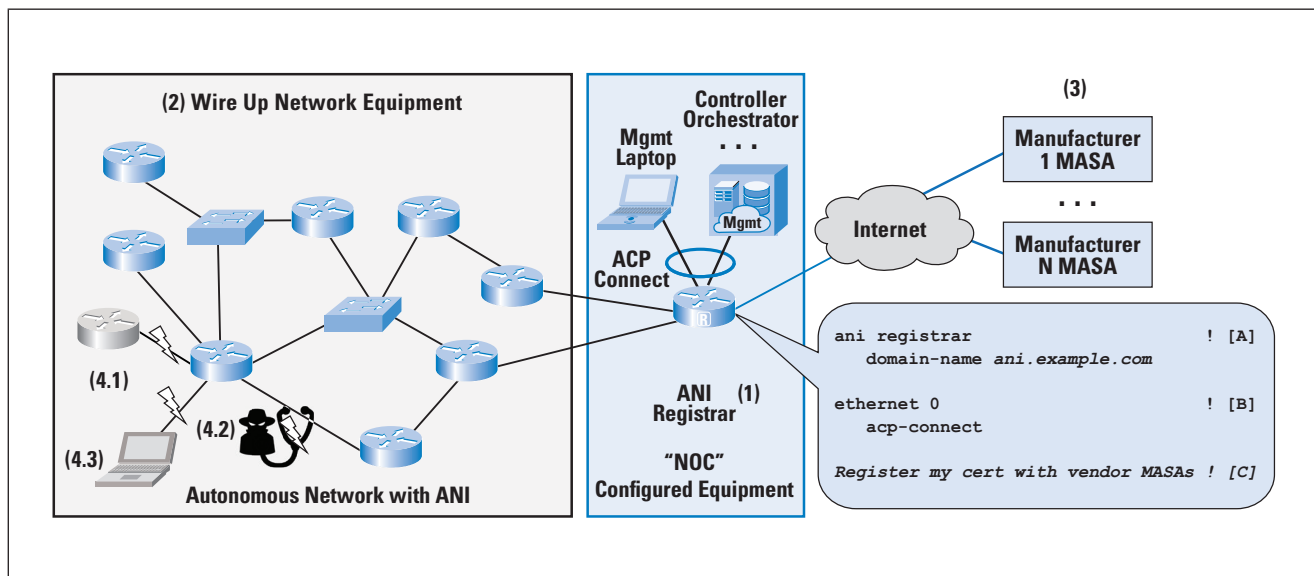
Instead of jumping directly into an explanation of how the ANI works, we first give a simple example of what the operator experience of a simple autonomic network could be.

In Figure 1, an operator wants to deploy a new network of devices such as routers and switches, namely those in the box labeled (2). These devices may be scattered across different physical locations, such as different offices or buildings. The actual reception of the new, factory-fresh equipment, unpacking, and physical attachment to pre-existing links may be performed in different locations by personnel who need to know only how to connect power and network cables accurately.

In contrast, without autonomic solutions, this process is very complex, insecure, and error-prone, and the description of all the challenges experienced would be much longer than this article. The challenges may be as simple as connecting a new device into a wrong Ethernet port, whereas any port would work for autonomic bootstrap. An operator must often ask the local installer to repeatedly power-cycle a device to activate a new or fixed configuration, a process that will be automatic in the ANI. In the worst case, the operator must ask the local installer to perform complex actions such as connecting a laptop to the device and configuring obscure and badly documented features. This situation can result in bizarre telephone interactions such as the operator asking the installer "Please take a photo of that screen and message it to me."

To avoid this situation, many device installations nowadays are done by staging. The device is first shipped to a central location where expert operators pre-configure and secure it on a trusted network, and then it is shipped again to the final deployment location.

Figure 1: An Example Autonomic Network



This process is more secure and more predictable, but it is a lot more expensive and slower. Eliminating the need for staging is hence one of the main advantages of the autonomic bootstrap process.

With the ANI, the operator only sets up a seed router—called the *ANI registrar*—for example in a *Network Operations Center* (NOC). The rest is fully automatic and secure, with local installation of new equipment by less-expert personnel (“plug in power cable, plug data cable into any free Ethernet port”). The NOC setup consists of only three simple steps:

- A. Set up the router labeled (1) as the registrar and assign a name to the ANI.
- B. Configure some local port(s) to provide link-layer access to the ANI, to connect management equipment such as a laptop for manual access or an SDN controller.
- C. Register the certificate of the registrar with the *Manufacturer Authorized Signing Authority* (MASA) services of the vendors whose routers and switches are being used in the new network. (We will soon describe what that registration does.)

Before this seed setup is in place, you may physically interconnect new routers or switches (2), but they will not do anything. When they have connectivity to a configured registrar, they will automatically form an ANI as follows:

Each new ANI device (at that stage called a *pledge*) automatically obtains a connection with the ANI registrar and attempts to enroll, receiving an ANI certificate so that it can participate. But the registrar first needs to prove to the ANI device that it is its “owner.” To do that, the registrar communicates (for example over the Internet) with the MASA of the vendor of that device.

That MASA has the information that this pledge is actually owned by this registrar's network and returns a security voucher that the registrar can present to the pledge, such that the pledge may now trust the registrar and therefore accepts an ANI certificate from the registrar. This process runs completely automatically without any further hand holding or configuration. This part of the ANI is known as *Bootstrapping Remote Secure Key Infrastructure* (BRSKI)^[10] (pronounced "Brewski").

After a new device is enrolled with an ANI certificate, it begins to establish a secure *Autonomic Control Plane* (ACP) connection with all its neighbors, authenticated and authorized mutually by ANI certificates of the device. This step too happens without further hand-holding or configuration. ACP connectivity is always established or re-established between any neighboring ANI routers or switches, regardless of any change in topology. It cannot be affected by faulty operator or SDN configuration of these devices. The goal of the ACP is quite simple: *If there is a physical path to a router or switch, the ACP will automatically provide encrypted and authenticated IPv6 connectivity to it that an operator cannot remove or misconfigure.* This function is exactly the type needed to avoid operational breakdowns such as [12].

Assume all devices were physically connected to each other as shown in Figure 1 and the ANI registrar is connected last (after it was configured). As a result, within minutes, all the devices will have run through BRSKI and set up the ACP. As a result, the network operator now has secure IP connectivity over the ACP from the management laptop and SDN controller to all ANI devices and can configure them manually or through SDN automation using this connectivity. Each ANI device has a permanent and private IP address within the ANI that does not change, even if the device is physically moved in the network.

How is this procedure different from 30-year-old Ethernet technology? Surely you can simply buy a set of inexpensive Ethernet switches, interconnect them, attach a configuration system at one point, and have achieved the same thing?

Indeed, the simplicity of operating Ethernet networks was an inspiration for the ANI, but beyond that, the ANI is fundamentally different. The ANI is above all secure, whereas the default behavior of traditional switches is not. An ANI device can join the ANI only if the operator actually owns it, as cryptographically certified by its manufacturer's MASA, for example via sales records, meaning that a stolen device cannot be enrolled for the ANI in another network. It also means that a device not belonging to this network operator (4.1) cannot be enrolled in this ANI network to launch an attack. To be clear, the operator has not relinquished any control or authority to the manufacturer by this process; only the operator decides which devices may attach to the network and what they may or may not do. The manufacturer's only role is to certify that each device is genuine.

All ACP traffic is encrypted hop-by-hop; therefore, an attacker cannot snoop or spoof any management traffic that uses the ACP, including any legacy unencrypted management protocol (4.2).

Lastly, ANI devices, even after having formed the ACP, are still unconfigured, ideally meaning that they should behave like current unconfigured routers: there is nothing running that could provide undesirable network connectivity to any hosts that attach, like some insecure or malicious laptop (4.3). Such an attached device would get no connectivity whatsoever. As a result, there is never a window of opportunity for attackers to impair unprotected equipment. Instead, the NOC has all the time it needs to remotely provision the devices. In later stages, such provisioning will occur autonomically, as we shall see.

Compared to many other zero-touch solutions, the ANI does not focus only on so-called *day-0/day-1* behavior up until the network is operational. Instead its services last through the whole life cycle. The ANI provides automated certificate renewal for all ANI devices to maintain and refresh its security model. The ACP protects any network OAM traffic that uses it. By its use of hop-by-hop encryption it also continuously protects the whole network and attached OAM equipment from traffic injection or spoofing attacks.

The use of the MASA service is one of the crucial benefits of the ANI process to enable reliable and secure device deployment without prior staging. Without a MASA, if an unconfigured device is connected to an unintended or hostile network, systems that use its default credentials can easily “kidnap” it. Furthermore, an attacker could then intercept the enrollment process in order to gain access to the whole network. For a network connection to become hostile, it is often sufficient for some virus-impaired device (such as a PC) to be on the same LAN or for the attacker to have impaired other network services such as the *Domain Name System* (DNS). Using a MASA to restrict access to cryptographically authorized devices closes off this avenue of attack.

Nevertheless, the MASA concept has raised concerns over the extent of control or observation by the manufacturer. In fact, the MASA can do neither. It can only generate cryptographic vouchers to inform the device and the person who owns it, thereby precluding configuration by anyone else. Manufacturers can operationalize this service in many ways, according to their customers’ requirements. The workflow described previously, where the owner communicates with the MASA during the enrollment of the device into its owner’s network, is just the simplest option for many owners because it offloads the difficult steps onto the manufacturer.

When a distinguished but elderly scientist states that something is possible, he is almost certainly right. When he states that something is impossible, he is very probably wrong.^[14]

Technical Outline of the ANIMA Model

As always in network management, literally thousands or millions of details cannot be standardized, or even described centrally. What we can do is define a model, a platform, and a toolkit, just as the *Simple Network Management Protocol* (SNMP) and the *Network Configuration Protocol* (NETCONF) have done in the past.

The main terminology we will use is the following. More details about these terms is available in RFC 7575^[1] and RFC 8993^[8]:

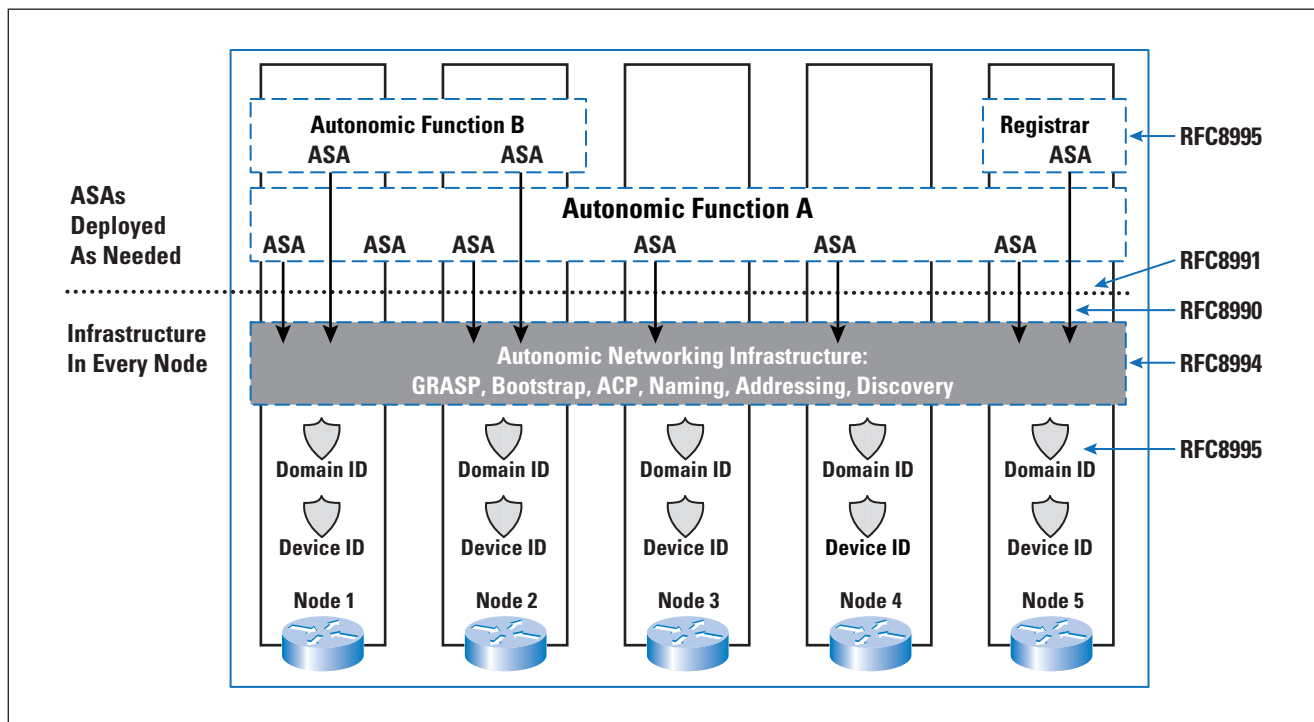
- *Autonomic Function*: A specific self-managing feature or function
- *Autonomic Service Agent* (ASA): An agent that implements an autonomic function, in part (for a distributed function) or whole
- *Autonomic Node*: A node that embodies autonomic functions
- *Autonomic Control Plane* (ACP): A self-configuring, fully secure, virtual network used for all autonomic messaging

The main items in the model follow:

- Bootstrapping and trust infrastructure^[10]. This item covers how nodes are authenticated and securely admitted to an autonomic network, and how they establish mutual trust.
- *Secure Autonomic Control Plane* (ACP)^[9]. This part is an automatically constructed and encrypted virtual network that contains only authenticated nodes that rightfully belong to a particular autonomic domain.
- Discovery for autonomic nodes: This item is a mechanism by which nodes attached to the ACP can discover each other. In practice, discovery occurs at a finer grain than nodes, because it really operates at the level of the capabilities and objectives of a node.
- Negotiation and synchronization for autonomic nodes: After nodes have discovered each other, they can synchronize data between themselves, or actively negotiate parameters and resources.
- Autonomic functions operate by negotiating and synchronizing data with their peers in other nodes, and by directly configuring manageable devices in their own scope.
- Discovery, synchronization, and negotiation proceed by use of the *GeneRic Autonomic Signaling Protocol* (GRASP)^[5].
- *Autonomic Service Agents* (ASAs) are composed of one or more autonomic functions, typically using GRASP via an *Application Programming Interface* (API)^[6].
- Centrally defined policy or configuration rules may be obtained by an ASA via GRASP synchronization, or if appropriate by conventional methods such as an interface to NETCONF or *Domain Name System Service Discovery* (DNS-SD).

Figure 2 shows an outline of the model as a whole, described in detail in RFC 8993^[8].

Figure 2: Layered Model of Network with Autonomic Functions



The only way of discovering the limits of the possible is to venture a little way past them into the impossible.^[14]

Some Details of Self-configuring Security

ANIMA does not attempt a monolithic bootstrap of a network from a predefined configuration. Instead, it proceeds step-by-step, and security comes first. The first stage of creating a secure autonomic control plane is bootstrapping a suitable key infrastructure that covers all the nodes that will constitute the ACP. This process is done, as previously described, by BRSKI.^[10] The process uses manufacturer-installed X.509 certificates (in IEEE 802.1AR IDevID format), in combination with a MASA. The network administrator decides which devices are authorized to join the network (for example, by serial number), but relies on the manufacturer to validate the certificate of each device whenever the device attempts to join the network via a local “join proxy.” These proxies all use a single “domain registrar” node that mediates the authorizing service. The join proxies themselves join the network by the same process; a GRASP mechanism is used for joining nodes (known as *pledges*) to find proxies, and for proxies to find each other and the registrar. Only the registrar needs to be configured in advance.

The ACP forms itself among pledges as soon as the pledges have completed their BRSKI enrollment. It is best described as a *Virtual Routing and Forwarding* (VRF) instance. It is based on a virtual router at each node, consisting of a separate IPv6 forwarding table to which the virtual interfaces of the ACP are attached, and an associated IPv6 routing table separate from the data plane.

Packet transmission is visible only as IPv6 link-local packets, encapsulating the autonomically created overlay network. This choice was made to ensure that there is no dependency on any pre-existing data plane (either IPv4 or IPv6), because autonomic functions must be able to operate *even if the normal data plane and normal routing are broken*. Even then, the ACP provides a secure channel to reach each node for (re-)configuration, without requiring a physically isolated console port. To start the ACP, all that is required is for each node to create its own IPv6 link-local address on each physical interface, as any modern network device does by default. The VRF consists of point-to-point IPv6 links and is secured using *Internet Protocol Security* (IPsec) with *Internet Key Exchange Protocol Version 2* (IKEv2) or *Datagram Transport Layer Security* (DTLS). From the viewpoint of autonomic service agents, the ACP uses an automatically generated IPv6 Unique Local Address prefix, and it uses *Routing Protocol for Low-Power and Lossy Networks* (RPL) internally. Like BRSKI, the ACP bootstraps itself, starting with a GRASP-based discovery process.

The security that the ANI itself requires is a simple but effective based “group-walled-garden” model for *Private Key Infrastructure* (PKI). It provides strong protection against intruders because of its certificate-based model with automated renewals. It also provides for simple ejection of impaired nodes through certificate revocation, certificate status verification, or short-lived certificates. Further levels of security are easily added when necessary. For example, the ANI itself already uses the common certificate-derived role-based security that distinguishes registrars from other nodes, so that no arbitrary impaired node can overtake the domain by acting as a fake registrar. You can expand such role-based security to other crucial roles in autonomic functions.

Of course, it would be naive to assume that, even with this key infrastructure and encrypted network, no malicious device, code, or user will ever penetrate the autonomic system. A malicious ASA could, for example, attempt a *Denial of Service* (DoS) attack within the ACP. The ANI platform provides services such as authentication, confidentiality, credential management, connectivity, and discovery to ASAs. An interesting analogy is *Transport Layer Security* (TLS), which provides authentication and confidentiality to web services. However, TLS cannot prevent the web services themselves from being untrustworthy, for example by breaking expectations of confidentiality by selling user data. In the same way, ASAs need to be intrinsically trustworthy on their own, regardless of whether they use the ANI. All legitimate ASAs should be designed to take appropriate precautions, and a watchdog ASA could be implemented to detect suspicious activity.

After the secure control plane has configured itself, the next stage is to bootstrap connectivity for network management. When this connectivity is achieved, conventional mechanisms (such as an SDN controller) can already reliably and securely reach remote nodes and configure them safely without risk of cutting themselves off.

In addition, fully autonomic management mechanisms (that is, ASAs) can start up. To understand how this process works, we first need to add more details about the GRASP protocol.

GRASP

The *GeneRic Autonomic Signaling Protocol* (GRASP)^[5] is used for signaling between ASAs, including special-purpose mini-ASAs that support BRSKI (discovery of join proxies and the domain registrar) and ACP creation (discovery of ACP neighbors). Readers will notice that these operations must take place *before* ACP security is in place, so they use a highly restricted subset of GRASP that is limited to specific link-local operations.

After that, GRASP runs over the ACP to guarantee security, so there are no restrictions on allowed operations and any two ASAs in the local domain may trust and communicate with each other. GRASP provides discovery, flooding, synchronization, and negotiation mechanisms for the objectives that ASAs support.

Rather than being a traditional type-length-value protocol, GRASP messages use *Concise Binary Object Representation* (CBOR), which provides an extensible data model derived from *JavaScript Object Notation* (JSON), but with a simple and efficient binary encoding. The flexibility of CBOR enables GRASP to accommodate a very wide range of data types, with protocol elements often mapping directly into various high-level language representations.

The word “objective” has a special meaning in GRASP. It is a data structure whose main contents are a *name* and a *value*. An objective occurs in three contexts: *discovery*, *negotiation*, and *synchronization*. A single ASA may support multiple independent objectives.

The *name* of an objective is simply a unique string describing its purpose.

The *value* consists of a single configurable parameter or a set of parameters of some kind. The parameter(s) apply to a specific service, function, or action. They may in principle be anything that can be set to a specific logical, numerical, or string value, or a more complex data structure. Basically, an objective is defined in the way that best suits its application; that is the great advantage of CBOR encoding. If desired, for example, the *value* of an objective could be expressed in the JSON data model. When an objective is shared between ASAs by flooding, synchronization, or negotiation, each ASA will maintain its own copy of the objective and its latest value.

GRASP messages allow for *discovery* of an ASA that handles a given objective name; *flooding* a given objective to all ACP nodes (the simplest form of synchronization); *synchronization* of the value of a given objective between two peer ASAs; and *negotiation* of the value of a given objective with a peer ASA.

An API for GRASP has been defined^[6] and implemented as part of a *Python* 3 prototype, making it very easy to implement demonstration ASAs in *Python*. A partial GRASP implementation has also been made as part of an ACP implementation in the *Rust* language.

Talking to the NOC

As noted previously, a key requirement for the success of ANIMA is smooth integration with existing network management tools and in particular with NOCs. To this end, an integration mechanism has been documented.^[4] The simplest approach is for trusted edge devices in the ACP to “leak” the (otherwise encrypted) ACP natively to certain network management hosts, presumed to be well secured. These edge devices would act as default routers to those management hosts and provide them with IPv6 connectivity into the ACP. A more complex approach would allow the management hosts simultaneous connectivity into the ACP and the traditional data plane.

A related issue is that if the NOC uses *DNS Service Discovery* (DNS-SD) to announce management services to managed nodes, these announcements will not be automatically available in the ACP, which for security reasons will not have routed access to the data plane where the DNS is available. This situation again can be solved by a trusted edge device that obtains service information from DNS-SD and redistributes it within the ACP, possibly by the GRASP flooding mechanism. For example, the information for a service named *syslog* could be flooded in a GRASP objective named *SRV.syslog*. Here, the flexibility of CBOR encoding is of great value because a JSON-like structure of service data is common.

Extending that point, since GRASP easily conveys JSON (or practically any other format), it is possible to integrate ASAs communicating via GRASP into almost any part of an existing network management system. For example, an ASA acting as a NETCONF client could retrieve YANG documents from a NOC database via GRASP and the ACP.

Autonomic Function Example 1: Address Management

A use case that has been fully defined is a GRASP-based mechanism for managing and assigning IP address prefixes.^[7] Firstly, we define two GRASP objectives for IPv4 or IPv6 prefix management at the edge of large-scale *Internet Service Provider* (ISP) networks. The first objective can be represented thus (in a simplified form):

```
["PrefixManager", [IP_version, prefix_length, prefix]]
```

and the second as:

```
["PrefixManager.Params", parameter_info]
```

The first objective will be used in GRASP negotiations between two “prefix manager” ASAs in nodes that need to delegate address space to subsidiary routers (using standard IPv6 prefix delegation), when one node is short of spare prefixes and the other one has an adequate pool of unused prefixes. If negotiation succeeds, prefixes will be transferred from the pool of one ASA pool to the other ASA pool. If negotiation fails, the ASA that is short of prefixes will use GRASP discovery to find another ASA that can help it. Each participating ASA will require persistent storage to manage its own address pool and to survive power outages or other failures such as network partitions. This feature will completely obviate any need for human management of an ISP’s distributed pool of prefixes, beyond initially configuring the maximum pool in one place.

The second objective may be flooded to all “prefix manager” ASAs to convey relevant policy, which can be enforced during prefix delegation by individual agents. For example, if the flooded parameter information is as follows:

```
[
  [{"role", "A"}, {"prefix_length", 34}],
  [{"role", "B"}, {"prefix_length", 44}],
  [{"role", "C"}, {"prefix_length", 56}]
]
```

...it would mean that devices of type A are allowed to receive IPv6 prefixes of length 34 bits, and so on.

You could use this mechanism in a variety of ways. One use case is where the three roles previously discussed correspond to three functions in an IP Radio Access Network: Radio Network Controller Site Gateways, Aggregation Site Gateways, and Cell Site Gateways. These devices will determine their own roles, and then select the prefix length they are allowed to request and offer to each other accordingly. Only central actions are to define the policy to be flooded out and to assign the operator’s total address space to a single device that will progressively delegate it to gateways that request prefixes.

This example illustrates that GRASP’s use of CBOR and its easy representation of JSON-like formats gives it great expressiveness and flexibility. While much work remains to be done on individual autonomic functions, the ANI and GRASP provide a solid and flexible foundation for further development.

Autonomic Function Example 2: Automating IP Multicast

One common interesting challenge for writing distributed autonomic service agents is solving problems that require decisions about the network topology—in a distributed fashion.

A simple example is automating deployment of a service such as IP Multicast, which needs to determine a small set of designated rendezvous routers, where a key requirement is their location balanced between the center and the edges of a network.

Using the ANI and GRASP, it is practical to build such distributed algorithms, for example using common criteria, such as calculation of one's own average path length as an indicator of centrality, and then running a distributed election algorithm that accounts for this and other criteria such as node performance and speed of attachment links to elect a few top contenders for the role, which then auto-configure the service and their precedence in it.

Autonomic Function Example 3: Automatic Protocol Security

We will end by considering an important early operational role for distributed autonomic behavior. That could start soon with very pragmatic incremental in-network automation, perhaps developed by operators as simple scripts in a scripting language such as *Python* or *Tcl* that can run locally on routers.

Consider an existing network where basic services are already running, for example, IPv4 and/or IPv6 addressing and routing. A software upgrade to the routers that adds support for the ANI could be installed, without affecting any of the pre-existing configuration and services. One of the most desirable services is protocol security, for example in routing protocols such as OSPF, IS-IS, and many others.

Most protocols have their own security mechanism and/or keying material requirements. However, security is often not configured because there is no automated key management, including key roll-over and revocation. Without good automation of key management, either networks fail to enable protocol security, or operators set up a single, network-wide password that is never changed. With the ANI, automation of such functions becomes much simpler, by using GRASP, running securely inside the ACP.

With this information in mind, you could easily write a *Python* or *Tcl* script using the GRASP API to auto-configure routing protocol security:

- *Discover* ANI neighbors on links that use the same routing protocol.
- *Generate* a random key.
- *Negotiate* the key with a neighbor.
- *Configure* a routing protocol key locally on the router.
- Periodically wake up, renegotiate, and configure a new key.
- Take suitable action if a neighbor disappears or re-appears.

Some protocols may not even have security included in the protocol itself, for example *Protocol Independent Multicast* (PIM). Instead, you need to secure packets via *IPsec Security Associations* (SAs). For those protocols, the previous script would then auto-configure the IPsec SA instead of an in-protocol key parameter. Such scripts are, of course, autonomic service agents by another name.

In summary, GRASP with ANI can solve the recurring core problems of in-network automation between routers:

Q: How do I communicate with a peer (link-local or across other routers) without having any configured IP connectivity?

A: ACP provides this connectivity automatically with no human intervention.

Q: How do I discover what peers with what type of services are available (especially when not link-local)?

A: GRASP discovers the peers.

Q: Should I trust these peers?

A: Your trust comes from the ANI certificate used for the ACP. No nodes that have not been registered for the domain and authenticated by their manufacturer can join.

Q: How can I avoid re-inventing a new protocol to coordinate with peers?

A: Use GRASP.

Securing existing protocols is only one example where you can use ANIMA immediately. Many or all the benefits apply equally to any other in-network function with similar issues: establishing and adjusting *Quality of Service* (QoS) and other policies; auto-configuring decentralized protocol instances; monitoring, fault isolation, and troubleshooting; and even auto-configuring the most basic user network configuration, such as IP prefix distribution as in the previous example. When completely new services are required, ASAs should be developed in languages best suited for such a task. This immediate applicability to real-world problems provides a significant deployment incentive.

Summary and Conclusion

The ANI is a foundation for network automation and it serves two purposes:

- For existing network OAM designs it provides core functions to more easily build and deploy networks with secure, resilient network management. ANI provides automated public key deployment and renewal and zero-touch auto-configured in-band network management connectivity that is protected from being brought down by operator or network management tool errors.
- For ongoing further automation of network OAM (with or without an ultimate goal of fully autonomic networking), the ANI provides fundamental functions to build distributed, in-network automation agents (ASA) without having to re-implement their core dependencies each time: security, mutual trust, connectivity, and network-wide and peer-to-peer common signaling (via GRASP).

As a system, ANI may look overwhelming at first with its large set of constituent components (buzzword bingo), but it is fundamentally a very pragmatic approach, with the goal of making network complexity self-managing.

- *The basis of ANI is a set of long-term, well-known, and widely-used protocol components:* IPv6, X.509, IPsec, DTLS, RPL, CBOR, etc.
- The core innovations of ANI are built on top of this foundation: BRSKI, Voucher, MASA on top of X.509, ACP on top of IPsec, DTLS and RPL, and GRASP on top of CBOR.
- *ANI is highly modular:* All components are defined to be fully reusable individually or in concert. Adopt and deploy only the subset you need.

Any sufficiently advanced technology is indistinguishable from magic.^[15]

References and Further Reading

- [1] Alex Clemm, Michael Behringer, Sheng Jiang, Max Pritikin, Laurent Ciavaglia, Steinthor Bjarnason, and Brian Carpenter, “Autonomic Networking: Definitions and Design Goals,” RFC 7575, June 2015.
- [2] Michael Behringer, Sheng Jiang, and Brian Carpenter, “General Gap Analysis for Autonomic Networking,” RFC 7576, June 2015.
- [3] Toerless Eckert, Max Pritikin, Kent Watsen, and Michael Richardson, “A Voucher Artifact for Bootstrapping Protocols,” RFC 8366, May 2018.
- [4] Toerless Eckert and Michael Behringer, “Using an Autonomic Control Plane for Stable Connectivity of Network Operations, Administration, and Maintenance (OAM),” RFC 8368, May 2018.
- [5] Carsten Bormann, Brian Carpenter, and Bing Liu, “GeneRic Autonomic Signaling Protocol (GRASP),” RFC 8990, May 2021.
- [6] Brian Carpenter, Bing Liu, Wendong Wang, and Xiangyang Gong, “GeneRic Autonomic Signaling Protocol Application Program Interface (GRASP API),” RFC 8991, May 2021.
- [7] Sheng Jiang, Zongpeng Du, Brian Carpenter, and Qiong Sun, “Autonomic IPv6 Edge Prefix Management in Large-Scale Networks,” RFC 8992, May 2021.
- [8] Michael H. Behringer, Brian Carpenter, Toerless Eckert, Laurent Ciavaglia, and Jéferson Campos Nobre, “A Reference Model for Autonomic Networking,” RFC 8993, May 2021.
- [9] Toerless Eckert, Michael H. Behringer, and Steinthor Bjarnason, “An Autonomic Control Plane (ACP),” RFC 8994, May 2021.

- [10] Max Pritikin, Michael Richardson, Toerless Eckert, Michael H. Behringer, and Kent Watsen, “Bootstrapping Remote Secure Key Infrastructure (BRSKI),” RFC 8995, May 2021.
- [11] Brian Carpenter, “Autonomic Networking,” *IETF Journal*, 2014, <https://www.ietfjournal.org/autonomic-networking/>
- [12] FCC, “June 15, 2020 T-Mobile Network Outage Report, A Report of the Public Safety and Homeland Security Bureau Federal Communications Commission,” PS Docket No. 20-183, October 2020.
<https://docs.fcc.gov/public/attachments/DOC-367699A1.docx>
- [13] Alex Clemm, Max Pritikin, Michael H. Behringer, and Steinthor Bjarnason, “A Framework for Autonomic Networking,” Internet Draft, work in progress, October 2013.
<https://datatracker.ietf.org/doc/html/draft-behringer-autonomic-network-framework-01>
and
Max Pritikin, Michael H. Behringer, and Steinthor Bjarnason, “Bootstrapping Trust on a Homenet,” Internet Draft, work in progress, February 2014.
<https://datatracker.ietf.org/doc/html/draft-behringer-homenet-trust-bootstrap-02>
- [14] Arthur C. Clark, “Hazards of Prophecy: The Failure of Imagination,” in *Profiles of the Future* (1962).
- [15] Arthur C. Clark, *Profiles of the Future* (revised edition, 1973).

MICHAEL BEHRINGER worked for 18 years at Cisco, where starting in 2010 he led the Autonomic Networking project. Since 2017 he has been working as an independent consultant. E-mail: michael.h.behringer@gmail.com

CARSTEN BORMANN likes bringing the Internet to odd places. Honorary professor for Internet Technology at the Universität Bremen, his research interests are in protocol design and networking system architectures. He is behind many of the IETF's *Internet of Things* efforts, including *Constrained RESTful Environments* (CoRE) and the *Constrained Application Protocol* (CoAP), the *Concise Binary Object Representation* (CBOR), and the *Concise Data Definition Language* (CDDL). He co-chairs the Thing-to-Thing Research Group (T2TRG) in the *Internet Research Task Force* (IRTF). He has authored and co-authored 43 Internet RFCs. E-Mail: cabo@tzi.org

BRIAN E. CARPENTER, M.Sc., Ph.D., is an Honorary Professor at the University of Auckland. Previously he worked in networking at IBM and CERN. He has chaired the IETF, the Board of the Internet Society, and the Internet Architecture Board. E-mail: brian.e.carpenter@gmail.com

TOERLESS ECKERT is co-chair and liaison contact of the *Autonomic Networking Integrated Model and Approach* (ANIMA) Working Group of the IETF. He worked for 18 years in the Cisco IOS development group at Cisco Systems UK and USA, where he started to develop the Autonomic Networking architecture. Since 2016, he is a Distinguished Engineer at Futurewei USA. Email: tte@cs.fau.de

SHENG JIANG, Ph.D., is a Distinguished Engineer of the Network Technology Laboratory, Huawei Technologies. He co-chairs the *Autonomic Networking Integrated Model and Approach* (ANIMA) Working Group of the IETF and has authored 28 RFCs. E-mail: jiangsheng@huawei.com

YIZHOU LI, M.Sc., is a Principal Engineer in Network Technology Lab, Huawei Technologies. Previously she worked in networking at Singapore Telecom. She is the Secretary of the *Network Virtualization Overlays* (NVO3) Working Group of the IETF. E-mail: liyizhou@huawei.com

JÉFERSON CAMPOS NOBRE, M.Sc., Ph.D., is assistant professor at the Institute of Informatics, *Federal University of Rio Grande do Sul* (UFRGS), Brazil. He is co-chair of the IETF-LAC Task Force from the *Latin American and Caribbean Network Operators Group* (LACNOG) and co-secretary of the *Network Management Research Group* (NMRG) of the *Internet Research Task Force* (IRTF). He has been involved in Autonomic Networking research since 2007. E-mail: jcnobre@inf.ufrgs.br

MICHAEL C. RICHARDSON, B.Sc. Physics/Computer Science, is an open source and open standards consultant. An autodidact, he wrote mail transfer agents as a teenager, and in the 1990s, after failing at high-energy physics, found his calling designing and building embedded networking products in the security sector. E-mail: mcr@sandelman.ca

Check your Subscription Details!

If you have a print subscription to this journal, you will find an expiration date printed on the back cover. For several years, we have “auto-renewed” your subscription, but now we ask you to log in to our subscription system and perform this simple task yourself. Make sure that both your postal and e-mail addresses are up-to-date since these are the only methods by which we can contact you. If you see the words “Invalid E-mail” on your copy this means that we have been unable to contact you through the e-mail address on file. If this is the case, please contact us at ipj@protocoljournal.org with your new information. The subscription portal is located here: <https://www.ipjsubscription.org/>

A Survey on Securing Inter-Domain Routing

Part 1 – BGP: Design, Threats, and Security Requirements

by Geoff Huston, APNIC

The *Border Gateway Protocol* (BGP) is the inter-domain routing protocol on the Internet, and after some 30 years of operation, BGP is now one of its more venerable core protocols. One of the major ongoing concerns related to BGP is its lack of effective security measures, and as a result the routing infrastructure of the Internet continues to be vulnerable to various forms of attack.

In Part 1 of this study, we will look at the design of BGP, the threat model, and the requirements from a security framework for BGP. In Part 2 we will look at the various proposals to add security to the routing environment and also evaluate the current state of the effort in the *Internet Engineering Task Force* (IETF) to provide a standard specification of the elements of a secure BGP framework.

Introduction

The Internet is a decentralised collection of interconnected component networks (autonomous systems). These networks are composed of *end hosts* (who originate and/or receive IP packets, and are identified by IP addresses) and active forwarding elements (routers) whose role is to direct IP packets as they pass through the network. The routing system is responsible for propagating the relative location of IP addresses to each routing element, so that routers can make consistent and optimal routing decisions in order to pass a packet from its source to its destination. Routing protocols are used to perform this information propagation.

The routing system of the Internet is divided into a two-level hierarchy. One level is *intra-domain* routing, which the set of autonomous routing systems operating within each component network use. The other level is a single *inter-domain* routing system that maintains the inter-network connectivity information that straddles these component networks. A single inter-domain routing protocol, BGP^[1] has provided inter-domain routing services for the disparate component networks on the Internet since the late 1980s.^[2] Given the central role of routing in the operation of the Internet, BGP is one of the critical protocols that provide essential coherence to the Internet.

The underlying distributed distance vector computations of BGP rely heavily on informal trust models associated with information propagation to produce reliable and correct results. You could liken them to a hearsay network—information is flooded across a network as a series of point-to-point exchanges, with the information being incrementally modified each time it is exchanged between BGP speakers. The design of BGP was undertaken in the relatively homogeneous and mutually trusting environment of the early Internet.

Consequently, its approach to information exchange was not designed primarily for robustness in the face of various forms of negotiated trust or overt hostility on the part of some routing actors.

Hostile actors are a fact of life in today's Internet. It's quite reasonable to characterise today's Internet environment as one where trust must be explicitly negotiated rather than assumed by default. This environment is no longer consistent with the inter-domain trust framework that BGP originally assumed. The BGP mutual trust model involves no explicit presentation of credentials, no propagation of instruments of authority, nor any reliable means of verifying the authenticity of the information being propagated through the routing system. Hostile actors can attack the network by exploiting this trust model in inter-domain routing to their own ends.

An attacker can easily transform routing information in ways that are extremely difficult for any third party to detect. For example, false routing information may be injected, valid routing information removed, or information altered to cause traffic redirection.^[3,4,5] You can use this approach to prevent the correct operation of applications, to conduct fraudulent activities, and to disrupt the operation of part (or even all) of the network in various ways. The consequences range from relatively inconsequential (minor degradation of application performance due to sub-optimal forwarding paths) through to catastrophic (major disruption to connectivity and comprehensive loss of any form of cohesive Internet). To resist this subversion of integrity of routing information, each BGP speaker must have:

- Sufficient information at hand to verify the authenticity and completeness of the information being provided to it via the inter-domain routing system, and
- The ability to generate authoritative information such that other BGP speakers may verify the authenticity of information that this speaker is passing into the inter-domain routing system.

A key question is whether we can add further information into the inter-domain routing environment such that attempts to pervert, remove, or withhold routing information may be readily and reliably detected. Any proposed scheme must also be evaluated for its impact on the scaling properties of BGP.

To ground any such evaluation of BGP, it's useful to briefly review the design of the BGP protocol.

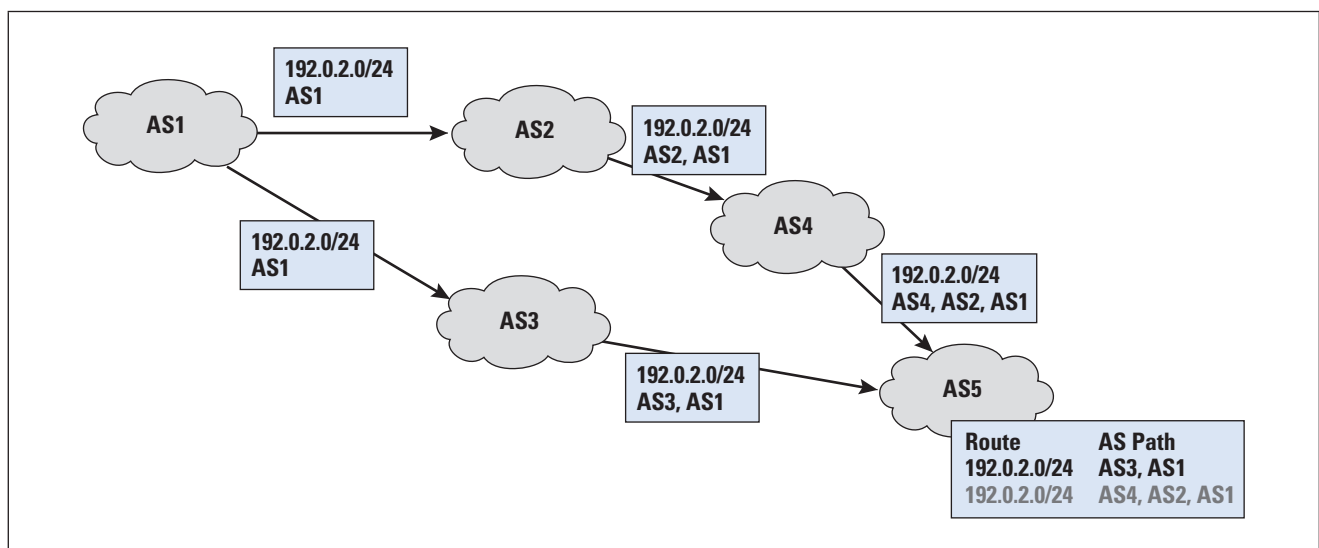
The Design of BGP

BGP underwent numerous refinements over its early operational life. The protocol was originally described in RFC 1105 in June 1989,^[6] allowing the inter-domain architecture of the Internet to move on from a constrained architecture of *core* and attached *stub* domains into a framework of peer routing domains without any central core.

A refinement to this protocol, BGP-2, was described in RFC 1163 in June 1990,^[7] and a further refinement, BGP-3, was described in RFC 1267 in October 1991.^[8] The current version, BGP-4, was first deployed within the Internet in 1993. The RFC describing this protocol, RFC 1771,^[9] was published in March 1995, and was subsequently refined with the publication of RFC 4271 in January 2006.^[1] The core protocol has been stable for some years now, although further refinement has been undertaken through the use of negotiated capabilities undertaken at BGP session startup.

BGP is an instance of what we commonly refer to today as a *Bellman-Ford Distance Vector* routing algorithm.^[10,11] This algorithm allows a collection of connected devices (BGP speakers) to each learn the relative topology of the connecting network. Its basic approach is very simple: each BGP speaker tells all its other neighbours about what it has learned if the new learned information alters the local view of the network. This scenario is a lot like a social rumour network, where everyone who hears a new rumour immediately informs all their friends. BGP works in a very similar fashion: each time a neighbour informs a BGP speaker about reachability to an IP address prefix, the BGP speaker compares this new reachability information against its stored knowledge that it gained from previous announcements from other neighbours. If this new information provides a “better” path to the prefix, then the local speaker moves this prefix and associated next-hop forwarding decision to the local forwarding table and informs all its immediate neighbours of a new path to a prefix, implicitly citing itself as the next hop. BGP keeps track of the propagation of route advertisements across the inter-domain space by recording the sequence of network *Autonomous Systems* (ASs) that propagate the route in a route attribute called the *AS Path*. A “better” route is one with a shorter AS path, and a loop is detected when a BGP speaker sees its own AS in the received AS Path (Figure 1).

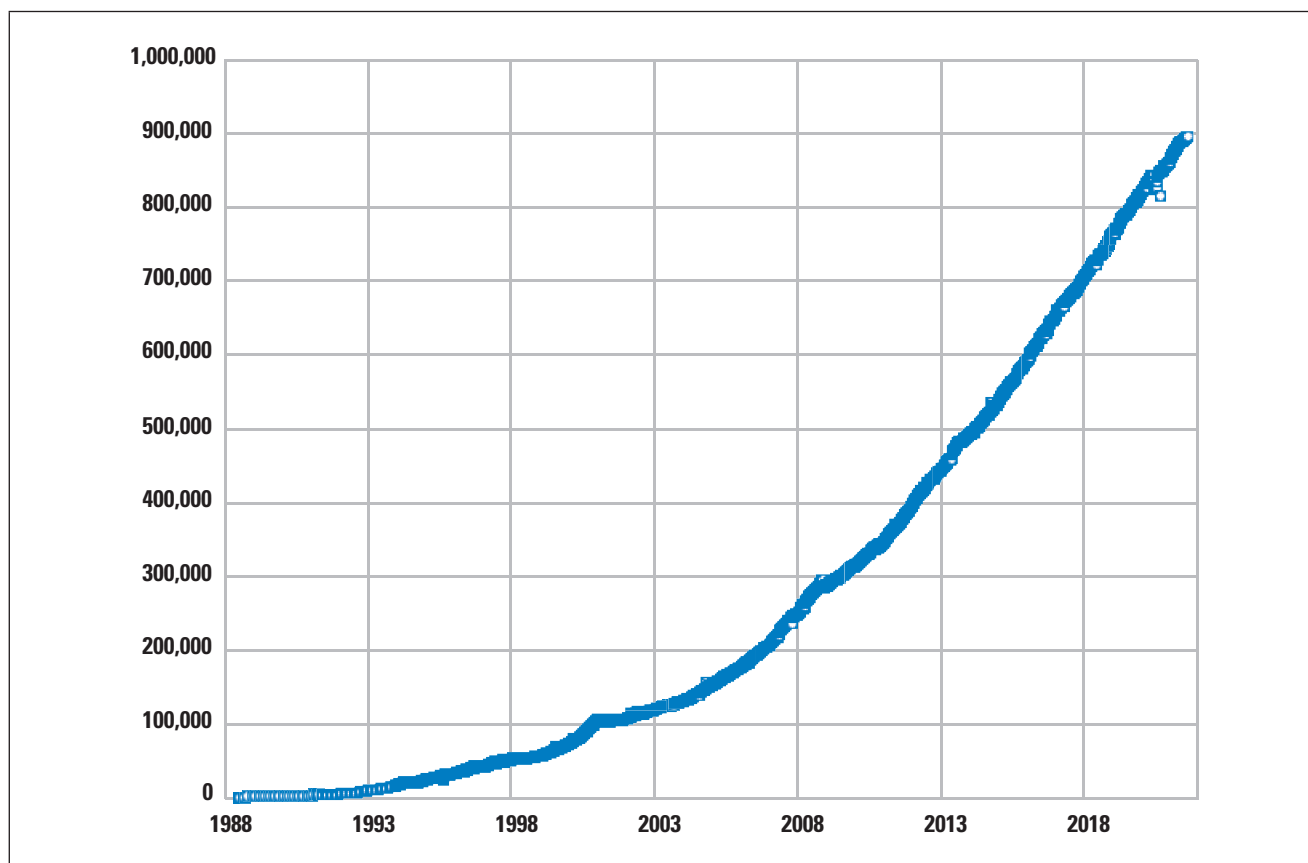
Figure 1: The Propagation of a route in BGP



In addition, there is a withdrawal mechanism, where a BGP speaker determines that it no longer has a viable path to a given prefix, in which case it announces a *withdrawal* to all its neighbours. When a BGP speaker receives a withdrawal, it stores the withdrawal against this neighbour. If the withdrawn neighbour happened to be the currently preferred next hop for this prefix, then the BGP speaker examines its per-neighbour data sets to determine which stored announcement represents the best path from those that are still extant. If it can find such an alternative path, it copies this path into its local forwarding table and announces this new preferred path to all its BGP neighbours. If there is no such alternative path, it announces a withdrawal to its neighbours, indicating that it no longer can reach this prefix.

Across the deployment lifetime of BGP-4, the IPv4 Internet has grown from an average of 20,000 distinct routing entries in 1993 to almost 1 million routing entries in 2021.^[12] Figure 2 shows the growth of the size of the Internet IPv4 routing table over time.

Figure 2: Internet IPv4 Routing Table Size, from [12]



BGP and TCP

BGP is not a link-level topology maintenance protocol. It assumes the existence of a relatively robust IP forwarding environment at the link level between BGP peers. This assumption has allowed BGP to use the *Transmission Control Protocol* (TCP) as a reliable transport protocol to support the transactions of the protocol across a BGP peer session.

TCP manages reliable message delivery and flow control between the BGP peers and allows BGP to operate across end-to-end connections whether they reside on the same subnet or across the Internet. There is no requirement for BGP speakers to be connected on a common media connection, and the choice of TCP allows this flexibility of connectivity by requiring only that a BGP peering session is supported by an IP network.

The TCP stream is divided into messages using BGP-defined markers, where each message is between 19 and 4096 octets long, extensible to 65,535 octets.^[11] The use of a reliable transport service implies that BGP itself need not explicitly confirm receipt of protocol messages, removing much of the protocol overhead seen in other routing protocols that sit directly on top of a media-level connection. There are no message identifiers, no message number initiation protocols, no explicit acknowledgement of messages, nor any provision to manage lost, reordered, or duplicated messages. TCP handles all of that. The use of a reliable transport protocol also obviates the need for BGP to periodically refresh the routing state by automatically reflooding the entire routing information set between BGP speakers. After the initial exchange of routing information, a pair of BGP routers exchange only incremental changes to routing information.

BGP Messages

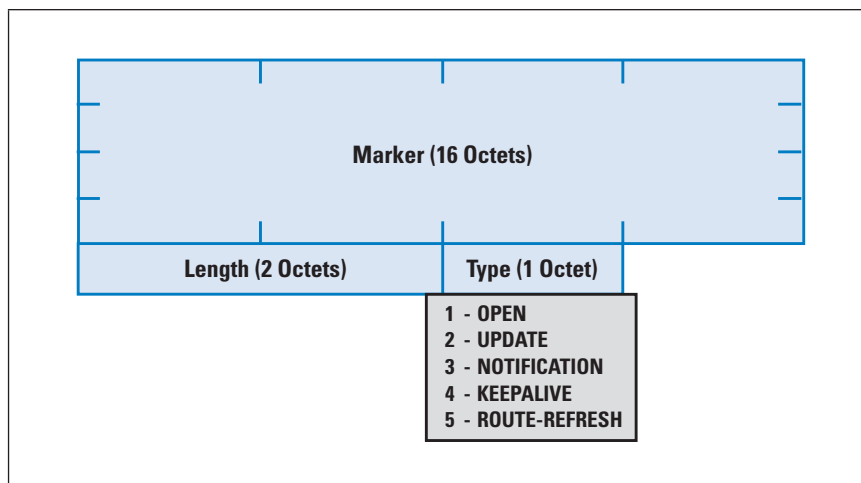
Because TCP is a *stream protocol* rather than a *record-oriented protocol*, BGP uses record marking within the TCP stream to delineate logical protocol units, or *messages* with a 16-byte marker as the BGP message delimiter. A 2-byte length and a 1-byte type field follow the marker, making the minimum BGP message size 19 bytes. The repertoire of defined messages follows:

- An OPEN message to start a BGP session
- An UPDATE message to exchange reachability information
- A NOTIFICATION message, which is used to convey a reason code prior to termination of the BGP session
- A KEEPALIVE message, used to confirm the continued availability of the BGP peer
- A ROUTE-REFRESH request message to request a resend of the routing information.

Figure 3 on the next page shows the common format of BGP messages.

BGP uses an explicit OPEN message to commence a BGP peering session. This message exchange confirms the identity of the BGP speakers and includes the option for a capability negotiation to understand what optional or extended capabilities each BGP speaker supports. A session is active only when both BGP speakers have sent their OPEN messages and neither has rejected the other's offered capabilities through a NOTIFICATION response.

Figure 3: BGP Common Header Message Format



When the session is active, BGP operates via the exchange of UPDATE messages. Each UPDATE message contains a set of address prefixes that are unreachable (withdrawals), followed by a set of common route object attributes and a set of address prefixes that share this set of attributes (announcements). The withdrawn prefixes are those prefixes where the local BGP speaker sees no reachability, and now wants to withdraw a previous advertisement of reachability. No routing attributes are associated with these withdrawn prefixes.

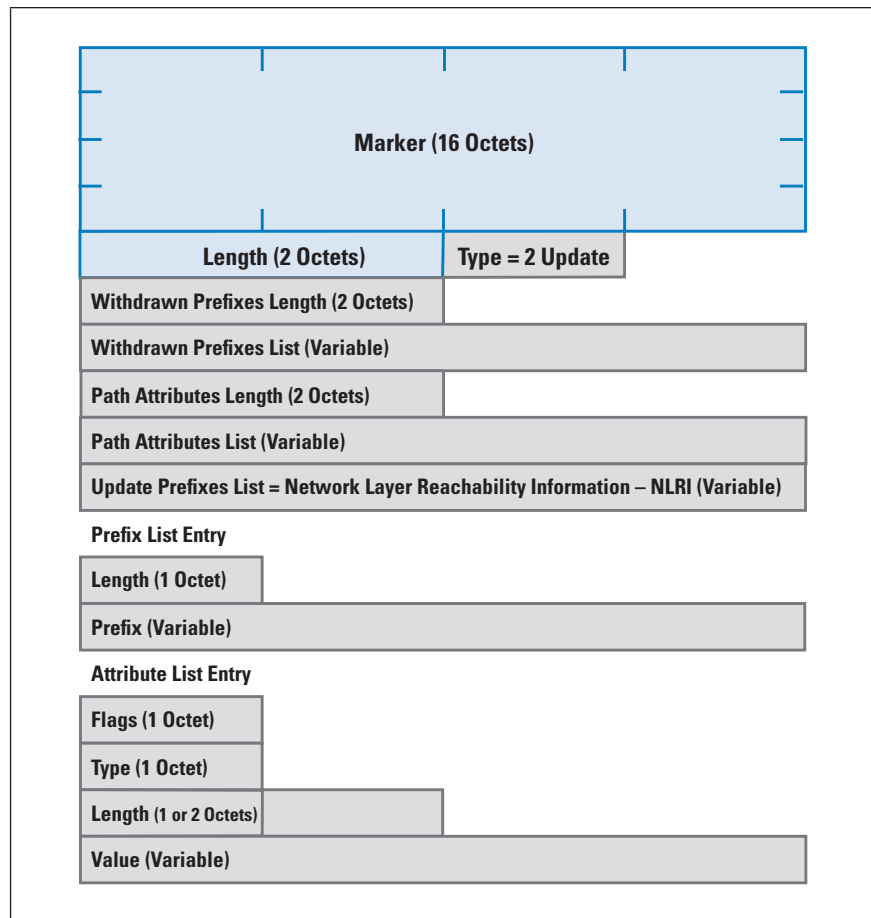
The announced prefixes are those prefixes where the local BGP instance has an updated view of the reachability of a prefix that was previously withdrawn or unannounced or has an updated view of the routing attributes of the locally selected “best” route for a prefix. BGP may group multiple prefixes together in a single UPDATE message but can do so only if all the updated prefixes share a common set of attributes. Within an UPDATE message, the withdrawn prefix set or the announced prefix set may be empty, but not both. Figure 4 on the following page shows the layout of the BGP UPDATE message.

AS Path Attribute

BGP binds together the concept of network address blocks and autonomous systems into a path vector-based routing technology. Every route object represented within a BGP-4 route database contains an address prefix and an associated path vector of AS values. BGP does not indicate the precise path a packet should follow within an AS, nor does it maintain a complete map of the topology of the Internet at a link-by-link level. BGP uses a level of abstraction that views the Internet as a set of per-AS routing domains, and the role of BGP is to maintain a routing map of the network at this AS level, associating every reachable address prefix with an AS transit path from the current location to the originating AS of the address prefix.

One of the most important route object attributes in BGP is the AS Path attribute of UPDATE messages that contain announced routes.

Figure 4: BGP UPDATE
Message Format



As address prefix reachability information traverses the Internet in the form of individual route objects in BGP, this BGP routing information is augmented by the list of autonomous systems that have processed this route information thus far, forming the AS Path attribute of a route object. Each BGP speaker adds its own AS value to the AS Path attribute of the route object when passing the route object through an *External BGP* (eBGP) session.

This AS Path attribute allows straightforward suppression of the looping of routing information, using the simple algorithm that a local AS will reject any forwarded route object that already contains its own AS in the AS Path attribute. Also, the length of the AS Path vector forms the BGP route metric. A local BGP system, when attempting to select one from numerous potential route objects that refer to the same address prefix, will, in the absence of any local policy directive, prefer the route object with the shortest AS Path length.

In addition to undertaking the role of path metric and loop detector, the AS Path attribute serves as a versatile mechanism for policy-based routing, where a local AS can alter the default preferences for route selection based on local policy settings coupled with pattern-matching rules to be performed on the AS Path.

Withdrawals have no associated AS Path.

BGP Route Selection Process and Routing Policies

A BGP speaker may receive two or more announcements for the same address prefix from different peers. The “best” announcement is selected as the locally used announcement, and this announcement is the one that is announced to its BGP peers. BGP defines an ordered sequence of comparisons to determine which route object is selected by the local BGP speaker as the preferred route to use:

- Prefer the route object with the highest value for LOCALPREF attribute value.
- Prefer the route object with the shortest AS PATH attribute length.
- Prefer the lowest origin value.
- Prefer the lowest MULTI EXIT DISCRIMINATOR attribute value.
- Prefer the minimum *Interior Gateway Protocol* (IGP) cost to the NEXT HOP address given in the route object.
- Prefer eBGP over *Interior BGP* (iBGP)-learned routes.
- If using iBGP, prefer the lowest *BGP Identifier* value.

Although network administrators usually employ routing policies depending on their needs,^[14,15] within the generic BGP route selection process the highest-priority selection rule is that a route for a more specific address prefix is to be preferred over that of a covering prefix.

The BGP Threat Model

One approach to providing a taxonomy for threats in routing in general, and BGP in particular, is to view a BGP peer session as a conversation between two BGP speakers and pose numerous questions relating to this conversation. These questions include:

- *How do we talk?* The manner in which the BGP session between the BGP speakers is secured such that the conversation is not altered, disrupted, or hijacked and is protected from unauthorised eavesdropping
- *Whom am I talking to?* Verification of the identity of the other party and verification that they are authorised to speak for the routing entity that they purport to represent.
- *What are you saying?* Verification of the authenticity and completeness of the routing information being passed in the BGP session.
- *Why should I believe you?* Verification that the routing information represents the current state of the forwarding system.
- *How recent is your information, and is it still valid?* Verification of how long routing information is valid and whether the information is still current.

You can further deconstruct each of these security questions to a set of specific objectives, as well as recognise a set of specific threats.

Securing a BGP Session

A BGP session between two BGP speakers is assumed to have some level of integrity at the session transport level.

BGP assumes that the messages one party sends are precisely the same messages the other party receives, and that the messages have not been altered or reordered, have not had spurious messages added into the stream, or have messages removed from the conversation stream in any way, and given that BGP uses a TCP transport session, some of these assumptions are reasonable but others less so.

As with any long-held TCP session, a BGP peer session is vulnerable to eavesdropping, spurious session reset, session capture, message alternation, and *Denial-of-Service* (DoS) attacks, all through what we might think of as conventional TCP attack vectors.

The threat at the BGP level is that a third party may attempt to break into the TCP session as an interception attack in the middle, and thereby alter the BGP message flow between the two end points. One form of threat is by injection, where the attacker injects spurious messages into the BGP session. Direct on-the-wire interception allows the attacker to have knowledge of the TCP sequence numbers, thereby making injection a trivial task. Even if the attacker is not able to intercept or eavesdrop the BGP session, it is still possible to attempt to guess the current sequence number.

While this guessing is often impractical in the case of injecting data into the session, if all that is to be injected is a TCP Reset, then the sequence number guess only has to sit within the current TCP window in order to be recognised as a valid reset TCP message.^[16] Another form of threat is by active intermediation, where the attacker sits on the connection between the two BGP speakers and intercepts all traffic in both directions. In this case, the attacker has complete control of the BGP message stream and can perform any form of message alteration. A variation of this form of threat is *session hijacking*, where the third party intrudes upon an active BGP session and injects its own traffic into the message stream—and that traffic allows the third party to take over the session and masquerade as one of the parties to the BGP session. Because timing is important in the overall performance of BGP, another form of attack at the session level is to delay messages. Although the content of the messages is unaltered, the implicit timing signals within the message stream are altered by this form of intervention, potentially causing the local BGP speaker to behave differently and fall out of sync with its routing peers.

Another form of attack is a *replay attack*, where older BGP messages are replayed into a hijacked TCP session. One form of this replay attack could be to replay a pair of messages that withdraw and then announce the same address prefix.

Route Flap Damping (RFD)^[17,18] is a widespread defensive BGP configuration that monitors the frequency of BGP updates for a given prefix from each peer, and if the update rate exceeds a locally set threshold, the advertisement of this prefix by the peer will be locally suppressed for a damping interval. The replay of updates could be used to trigger an RFD response in the remote BGP speaker.^[19] If a route is fully dampened through RFD, the BGP speaker will not advertise updates for this prefix for a damping interval (commonly 60 minutes), possibly causing a route disruption within that time frame. Another form of replay attack is to replay a route advertisement for a previously withdrawn prefix, possibly in conjunction with some form of prefix hijack attack.

Another form of threat is withholding traffic. BGP uses keepalive timers to determine remote end “liveness.” By intercepting and withholding all messages for the hold-down timer interval, a third party can force the BGP session to be terminated and reset. This action causes the entire route set to be re-advertised upon session resumption so that repeated attacks of this form can be an effective form of denial of service for BGP.

It is also possible to undertake a *saturation attack* on a BGP speaker by sending it a rapid stream of invalid TCP packets. In this case, the processing capability of the BGP speaker is put under pressure, and the objective of the attack is to overwhelm the BGP speaker and cause the BGP session to fail and be reset. This type of attack is particularly problematic if the BGP session uses *Message Digest 5* (MD5) or *Internet Protocol Security* (IPsec) as session protection protocols, because the cryptographic function overhead also applies to the injected packets, increasing the processing overhead on these spurious injected packets.

The underlying aspect of the BGP protocol is that BGP itself has no enforced minimum level of message protection. BGP messages are, by default, placed into the TCP stream without encryption or additional message wrapping of message sequencing. Any threat that is applicable to long-held TCP sessions applies to this default mode of BGP operation.

Verifying BGP Identity

BGP sessions commence by passing the local AS to the remote end of the session in the BGP OPEN message and receiving the AS of the remote end in the received OPEN message. BGP itself does not verify these asserted AS identities, and it is theoretically possible for a remote party to masquerade itself as another AS and assert an identity in BGP that the other party cannot directly verify by, and neither can any third party that subsequently receives this routing information. Most BGP implementations provide a level of protection against this threat by applying a constraint that the local BGP speaker will initiate a peer session only with a configured remote IP address, and reject all other TCP connection attempts.

Furthermore BGP will not complete the BGP OPEN message exchange if the AS in the OPEN message does not match the AS number associated with the remote end IP address in the configuration.

This approach places a heavy reliance on the out-of-band process of BGP configuration, and if an attacker can compromise or take control over BGP equipment connected to the Internet or use social engineering to convince a network administrator to configure incorrect information into the BGP configuration, then it is possible to masquerade as a different party in BGP and potentially inject incorrect information into the routing system.

The real question here is: “Are you really who you claim to be?” Here it is necessary for the BGP speaker to be able to confirm the validity of the peer claim that it is speaking for an AS.

Verifying BGP Information

The objective here is to verify the authenticity and completeness of the routing information being passed in the BGP session. The intention of BGP is that a local BGP speaker provides to all its BGP peers a complete feed of its locally selected route objects.

When a session is opened with a remote BGP speaker, the local BGP instance believes everything it is told without further qualification. The threat is that a BGP peer can deliberately feed false information to the local BGP instance, which BGP itself will be unable to detect as false. The false information could be in the form of suppression of routing information, or alteration of the route object that is being passed, or the invention of spurious route objects. The BGP speaker could be asserting that an AS Path is genuine when it reflects an artificial path, or that it has the authority to originate an advertisement for a prefix when, in fact, no such authority exists.

A BGP speaker may preserve all the attributes of a route object, but alter the prefix set to be the equivalent collection of more specific prefixes. The deliberate alteration of routing information can cause the local BGP instance to make an incorrect choice of a local best path and also cause the local BGP instance to propagate this incorrect information to its neighbours.

Not only could the BGP speaker be passing incorrect attributes for an address prefix in order to bias the local route selection process, but it also could be providing incorrect information regarding the prefix itself. The prefix that is the subject of the route object could be a prefix that has never been allocated and should not be legitimately routed, or the prefix could be an aggregate address prefix that spans both allocated and unallocated address space.

Prefix hijacking is a major threat to the integrity of the BGP routing. The fundamental weakness here is that BGP provides no explicit means of verifying the authenticity of the address prefixes that are listed in a BGP UPDATE message, nor the authenticity of the attributes of the prefix, including the origination information and the AS Path vector. The threat here is that by deliberately altering this information, the local BGP speaker can be induced to make incorrect route selection decisions and thereby make incorrect forwarding decisions for IP traffic.

A known common problem illustrative of exploiting this vulnerability is operational misconfiguration,^[20] which could result in propagating more specific routes and other forms of route leakage, or withholding that may affect the routing decisions made by other BGP speakers. This form of verification of intentionality by a remote BGP speaker is far more challenging—while these forms of security mechanisms are intended to verify that the received information matches the original information that was passed into the routing system, they are incapable of verifying that such information is consistent with the true intent of the originator of the information.

Verifying Forwarding Paths

The overall intention of the BGP protocol is to distribute the current binding of address to location such that individual routers can make accurate judgements about how to populate their local forwarding tables and hence make optimal local decisions for each packet that passes along the shortest path to its ultimate destination.

BGP does not provide any ability for a local BGP speaker to validate that the route advertisements it receives from a BGP peer accurately represent the current state of the network forwarding system. The threat model here is that a bad actor in the routing system may make a different forwarding decision to that being advertised in the routing system.

This situation can represent a subversion of local policies, theft of carriage capacity, deliberate denial of service, or the potential to eavesdrop on a conversation or support the interception and alteration of application-level transactions. Even a completely secured control plane does not avert such vulnerabilities.^[21]

The Consequences of Attacks on the Routing System

The ability to alter the routing system provides a broad array of potential consequences.^[3] The consequences fall into numerous broad categories, which are briefly described here:

1. *The ability to eavesdrop.* The forwarding system can be altered so as to pass all traffic to a class of destination addresses through a certain path. This change allows the attacker to attempt to pass all such traffic through an eavesdropping location prior to conventional delivery. In such a case the parties may not be aware that an eavesdropping attack is taking place.

2. *Denial of service.* The simplest form of a DoS is where traffic to an address prefix is passed to a point where it is then discarded. Routing loops also are a form of DoS, where not only will the traffic to a destination address prefix never reach its intended destination, but the traffic will be held in the loop for the life of the packet *Time to Live* (TTL) field. For sufficiently short loops the potential exists for the loop to act as a link load amplifier, where the traffic on the loop is several times the traffic load being addressed to the affected destination address prefix.
3. *The potential to masquerade.* Subversion of routing allows sites to masquerade as other sites; the routing system misdirects the traffic to the masquerading site. The consequences of such an attack can vary from the specific, where a particular site is targeted, to the more generic, where authoritative *Domain Name System* (DNS) servers are the subject of the masquerading attack, and the DNS responses are believed to be authentic. In this case if the masquerading occurs at the root level of the DNS hierarchy, incorrect information can be provided to any query, allowing for the attack to then be extended to any site.
4. *The ability to steal addresses and obscure identity.* Routing an unallocated address is subtly different from routing an already allocated address. Here the consequence is not displacement of traffic forwarding to incorrect locations in the network, but the assertion of the existence of addresses and forwarding paths to those addresses that should not exist in the network in the first place. The consequence is the ability to use addresses on the network that have no allocation registration information associated with them, allowing the originator of the routing attack some degree of ease to mount an anonymous attack at the application level. Such forms of attack have been observed to be associated with SPAM and botnet controllers where anonymity of the attack coordinator is desired.

Security Requirements

The primary requirements for securing BGP are securing both the transmission of the data payload of the BGP protocol and the semantics of that payload.

The security requirements for transmission are such that the data that a BGP speaker receives can be cryptographically verified to have been sent by the BGP peer, the data is not a replay of previously transmitted data, and no data has been removed from the transmission.^[22]

There is no strict requirement for encryption of the BGP payload, because the routing information being exchanged is not intrinsically confidential to the two parties involved. The security requirements for the semantics of the payload concern specifically some selected fields (transitive attributes) of the BGP UPDATE message. The BGP speaker must be able to verify that the advertised prefix is valid, and that the originating AS has been duly authorised by the legitimate right-of-use holder for that prefix.

The BGP speaker should also be able to validate that the AS Path in the UPDATE represents a valid inter-AS transit path through the network in terms of inter-AS topology and AS transit policies, and that the prefix reachability information has been propagated along the reverse inter-AS Path.^[22]

It is noted that route withdrawals and nontransitive announcement attributes are local, and thus do not need to be transitively protected in a similar fashion to route origination and the AS Path attribute of announcements. You can adequately protect withdrawals and local attributes with BGP peer session protection.

The associated requirements for a secure inter-domain routing system include that the additional use of security credentials and verification of routing information should not alter the temporal properties of the BGP protocol, and that authentication of the security credentials should occur in the same time frame as the BGP message processing operation. It is also a requirement that piecemeal incremental deployment should be feasible.^[23,24,25] A secure operational mode should be a capability negotiation with each BGP peer, with the ability to support backward compatibility with those BGP peers that do not recognise such a capability. It seems to be a good idea to start deployment of BGP security on the most-connected nodes and incrementally deploy it towards least-connected nodes.

Additionally, it suggests the question: How does a party that uses security credentials deal with information arriving from a peer that does not use any security credentials? Having no security credentials does not necessarily mean that the information is wrong, of course. But importantly, in these piecemeal deployment scenarios there should be some incremental benefit of piecemeal deployment to those actors who choose to supply such security credentials and those who choose to validate routing information using these credentials.

A routing system, secure or otherwise, should never make route selections that include routing loops. It is preferred that in a fully secured environment a secure routing system would be able to converge on best paths that are either identical to or no worse than an unsecured BGP speaker would select, assuming that such paths can be validated in a secure environment. In an environment of partial adoption of secure routing systems, it is recognised that a BGP speaker may use local preference settings that prefer sub-optimal paths that have preferred security credentials over unsecured paths.

The trust model of routing appears to involve two forms of trust. The first is a trust environment related to the public network and the legitimacy of use of a public address and a public AS number. It is necessary to be able to verify that a particular party has the right to use these number resources in a public context. The closest fit in the form of a trust model for verification of this assertion of right of use is a public authority that can provide authoritative information on the distribution of these numbers.

This approach leads to a rooted hierarchy model of trust, where the trust anchor is this public authority.

The second form is a trust environment in private contexts, where the use of an address or AS number is bounded by a specific context of use, and the trust in an assertion of a right of use is one made in the context of this bounded environment. In this environment, there is no clear ability to use public authorities as a trust anchor, and other means of trust that may involve reputation, or web of trust concepts may be appropriate.

A general security approach to BGP should be able to encompass that diversity of deployment environments and the corresponding diversity of authority models.

Tools for Securing BGP

The vulnerabilities of BGP arise from four fundamental weaknesses in the BGP and inter-domain routing environment, including:

- No mechanism to protect the integrity, currency, and source authenticity of BGP messages
- No mechanism to verify the authenticity of an address prefix and an AS origination of this prefix in the routing system
- No mechanism to verify the authenticity of the attributes of a BGP UPDATE message
- No mechanism to verify that the local cache *Routing Information Base* (RIB) information is consistent with the current state of the forwarding table

The other observation about BGP security is that it appears that by far the most straightforward form of attack is to obtain control and configuration access to a deployed router and use this compromised platform as the base for launching attacks on the routing system. In the face of such an encompassing attack on the control instruments of the routing system, BGP session-level security needs to be placed in some perspective. It is not possible to prevent routers from attempting to generate false information as long as routers themselves are in a position to be compromised.

The consequent vulnerability on the routing system, as distinct from a narrower view of BGP, is that there is no mechanism that limits the extent to which a misbehaving routing element can make inaccurate claims about reachability in the routing system.

The Security Toolset for BGP Session Protection

The available tools for securing BGP start at the level of the BGP TCP session and encompass the tools that are used to protect TCP and the two ends of the TCP session.

The TCP protection mechanisms include the generalized TTL security mechanism,^[26,27] which is intended to limit the effective radius of potential attack on the session to hosts that lie on or within the worst-case hop-count radius between the two BGP speakers and host-level defences against TCP SYN attacks.^[28] In many ways, this form of defence is effective when using multi-hop BGP sessions in that the attacker cannot subvert the defence, but it still leaves the session vulnerable to any attacker that lies within the TTL radius.

You can get greater levels of session protection by using cryptographic protection. Over time the IETF has worked on three approaches to protect the BGP TCP through cryptographic protection. They include:

- The use of IPsec.^[29] IPsec has not been widely used for BGP sessions, and the reasons why relate to the complications for rekeying *Internet Key Exchange (IKE)/IPsec* sessions and the potential *Distributed Denial-of-Service (DDoS)* vector.^[30]
- The *TCP MD5 Signature Option*.^[31] Although the MD5 signature option has some potential weaknesses when compared with IPsec,^[29] MD5 is considered preferable to no form of TCP protection at all, particularly with respect to the TCP Reset injection attack. However, there are issues with re-keying a long-held session, and the BGP speakers probably need to use graceful restart mechanisms in conjunction with MD5 to perform a re-key of the session.
- The *TCP Authentication Option*,^[32] which the IETF has marked as a replacement for the earlier MD5 approach. The *TCP Authentication Option* supports stronger crypto algorithms compared to MD5. It uses a two-fold security approach that reduces the critical reliance on a user-configured key. This approach also allows the configuration of up to 64 keys for a session and provides a simple key coordination mechanism by giving the ability to change keys (move from one key to another) within the same connection without causing any TCP connection closure. By comparison, changing TCP MD5 keys during an established connection might cause a flap or restart in the connection, which in the context of BGP may have operational implications.

From time to time the topic of BGP over *Transport Layer Security (TLS)*^[33] is raised, and it is possible that sooner or later we might hear of BGP over *Quick UDP Internet Connections (QUIC)*.^[34,43] The salient question is one of balancing the additional burden of adding more transport choices to BGP implementations with the likely benefits that these additional choices may provide. As we've seen in the IPv6 transition and more recently in the increasing diversity of choices for encrypting DNS transactions, adding more options can offer just confusion and impede adoption instead of accelerating it.

However, the most important guideline in securing BGP sessions is to use multi-hop BGP and multi-access LAN sessions sparingly and preferably use a direct 1:1 channel connection when such a choice is available.

The Security Toolset for BGP Message Protection: RPKI

In addition to message integrity protection that transparent session-level protection mechanisms provide, the tools to provide protection of the integrity of BGP messages relate to the use of digital signatures to provide a set of credentials that allow relying parties to verify the correctness of the information carried as the message payload in BGP.

The reason for the use of digital signatures as opposed to an integrity check using some form of shared secret was obvious after the observation that the number and identities of all eventual recipients of the information are not known in advance, and non-repudiation is desirable.^[3] Verification of the contents of a message is not only a test of whether the message has been altered in any way during its transit between BGP speakers, but also a test of whether the message represents correct origination information and correct operation of the processing of the message during the message propagation (*authenticity*).

This requisite implies a need to establish a means of verification of information where the author of any security credentials relating to origination and propagation is not necessarily known to the relying party that is attempting to validate the information. This need typically invokes a form of validation that relies upon third-party transitive trust, where the relying party is attempting to build a testable chain of trust between its trust anchor and the party or action that is the subject of the verification operation. Conventionally, this requirement implies the use of some form of *Public Key Infrastructure* (PKI). In this case, we are not looking to use such a PKI to validate claims of identity, authority to perform a particular function, or some form of verifiable attribution. We need some form of mechanism to associate a public key with an IP address prefix or an AS number in a sense of functional control, where the certification authorities in this PKI are attesting that the certified subject has *functional control* of a collection of IP number resources (AS numbers and IP address prefixes). The associated certificate issuance practices are intended to support transitive trust in such attestations of association.

We have adopted a structure using X.509 public key certificates and a certificate extension that uses a canonical list of IP address resources and AS numbers^[35] as the foundation for this *Number Resource PKI* (RPKI).^[36] Verification of a digital signature entails a test of the authenticity and current validity of the associated certificate that describes the public key of the address or AS number holder in the context of a structured set of signed relationships between certificate issuers and subjects. In other words, the holder of the matching private key is the current functional controller of those IP addresses and AS number and can digitally sign authorities and attestations about such number resources on the basis of that functional control.

Given that the discourse of BGP messages is about address prefixes and AS numbers, the RPKI provides a solid foundation for digital signatures to be associated with various routing actions that are described in BGP messages.^[37] It does not attest in any way to the identity of these number resource holders.

Anchoring the model of authority and trust in the RPKI certificate structure has resulted in a framework where the issued certificates are aligned with the IP address and AS number allocation and assignment framework. If an Internet Registry has issued a set of IP addresses and AS numbers to an entity, then this registry would be able to publish a public key certificate that associates a private key provided by the entity with the IP resource set. Further allocations from a registry to a registry address holder would result in re-issuance of the certificate for the address holder with a larger resource set, while reduction in this set would result in both re-issuance and revocation of the previous certificate. This certificate framework would allow auditing of the certificate state by inspecting the registry contents of the Internet Registries, because the intention of this PKI is to mirror the overall state of the number registries with the set of issued certificates.

The RPKI is different from many other PKIs because the requirements related to adding digital signatures to the routing domain are different from many other PKI deployment environments. The common question that the PKI attempts to answer is: “Is this data authentic?” The data is signed with a digital signature, and the key used to generate that signature is described in a certificate. The validity of that certificate can be ascertained by using a collection of certificates and *Certificate Revocation Lists* (CRLs), such that a relying party may validate the data by using its local trust anchor(s) and constructing a *validation path* of issuer-subject chained certificates from a trust point to the digital signature. If this collection of certificates is bundled with the digital signature and the data itself, then the only data items that need to be distributed outside the data flow are the PKI trust anchors.

Distributing the RPKI Data Collection

When the RPKI is combined with a use case for the routing domain, we are looking at a design space that is somewhat atypical in the PKI world. For example, in the WebPKI the certificates that are passed between server and client in the initial exchange of a *Transport Level Security* (TLS) session are related to the particular domain name used in the TLS session being established.^[33] The critical distinction here between the secure client/server transactions using the WebPKI and the promulgation of routing information in the routing system using RPKI is that the routing system continuously presents the *entire* routing domain to each relying party. Each relying party needs to have continual access to the entire RPKI certificate and CRL collection, rather than the TLS practice of processing individual signatures and certificates on an as-needed basis.

This requirement for all participating entities to have access to all the RPKI data at all times poses a design challenge about how to manage this RPKI and use it in a routing protocol such as BGP.

A basic approach here is for each Internet Registry to publish its certificate products in its own publication point. This paradigm is analogous to the pre-*Content Delivery Network* (CDN) model of web content publication, where each element is independently published. Of course, in this case while publication is easy, the onus is shifted to the relying party client, or BGP validator, who has to assemble a local cache of all RPKI signed data. It becomes the task of clients of the RPKI to maintain a local cache of the entire RPKI by continuously sweeping across these publication points looking for, and retrieving, changes, and validating all such signed objects as they are received.

At this point BGP updates could be passed to this local RPKI engine, and the data in the update can be compared against the validated information contained in the local RPKI cache. If RPKI validation was performed at the point of acceptance into the local cache (that is, discarding all RPKI products that cannot be validated within the framework of the RPKI validation procedures), then you could verify the route information against the assembled (and validated) crypto data without a high on-demand crypto processing overhead. An alternative approach is to express the validation outcomes from the local RPKI cache as a filter list. If this list were maintained on a router, then the overheads in passing route objects through such a filter would be little different from the many other routing policy maps used in operational configurations.

The drawback in this distributed approach is the need for these clients to constantly sweep all the RPKI publication points to ensure that their local cache is up-to-date. The meaning of “up-to-date” is relative here, but it is worth remembering that the average time to propagate a BGP update across the global Internet depends on the average AS path length (around 4 to 5 autonomous systems on average at present) and the interaction with the BGP *Minimum Route Advertisement Interval* (MRAI) timers. Whereas the worst case would be 300 seconds (assuming that the full MRAI delay would be applied on each eBGP session), the fastest case is well under a second. So how quickly should the local cache be populated to keep up with the propagation of routing information in BGP? Before leaping to a target time, it is also worth remembering the scaling question. With around 100,000 distinct ASs in the Internet routing system, today’s worst-case scenario is some 100,000 RPKI clients performing a sweep across 100,000 distinct RPKI publication points every few seconds (or even more frequently if the RPKI system is intended to be highly responsive).

In some ways, this scenario puts the load on the wrong side of the information distribution process. By making the relatively infrequent publication process one that involves a local action without any associated notification of a change, then the burden is shifted to the client set, that has to poll every publication point continuously just to ascertain if anything has changed. To put it as plainly as possible, this particular information distribution design is completely broken! If the client set is known in advance (such as is the case in the DNS in synchronising the information across primary and secondary authoritative services), you could use notification mechanisms. But in the case of the RPKI system, the publishers of authoritative information have no information as to who the clients who need to be notified of a change in the part of a *Certificate Authority* (CA) of the RPKI data collection even are. Hence, notification is not a viable option in this framework.

You can mitigate these relatively formidable scaling issues by changing the publication behaviour, in a manner analogous to the way in which CDNs have improved web performance by shifting content publication models to various permutations of anycast-related models of content replication. In the context of the RPKI, these permutations could entail the use of a smaller set of RPKI publication points that many RPKI certificate issuers share, or a reduction in the number of independent CAs who each publish their own products through the extensive use of *Registration Agents* (RAs). The information being published is signed, so there is no particular benefit to retrieving the data from any particular publication point. As long as the client can validate the data, the client can be assured that the data it has retrieved is most likely to be genuine, irrespective of the location used for the retrieval. It is possible to use third-party aggregators in such a role; these aggregators would take on the task of continuous monitoring of all RPKI CA publication points and publish an aggregated data set of all current RPKI data. You could take this model further into a *push* model by having clients register their interest in updates from the intermediary and allow the intermediary to send them information updates as they are received from the primary CA publication point sources.

Again, it must be noted that the information is signed, so the potential that the intermediary could alter the RPKI information is limited. The design gap in such mediated distribution approaches is to provide a mechanism for clients of these aggregated intermediaries to be assured that the collection the intermediary has provided is the entire collection of RPKI data, and any credible intermediary approach would need to explicitly address this problem of *information completeness*.

However, although these approaches reduce the load imposed on the RPKI clients by increasing the load on information publication, such aggregated publication models also create critical points of concentration of routing data, and a sustained denial-of-service attack against such aggregate publication points could significantly affect the routing system as the local RPKI caches lose currency and coherency.

These approaches have their own strengths and risks. Highly distributed publication models impose undue costs on clients because the clients need to maintain an aggregate and current data collection in their local cache. Aggregate data-publishing models relieve load from clients but have some unresolved issues in terms of assured completeness of the aggregated data collection; they also run the risk of creating new points of vulnerability in terms of the consequence of DOS attacks launched against these aggregated publication points.

The current RPKI operational framework that is used in the *Route Origination Validation* (ROV) tool^[38] uses the approach of an *out-of-band* RPKI *pull* system together with some use of aggregated RPKI publication points. The local cache currency performance level of an RPKI client is phrased in units of minutes rather than seconds, and the overall system operates at a level of coherency that is at a time scale of hours rather than minutes. The initial design of this RPKI distribution system is for each client to operate autonomously and maintain a local cache to keep synchronised with the current state of all the RPKI publication points using the *rsync* protocol^[39] together with the concept of a manifest.^[40] This manifest allows a client to ensure that it has retrieved the entirety of the data available at each RPKI publication point. The *rsync* protocol was subsequently found to be a poor choice for this role,^[41] and these days the *RPKI Repository Delta Protocol* (RRDP) is the preferred RPKI repository synchronisation tool.^[42]

In terms of the application of the RPKI to the BGP environment, we should ask an obvious question here: If the intent of the flooding system is to provide a reliable and efficient way to flood current information to all clients, then why not just use BGP itself? BGP is an Internet-wide information flooding protocol using a *push*-based approach that is intended to ensure that all BGP speakers have a consistent and current collection of reachable route objects. If the set of clients that want to maintain an up-to-date synchronised local cache is isomorphic to the set of BGP speakers, then adding a BGP message payload type in the same manner that *Address Family Indicator/ Subsequent Address Family Indicator* (AFI/SAFI) indicators are already used in Multi-Protocol BGP today seems only logical.

Part of the reason why the RPKI has had to re-invent this particular wheel of reliable flooding lies in the strictures imposed on the standardisation effort in the IETF, where the *Secure Inter-Domain Routing* (SIDR) Working Group was constrained from proposing changes to the BGP protocol itself.

In retrospect, this constraint appears to have been a rather suboptimal and, in hindsight, extremely poor piece of guidance from the *Internet Engineering Steering Group* (IESG) at the time.

If we could contemplate changes to BGP, then one approach to the RPKI distribution tasks is to maintain the association of the validation material with the data, and in the context of the routing environment it would staple a collection of certificates (and CRLs) to each route object. In a sense this approach would attempt to reproduce the TLS model in BGP, where each prefix being updated would have a subset of the RPKI certificates stapled to the update that would permit an associated signed attestation to be validated within the framework of the RPKI. This method is not without additional impositions, and it would impose costs on the operation of the BGP protocol and BGP speakers. Stapling crypto credentials to BGP updates would bloat both the volume of stapled data (through the use of long validation chained paths and long-term certificate issuance policies which, in turn, create extended CRL lists) and the amount of crypto processing of these stapled digital credentials. There would be a significant level of the retransmission of certificates on a pair-wise basis in such a system if the protocol were to bundle the entire RPKI validation chain data with every routing protocol update. The validation processing load would also likely be beyond the processing capabilities of most routers, and there are considerations of the maximum message size in the BGP protocol itself (which, until RFC 8654,^[13] published in October 2019, was 4,096 octets), which limited the amount of attached data that can be placed into BGP.

None of these issues is intractable, and many proposals have been made to attempt to optimise such additional loads and processing demands. We will look at some of these proposals in Part 2 of this survey.

Coming in Part 2

In Part 2, we will take these various requirements and tools and look at the various proposals that have been published for securing BGP. We will also evaluate the current state of the effort in the IETF to standardise a secure BGP Framework.

References

- [1] Yakov Rekhter, Susan Hares, and Tony Li, “A Border Gateway Protocol 4 (BGP-4),” RFC 4271, January 2006.
- [2] Yakov Rekhter, “Experience with the BGP Protocol,” RFC 1266, October 1991.
- [3] Sandra Murphy, “BGP Security Vulnerabilities Analysis,” RFC 4272, January 2006.
- [4] Abbie Barbir, Sandra Murphy, and Yi Yang, “Generic Threats to Routing Protocols,” RFC 4593, October 2006.

- [5] Hitesh Ballani, Paul Francis, and Xinyang Zhang, "A study of prefix hijacking and interception in the Internet," *ACM SIGCOMM Computer Communications Review*, Volume 37, No. 4, October 2007.
- [6] Kirk Lougheed and Yakov Rekhter, "Border Gateway Protocol (BGP)," RFC 1105, June 1989.
- [7] Kirk Lougheed and Yakov Rekhter, "Border Gateway Protocol (BGP)," RFC 1163, June 1990.
- [8] Kirk Lougheed and Yakov Rekhter, "Border Gateway Protocol 3 (BGP-3)," RFC 1267, October 1991.
- [9] Yakov Rekhter and Tony Li, "A Border Gateway Protocol 4 (BGP-4)," RFC 1771, March 1995.
- [10] Richard Bellman, "On a routing problem," *Quarterly of Applied Mathematics*, Volume 16, No.1, April 1958.
- [11] Lester Randolph Ford, Jr., "Network Flow Theory," RAND Corporation, Paper P-923, August 1956.
<https://apps.dtic.mil/sti/pdfs/AD0422842.pdf>
- [12] Geoff Huston, *The BGP Report*,
<https://bgp.potaroo.net>
- [13] Keyur Patel, David Ward, and Randy Bush, "Extended Message Support for BGP," RFC 8654, October 2019.
- [14] Tim Griffin and Geoff Huston, "BGP Wedgies," RFC 4264, November 2005.
- [15] Feng Wang and Lixin Gao, "On Inferring and Characterizing Internet Routing Policies," in *IMC '03: Proceedings of the 3rd ACM SIGCOMM Conference on Internet Measurement*, October 2003.
- [16] Mitesh Dalal, Randall R. Stewart, and Anantha Ramaiah, "Improving TCPv2's Robustness to Blind In-Window Attacks," RFC 5961, August 2010.
- [17] Curtis Villamizar, Ravi Chandra, and Ramesh Govindan, "BGP Route Flap Damping," RFC 2439, November 1998.
- [18] RIPE Routing Working Group Recommendations on Route Flap Damping, January 2013.
<https://www.ripe.net/publications/docs/ripe-580>
- [19] Kotikalapudi Sriram, Doug Montgomery, Oliver Borchert, Okhee Kim, and D. Richard Kuhn, "Study of BGP Peering Session Attacks and Their Impacts on Routing Performance," *IEEE Journal on Selected Areas in Communications*, Volume 24, No. 10, Oct. 2006.

- [20] Ratul Mahajan, David Wetherall, and Tom Anderson, "Understanding BGP Misconfiguration," in *SIGCOMM '02: Proceedings of the 2002 Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications*, 2002.
- [21] Sharon Goldberg, Shai Halevi, Aaron D. Jaggard, Vijay Ramachandran, and Rebecca Wright, "Rationality and Traffic Attraction: Incentives for Honest Path Announcements in BGP," *ACM SIGCOMM Computer Communications Review*, Volume 38, No. 4, August 2008.
- [22] Blaine Christian and Tony Tauber, "BGP Security Requirements," November 2008, Internet Draft, work in progress, November 2008, **draft-ietf-rpsec-bgpsecrec-10**
- [23] Xinming He, Christos Papadopoulos, and Pavlin Radoslavov, "A framework for incremental deployment strategies for router-assisted services," in *Proceedings IEEE INFOCOM 2003: Twenty-Second Annual Joint Conference of the IEEE Computer and Communications Societies*, Volume 2, March 2003.
- [24] Martin Suchara, Ioannis Avramopoulos, and Jennifer Rexford, "Securing BGP Incrementally," in *CoNEXT '07: Proceedings of the 2007 ACM CoNEXT Conference*, December 2007.
- [25] Jennifer Rexford and Joan Feigenbaum, "Incrementally-Deployable Security for Interdomain Routing," in *CATCH '09: Cybersecurity Applications & Technology Conference for Homeland Security*, March 2009.
- [26] Vijay Gill, John Heasley, and David Meyer, "The Generalized TTL Security Mechanism (GTSM)," RFC 3682, February 2004.
- [27] Vijay Gill, John Heasley, David Meyer, Pekka Savola, and Carlos Pignataro, "The Generalized TTL Security Mechanism (GTSM)," RFC 5082, October 2007.
- [28] Wesley M. Eddy, "TCP SYN Flooding Attacks and Common Mitigations," RFC 4987, August 2007.
- [29] Karen Seo and Stephen Kent, "Security Architecture for the Internet Protocol," RFC 4301, December 2005.
- [30] Brian Weis, "Why IPsec and BGP don't play well together in real networks," Security Area Working Group presentations, IETF 66, July 2006.
<https://www.ietf.org/proceedings/66/slides/saag-2.pdf>
- [31] Andy Heffernan, "Protection of BGP Sessions via the TCP MD5 Signature Option," RFC 2385, August 1998.
- [32] Ronald P. Bonica, Allison Mankin, and Joe Touch, "The TCP Authentication Option," RFC 5925, June 2010.
- [33] Eric Rescorla, "The Transport Layer Security (TLS) Protocol Version 1.3," RFC 8446, August 2018.

- [34] Jana Iyengar and Martin Thomson, “QUIC: A UDP-Based Multiplexed and Secure Transport,” RFC 9000, May 2021.
- [35] Charles Lynn, Karen Seo, and Stephen Kent, “X.509 Extensions for IP Addresses and AS Identifiers,” RFC 3779, June 2004.
- [36] Geoff Huston, Robert Loomans, and George Michaelson, “A Profile for X.509 PKIX Resource Certificates,” RFC 6487, February 2012.
- [37] Matt Lepinski, Andrew Chi, and Stephen Kent, “Signed Object Template for the Resource Public Key Infrastructure (RPKI),” RFC 6488, February 2012.
- [38] Geoff Huston and George Michaelson, “Validation of Route Origination Using the Resource Certificate Public Key Infrastructure (PKI) and Route Origin Authorizations (ROAs),” RFC 6483, February 2012.
- [39] Andrew Tridgell and Paul Mackerras, “The rsync algorithm,” 1996.
<https://rsync.samba.org/>
- [40] Matt Lepinski, Stephen Kent, Geoff Huston, and Rob Austein, “Manifests for the Resource Public Key Infrastructure (RPKI),” RFC 6486, February 2012.
- [41] George Michaelson, and Byron Ellacott, “rsync Considered Inefficient and Harmful,” IETF 89, March 2014.
<https://www.ietf.org/proceedings/89/slides/slides-89-sidr-6.pdf>
- [42] Rob Austein, Tim Bruijnzeels, Bryan Weber, and Oleg Muravskiy, “The RPKI Repository Delta Protocol (RRDP),” RFC 8182, July 2017.
- [43] Geoff Huston “A Quick Look at QUIC,” *The Internet Protocol Journal*, Volume 22, No. 1, March 2019.
- [44] Tom Strickx and Celso Martinho, “Understanding How Facebook Disappeared from the Internet,” *The Cloudflare Blog*, October 4, 2021.
<https://blog.cloudflare.com/october-2021-facebook-outage/>
- [45] Mark Handley, “Why the Internet only just works,” *BT Technology Journal*, Volume 24, No. 3, July 2006.

GEOFF HUSTON, B.Sc., M.Sc. A.M., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Kjetil Aas	Gareth Bryan	Joan Marc Riera	Geert Jan de Groot	Anders Marius
Fabrizio Accatino	Stefan Buckmann	Duocastella	Christopher Guemez	Jørgensen
Michael Achola	Caner Budakoglu	Pedro Duque	Gulf Coast Shots	Merike Kao
Martin Adkins	Darrell Budic	Holger Durer	Sheryll de Guzman	Andrew Kaiser
Melchior Aelmans	BugWorks	Mark Eanes	Rex Hale	Christos Karayiannis
Christopher Affleck	Scott Burleigh	Andrew Edwards	Jason Hall	Daniel Karrenberg
Scott Aitken	Chad Burnham	Peter Robert Egli	Darow Han	David Kekar
Jacobus Akkerhuis	Jon Harald Bøvre	George Ehlers	Handy Networks LLC	Stuart Kendrick
Antonio Cuñat Alario	Olivier Cahagne	Peter Eisses	James Hamilton	Robert Kent
Nicola Altan	Antoine Camerlo	Torbjörn Eklöv	Stephen Hanna	Jithin Kesavan
Shane Amante	Tracy Camp	Y Ertur	Martin Hannigan	Jubal Kessler
Marcelo do Amaral	Ignacio Soto Campos	ERNW GmbH	John Hardin	Shan Ali Khan
Matteo D'Ambrosio	Fabio Caneparo	ESdatCo	David Harper	Nabeel Khatri
Selva Anandavel	Roberto Canonico	Steve Esquivel	Edward Hauser	Dae Young Kim
Jens Andersson	David Cardwell	Jay Etchings	David Hauweele	William W. H. Kimandu
Danish Ansari	John Cavanaugh	Mikhail Evstiounin	Marilyn Hay	John King
Finn Arildsen	Lj Cemerar	Bill Fenner	Headcrafts SRLS	Russell Kirk
Tim Armstrong	Dave Chapman	Paul Ferguson	Hidde van der Heide	Gary Klesk
Richard Artes	Stefanos Charchalak	Ricardo Ferreira	Johan Helsingius	Anthony Klopp
Michael Aschwanden	Greg Chisholm	Kent Fichtner	Robert Hinden	Henry Kluge
David Atkins	David Chosrova	Armin Fisslthaler	Asbjørn Højmark	Michael Kluk
Jac Backus	Marcin Cieslak	Michael Fiumano	Damien Holloway	Andrew Koch
Jaime Badua	Lauris Cikovskis	The Flirble Organisation	Alain Van Hoof	Ia Kochiashvili
Bent Bagger	Guido Coenders	Gary Ford	Edward Hotard	Carsten Koempe
Eric Baker	Brad Clark	Jean-Pierre Forcioli	Bill Huber	Richard Koene
Santosh Balagopalan	Narelle Clark	Susan Forney	Hagen Hultzs	Alexander Kogan
William Baltas	Horst Clausen	Christopher Forsyth	Kauto Huopio	Antonin Kral
David Bandinelli	Joseph Connolly	Andrew Fox	Kevin Iddles	Robert Krejčí
Benjamin Barkin-	Steve Corbató	Craig Fox	Mika Ilvesmaki	Mathias Körber
Wilkins	Brian Courtney	Fausto Franceschini	Karsten Iwen	John Kristoff
Feras Batainah	Beth and Steve Crocker	Valerie Fronczak	David Jaffe	Terje Krogdahl
Michael Bazarewsky	Dave Crocker	Tomislav Futivic	Ashford Jaggernaut	Bobby Krupczak
David Belson	Kevin Croes	Laurence Gagliani	Martijn Jansen	Murray Kuchera
Hidde Beumer	John Curran	Edward Gallagher	Jozef Janitor	Warren Kumari
Pier Paolo Biagi	André Danthine	Andrew Gallo	John Jarvis	George Kuo
Tyson Blanchard	Morgan Davis	Chris Gamboni	Dennis Jennings	Dirk Kurfuerst
John Bigrow	Jeff Day	Xosé Bravo Garcia	Edward Jennings	Darrell Lack
Orvar Ari Bjarnason	Julien Dhallenne	Osvaldo Gazzaniga	Aart Jochem	Andrew Lamb
Axel Boeger	Freek Dijkstra	Kevin Gee	Nils Johansson	Richard Lamb
Keith Bogart	Geert Van Dijk	Greg Giessow	Brian Johnson	Yan Landriault
Mirko Bonadei	David Dillow	John Gilbert	Curtis Johnson	Edwin Lang
Roberto Bonalumi	Richard Dodsworth	Serge Van Ginderachter	Richard Johnson	Sig Lange
Julie Bottorff	Ernesto Doelling	Greg Goddard	Jim Johnston	Markus Langenmair
Photography	Michael Dolan	Tiago Goncalves	Jonatan Jonasson	Fred Langham
Gerry Boudreaux	Eugene Doroniuk	Ron Goodheart	Daniel Jones	Tracy LaQuey Parker
L de Braal	Karlheinz Dölger	Octavio Alfageme	Gary Jones	Jose Antonio Lazaro
Kevin Breit	Joshua Dreier	Gorostiaga	Jerry Jones	Lazaro
Thomas Bridge	Lutz Drink	Barry Greene	Michael Jones	Rick van Leeuwen
Ilia Bromberg	Aaron Dudek	Jeffrey Greene	Amar Joshi	Simon Leinen
Václav Brožík	Dmitriy Dudko	Richard Gregor	Javier Juan	Robert Lewis
Christophe Brun	Andrew Dul	Martijn Groenleer	David Jump	Christian Liberale

Martin Lillepuu	Roberto Montoya	Andrew Potter	Timothy Schwab	Lorin J Thompson
Roger Lindholm	Charles Monson	Eduard Llull Pou	Roger Schwartz	Fabrizio Tivano
Link Light Networks	Andrea Montefusco	Tim Pozar	SeenThere	Peter Tomsu Fine Art
Sergio Loreti	Fernando Montenegro	David Raistrick	Scott Seifel	Photography
Eric Louie	Joel Moore	Priyan R Rajeevan	Yury Shefer	Joseph Toste
Adam Loveless	John More	Balaji Rajendran	Yaron Sheffer	Rey Tucker
Josh Lowe	Maurizio Moroni	Paul Rathbone	Doron Shikmoni	Sandro Tumini
Guillermo a Loyola	Brian Mort	William Rawlings	Tj Shumway	Angelo Turetta
Hannes Lubich	Soenke Mumm	Mujtiba Raza Rizvi	Jeffrey Sicuranza	Phil Tweedie
Dan Lynch	Tariq Mustafa	Bill Reid	Thorsten Sideboard	Steve Ulrich
David MacDuffie	Stuart Nadin	Petr Rejhon	Greipur Sigurdsson	Unitek Engineering AG
Sanya Madan	Michel Nakhla	Robert Remenyi	Fillipe Cajaiba da Silva	John Urbanek
Miroslav Madić	Mazdak Rajabi Nasab	Rodrigo Ribeiro	Andrew Simmons	Martin Urwaleck
Alexis Madriz	Krishna Natarajan	Glenn Ricart	Pradeep Singh	Betsy Vanderpool
Carl Malamud	Naveen Nathan	Justin Richards	Henry Sinnreich	Surendran Vangadasalam
Jonathan Maldonado	Darryl Newman	Rafael Riera	Geoff Sisson	Ramnath Vasudha
Michael Malik	Thomas Nikolajsen	Mark Risinger	Helge Skrivervik	Philip VENABLE
Tarmo Mämers	Paul Nikolich	Fernando Robayo	Terry Slattery	Buddy Venne
Yogesh Mangar	Travis Northrup	Gregory Robinson	Darren Sleeth	Alejandro Venera
Bill Manning	Marijana Novakovic	Ron Rockrohr	Richard Smit	Luca Ventura
Harold March	David Oates	Carlos Rodrigues	Bob Smith	Scott Vermillion
Vincent Marchand	Ovidiu Obersterescu	Magnus Romedahl	Courtney Smith	Tom Vest
Normando Marcolongo	Tim O'Brien	Lex Van Roon	Eric Smith	Vista Global Coaching
Gabriel Marroquin	Mike O'Connor	Alessandra Rosi	Mark Smith	& Consulting
David Martin	Mike O'Dell	David Ross	Tim Sneddon	Dario Vitali
Jim Martin	John O'Neill	William Ross	Craig Snell	Jeffrey Wagner
Ruben Tripiana Martin	Jim Oplotnik	Boudhayan	Job Snijders	Don Wahl
Timothy Martin	Packet Consulting	Roychowdhury	Ronald Solano	Michael L Wahrman
Carles Mateu	Limited	Carlos Rubio	Asit Som	Laurence Walker
Juan Jose Marin	Carlos Astor Araujo	Rainer Rudigier	Ignacio Soto Campos	Randy Watts
Martinez	Palmeira	Timo Rüter	Evandro Sousa	Andrew Webster
Ioan Maxim	Alexis Panagopoulos	RustedMusic	Peter Spekrijse	Tim Weil
David Mazel	Gaurav Panwar	Babak Saberi	Thayumanavan Sridhar	Jd Wegner
Miles McCredie	Manuel Uruena Pascual	George Sadowsky	Paul Stancik	Westmoreland
Brian McCullough	Ricardo Patara	Scott Sandefur	Ralf Stempfner	Engineering Inc.
Joe McEachern	Dipesh Patel	Sachin Sapkal	Matthew Stenberg	Rick Wesson
Alexander McKenzie	Alex Parkinson	Arturas Satkovskis	Adrian Stevens	Peter Whimp
Jay McMaster	Craig Partridge	PS Saunders	Clinton Stevens	Russ White
Mark Mc Nicholas	Dan Paynter	Richard Savoy	John Streck	Jurrien Wijnhuizen
Olaf Mehlberg	Leif Eric Pedersen	John Sayer	Martin Streule	Derick Winkworth
Carsten Melberg	Rui Sao Pedro	Phil Scarr	David Strom	Pindar Wong
Kevin Menezes	Juan Pena	Gianpaolo	Colin Strutt	Makarand Yerawadekar
Bart Jan Menkveld	Chris Perkins	Scassellati	Viktor Sudakov	Phillip Yialeloglou
Sean Mentzer	Michael Petry	Elizabeth Scheid	Edward-W. Suor	Janko Zavernik
William Mills	Alexander Peuchert	Jeroen Van Ingen	Vincent Surillo	Bernd Zeimet
David Millsom	David Phelan	Schenau	Terence Charles	Muhammad Ziad
Desiree Miloshevic	Derrell Piper	Carsten Scherb	Sweetser	Ziauddin
Joost van der Minnen	Rob Pirnie	Ernest Schirmer	T2Group	Tom Zingale
Thomas Mino	Marc Vives Piza	Philip Schneck	Roman Tarasov	Jose Zumalave
Rob Minshall	Jorge Ivan Pincay Ponce	Peter Schoo	David Theese	Romeo Zwart
Wijnand Modderman	Victoria Poncini	Dan Schrenk	Douglas Thompson	廖明沂.
Mohammad Moghaddas	Blahoslav Popela	Richard Schultz	Kerry Thompson	



Follow us on Twitter and Facebook

@protocoljournal



<https://www.facebook.com/newipj>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Internet
Society



Diamond Sponsors

Your logo here!

Ruby Sponsors



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

John Crain, Chief Security, Stability and Resilience Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Principal – Internet Technology, Policy, and Advocacy
The Internet Society

Dr. Jun Murai, Founder, WIDE Project
Distinguished Professor, Keio University
Co-Director, Keio University Cyber Civilization Research Center, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

April 2022

Volume 25, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

FROM THE EDITOR

In This Issue

From the Editor	1
Securing Inter-Domain Routing.....	2
Fragments	38
Thank You!	40
Call for Papers	42
Supporters and Sponsors	43

I have spent some time in recent months studying the history and development of the world-wide telephone network. Broadly speaking, the telephone network has evolved in two directions away from the traditional system of interconnected public and private telephone switches and their associated hard-wired telephones. First, starting in the mid-1980s we saw the introduction of mobile devices and networks, eventually leading to what we refer to as “smartphones” today. Secondly, many of the traditional telephone networks have been augmented or completely replaced by numerous systems that employ *Voice-over Internet Protocol* (VoIP) technologies. In spite of the differences in technologies, it is still possible to place and receive voice calls to *telephone numbers*, an addressing system that has proved remarkably resilient to growth and technological evolution since its introduction some 130 years ago. (The first commercial telephone exchange was installed in 1892 in La Porte, Indiana.) Unlike IP addresses, telephone numbers are not fixed-length, nor are they managed by a single global entity, but for such a system to work we do rely on a unique set of country codes and numerous interconnection agreements, thus there are some similarities to the way the Internet operates.

I have also been reading numerous recent postings to the “internet-history” e-mail list, operated by The Internet Society. If you’re interested in hearing from Internet pioneers such as Vint Cerf, Brian Carpenter, Noel Chiappa, Jack Haverty, and many others, this list is a great place to start. You can find further details here:

<https://elists.isoc.org/mailman/listinfo/internet-history>

In our previous issue, Geoff Huston presented Part 1 of “A Survey on Securing Inter-Domain Routing.” He described the design and operation of the *Border Gateway Protocol* (BGP), the threat model, and the requirements from a security framework for BGP. In this issue, Geoff concludes the survey by looking at the various proposals to add security to the routing environment and evaluates the current state of the effort in the *Internet Engineering Task Force* (IETF) to provide a standard specification of the elements of a secure BGP framework.

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

A Survey on Securing Inter-Domain Routing

Part 2 – Approaches to Securing BGP

by Geoff Huston, APNIC

The *Border Gateway Protocol* (BGP) is the inter-domain routing protocol of the Internet, and after some thirty years of operation it is now one of the more venerable of the core protocols on the Internet. One of the major ongoing concerns related to BGP is its lack of effective security measures, and as a result the routing infrastructure of the Internet continues to be vulnerable to various forms of attack.

In Part 1 we looked at the design of BGP, the threat model, and the requirements from a security framework for BGP. In Part 2 we will look at the various proposals to add security to the routing environment and also evaluate the current state of the effort in the *Internet Engineering Task Force* (IETF) to provide a standard specification of the elements of a secure BGP framework.

The approaches to securing BGP can be further classified in the same fashion as the security requirements: securing the operation of BGP and securing the integrity of the BGP data.

Securing the Operation of BGP Sessions

BGP uses a long-held *Transmission Control Protocol* (TCP) session, and you can use the same approaches to secure any TCP session^[1] in the context of a BGP session. These approaches fall into two categories: those that simply attempt to protect the TCP session from disruption via injection of spurious traffic, and those that also attempt to protect the TCP session from eavesdropping and alteration by encrypting the payload.

Generalized TTL Security Mechanism

The *Generalized TTL Security Mechanism* (GSTM), originally described in [2] and updated in [3], is based on the observation that the overall majority of BGP peering sessions are established between routers that are directly connected. The technique is to configure each BGP IP packet to be sent with a *Time To Live* (TTL) field value in the IP header of 255, and for the BGP receiver to discard all packets with an inbound TTL of less than a set threshold value. For a direct connection, the inbound TTL value should be 255, so the receiver can discard all inbound TCP packets within this session with a TTL of 254 or less.

The motivation for this approach is that spoofing of the TTL field in an IP header is challenging for an unassisted remote attacker. This TTL packet filter is a lightweight defensive measure intended to add some protection to the BGP session from efforts to intrude into the session using remote attacks.

You can use this GTSM approach for multi-hop BGP peer sessions, as well as directly connected BGP sessions, but it is not all that robust in terms of its security properties because of the additional variables introduced with TTL changes due to routing changes and the potential to mask the conventional TTL behaviour with tunnelling techniques.

TCP MD5 Signature Option

A more robust approach to protecting the TCP session is through the use of cryptographic protection of the TCP session. While these crypto approaches can be highly resilient to intrusion attempts, they also expose the BGP speaker to potential *Denial of Service* (DoS) attacks if the processing load of the cryptographic functions to detect bogus packets is sufficiently high. The target still has to process bogus packets just to ascertain that they are bogus.

The TCP MD5 Signature Option^[4] uses message authentication codes—which are a class of cryptographic hash algorithms applied to messages of arbitrary length that produce a *message digest* of the message—intended to protect the integrity of the message. The desired property of a message digest is that it is infeasible to generate two messages that have the same message-digest value, and equally infeasible to generate a new message that has a particular message digest value. The *Message-Digest 5* (MD5) algorithm^[5] is intended for digital signature applications where a message digest is generated over the combination of a message and a secret shared key value. The message and the digest value can be transmitted openly, and the receiver can use a local copy of the secret key and apply the message-digest algorithm to the combination of the received message and the key. If the digest value matches the received value, then the receiver can be assured that the message has not been altered in transit, and that the message was generated by a party who also has knowledge of the key.

The TCP MD5 Signature Option is a TCP extension where each TCP segment contains a TCP option that contains the 128-bit MD5 digest of the combination of the TCP pseudo header, the TCP segment payload excluding TCP options, and a connection-specific key. This combination establishes a cryptographically secure signature of the packet. Without knowing the key, it is very challenging to construct a TCP segment with a valid signature, and it is not readily possible to alter the packet without causing the signature to be invalidated. The receiver calculates the MD5 digest across the received data, using a locally held copy of the key, and rejects the segment if the digest value fails to match that provided in the packet. In the context of BGP, the TCP session is resistant to various forms of intrusion attack unless the attacker has knowledge of the shared secret key value. The TCP MD5 specification does not specify how the shared key is passed between the two BGP speakers, nor how the key value can be changed during the session. This latter problem is significant in that continued use of a key weakens its integrity, and it is conventionally advised that MD5 session keys be changed every 90 days or so in this type of use context^[6].

With a mechanism for in-band key change, this advice implies the need for a BGP session reset every 90 days or so, which is counter to conventional operational practice in BGP, where sessions are held up for as long as possible. Even with tools such as BGP *Graceful Restart*, deliberate BGP session resets are generally avoided in the operational community.

TCP Authentication Option

A somewhat different approach—the *TCP Authentication Option* (AO)^[7]—uses a *Message Authentication Field* in the place of the MD5 message digest, where the final bit of the length field of the option determines whether or not a key ID has been appended to the *Message Authentication Code* (MAC). The message-digest algorithm in this case is specified as HMAC-MD5-96, although you can use other algorithms if you configure them in advance. This approach relies on a similar form of out-of-band provisioning as the original MD5 approach, where each end of the conversation must configure a *TCP Security Association Database* before using this mechanism. This database contains a description of the supported TCP connections, the key set, the MAC algorithm, and MAC length.

IPSec

Internet Protocol Security (IPsec) is a suite of protocols that operate at the IP level of the protocol stack; these protocols secure all communications between two endpoints^[8]. The functionality of IPsec includes methods for protection of IP packet headers, methods for protection and encryption of IP payloads, and key management services that allow key rollover during long sessions. This implementation is one of public/private key cryptography, and it can ensure the confidentiality and integrity of all IP messages passed between two hosts. You can use IPsec to secure BGP sessions, and it provides greater levels of assurance than MD5 offers.

However, IPsec is not widely used in the public Internet for the purpose of securing BGP sessions^[7,9], and no generally accepted profile of IPsec for BGP has been standardised so far, with earlier efforts along these lines not progressing within the standards process. The perceived problem with IPsec relates to the complications for rekeying *Internet Key Exchange* (IKE)/IPsec sessions, and the observation that processing load to detect bogus packets is considerably higher with IPsec than with MD5. Using IPsec for BGP exposes a DoS attack where a stream of bogus IPsec packets directed at a BGP speaker may be capable of exercising the processor into a fully saturated mode of operation, causing degradation of other concurrent router functions.

More Options

As was observed in Part 1 of this survey, there are many alternatives here, including *Transport Layer Security* (TLS)^[80] and *Quick UDP Internet Connections* (QUIC)^[81], but more choice is not a substitute for better quality.

These session-level encryption approaches that applications use provide no better answer to dynamic rekeying, and they follow a now well-established Internet tradition of adding more options to divert attention from the observation that the common fundamental problems are inadequately addressed. The design goal of such application-level session approaches is protection for transient short-duration sessions, while the vulnerabilities associated with long-held BGP sessions are somewhat different.

The best advice today is that a combination of TCP AO and GTSM is as good as it gets at present. However, it's also highly desirable to avoid multi-hop BGP wherever possible and directly attach the two BGP speakers. That way reduces considerably the radius of potential eavesdroppers and attackers.

Securing the Integrity of Routing Information Passed in BGP

One of the earlier recognised works that addressed routing security was the 1988 study on *Byzantine Robustness* by Radia Perlman^[9]. If failure or malicious behaviour on the part of one or more entities in the system occurs, all correctly operating entities should reach a mutually consistent decision regarding the validity of each message in finite time. This study was in the area of link-state protocol design, and the work described a protocol that satisfied the properties for Byzantine Robustness. It categorised route validation in three approaches:

- *Bound* or just in time — validation occurs the same moment a route is announced, and appropriate measures are taken immediately. Credentials must be available immediately.
- *Unbound* or just in case — validation occurs only if a new router takes part in the system. Credentials are retrieved on arrival of this router.
- *Interrogative* or just too late — validation occurs sporadically, requesting validation or credentials from a remote system when necessary.

Although the link-state approach described in this paper does not exactly match the inter-domain routing environment, the concept of validation of routing information is a consistent theme in all BGP security architectures.

Subsequent work by Smith and Garcia-Luna-Aceves^[10,11], published in 1996, attempts to address session security by modifying the BGP protocol. This work proposed the protection of BGP control messages using message encryption at the BGP level, with session keys exchanged at BGP session establishment time. It also proposed the addition of a message sequence number to protect against replay attacks and message removal. This approach also proposed a predecessor path attribute that indicated the *Autonomous System* (AS) prior to the destination AS for the current route and proposed digitally signing all fixed fields in the UPDATE message. The predecessor attribute constructs a means of validation of the AS Path attribute.

These proposed changes to the BGP protocol required comprehensive adoption and deployment in order to be effective, because partial adoption would create gaps in any assurance that a predecessor attribute could provide. Their approach was similar to the earlier *Interdomain Routing Protocol* (IDRP) work^[12]. IDRP eschewed the use of TCP and included a reliable flow-controlled transport into the IDRP protocol, also including numerous message integrity protection options.

A contemporary proposal to the Smith and Garcia-Lunes-Aceves proposal for securing BGP was based on leaving the BGP protocol unchanged, but augmenting the BGP data flow with access to credential information. This additional information was intended to allow a BGP speaker to confirm the authenticity of origination information in BGP UPDATE messages by validating the binding of address prefixes to originating ASes^[13]. This proposal, *Network Layer Reachability Information* (NLRI) *Origin AS Verification*, used the *Domain Name System* (DNS) as the distribution mechanism for origination information, where a BGP speaker could perform a DNS query to validate the prefix size and authorised originating AS information contained in a BGP route object. Informally, it was intended to allow a DNS query to answer the question: “Which ASes have been authorised by the address holder to originate a route for this prefix?” The proposed framework assumed that the reverse DNS space was securely associated with the holder of the address prefix, and the DNS response was verifiable [using a *Domain Name System Security Extensions* (DNSSEC)-signed DNS record and DNSSEC validation^[14], presumably, although this work was contemporaneous with DNSSEC and did not use it in this proposal]. This proposal assumed that the performance of DNS queries was within the same order of timescale as the propagation of BGP messages within BGP. It also assumed that there was no circularity, where a DNS recursive resolver or authoritative name server that the BGP speaker used was located within an address prefix that was being validated prior to local acceptance of the route associated with that prefix.

The DNS delegation hierarchy would need to be precisely aligned to the address allocation framework, so that the zone administrator of each of these origination authentication zones was in fact the duly delegated holder of the addresses, and this alignment should, preferably, be capable of third-party validation. Meeting these requirements would create a digital signature hierarchy embedded in the DNS that would be aligned to the address allocation framework.

The *Internet Routing Registry* (IRR)^[82] pre-dates most other efforts in this space, dating back to the routing work of the early 1990s in the *Routing Arbiter* project that was part of the US *National Science Foundation Network* (NSFNET), and a project coordinated under the auspices of the RIPE IRR in Europe. The IRR objective was to provide a set of routing policy databases populated by the ASes themselves that described the addresses that they intended to announce in the routing system and the routing policies that they intended to apply to these announcements^[83].

The Routing Registry was a response to the need described in RFC 1787^[84] for improving global consistency by allowing providers to share routing policies. Each participating AS submits policy data, encoded using the *Routing Policy Specification Language* (RPSL)^[27,28]. Clients may use the registry to determine the stated policies for a particular AS, including what ASes (and possibly prefixes) are suitable for import or export, potentially using the data to populate filter sets on their BGP feeds. Additional information that an AS provides to the IRR could include policy concerning the configuration of BGP communities and the policy responses associated with particular community settings.

However, the utility of the IRR for securing routing is quite limited. First, the IRR does not provide information about current routes, only about potential routes. Some potential routes may be legal according to the IRR, but undesirable from a more global point of view. Next, the IRR has many security vulnerabilities concerning the integrity of registry contents and authorization of changes to the registry. There is no intrinsic authority model that constrains which party can publish data about addresses and ASes in an IRR. Moreover, some policy information concerning agreements between peering ASes is not intended for broader public distribution and the IRRs did not normally implement any form of limited disclosure rules. Efforts to improve the controls over the authority framework in registries and access frameworks^[85] never really gained traction. The IRR system is a misnomer, in that there is not a single IRR but many IRRs. The contents of these IRRs are not necessarily mutually consistent, and there is no clear way to resolve any such conflicts. Not only is there no authority model ensuring that only authorised parties may publish routing policy data about their own address prefixes and ASes, but there is also no way to describe the intended lifetime of the information. Old information that is no longer current or relevant sits alongside current information, and this current information sits along with contingency information that may never be actually used.

Although the overall approach of providing an out-of-band commentary on routing, enumerating all the cases of authorised (or valid) route objects has been a useful tool for many operational environments, IRR tools are only truly useful in the context of being able to detect and filter routing anomalies if the information is verifiable and authentic, current, and complete. In other words, IRRs are most useful if they are carefully and continuously managed, and the accuracy and usefulness of the information rapidly declines if the information in the registry is neglected. Our experience with IRRs suggests that it would be somewhat foolhardy to automatically apply IRR data to populate route filters, given the risks of incorrect outcomes—both positive and negative. In addition, although there have been good counter examples in some operational communities, the broader judgement for IRRs being capable of supporting a robust whole-of-Internet role for route integrity is somewhat negative^[86].

It looks like the common requirements in this space appear to relate to *authenticity*, *currency*, and *completeness*.

Digital signatures can provide strong assurance related to authenticity and currency of information, assuming that the enrollment practice that governs the authority to generate such signatures is robust. Given such a practice, the consequent observation is that whether or not this digital signature framework is placed into the DNS via a DNSSEC framework^[15] or into a framework of X.509 certificates and an associated *Public Key Infrastructure* (PKI) is, at one level, an isomorphic transform of the same information. The issue of the choice of DNS (and DNSSEC) or X.509 certificates (and certificate-based validation) is then an issue of the performance requirements of these systems.

Completeness is a more challenging requirement. The identification of invalid routing information in the partial adoption case of this approach is unclear. When a query to an information source has a negative response, it is unclear whether the route object that was the basis of the query is not valid (such as a bogus prefix or a bogus AS), or the database being queried is incomplete.

Let's now move forward in time to review some more recent proposals to secure BGP.

Secure BGP

Secure BGP (sBGP)^[16] offered a relatively complete approach to securing the BGP protocol by placing digital signatures over the address and AS Path information contained in routing advertisements, and defining an associated PKI for validation of these signatures.

sBGP defines the “correct” operation of a BGP speaker in terms of a set of constraints placed on individual protocol messages, including ensuring that:

- No protocol UPDATE messages have been altered in transit between the BGP peers
- The UPDATE messages were sent by the indicated peer
- The UPDATE messages contain more recent information than has been previously sent to this BGP speaker from the peer
- The UPDATE was intended to be received by this BGP speaker
- The peer is authorised to advertise information on behalf of the peer AS

In addition, for every prefix and its originating AS, the prefix must be a validly allocated prefix, and the prefix “right-of-use” holder must have authorised the advertisement of the prefix and the originating AS to advertise the prefix.

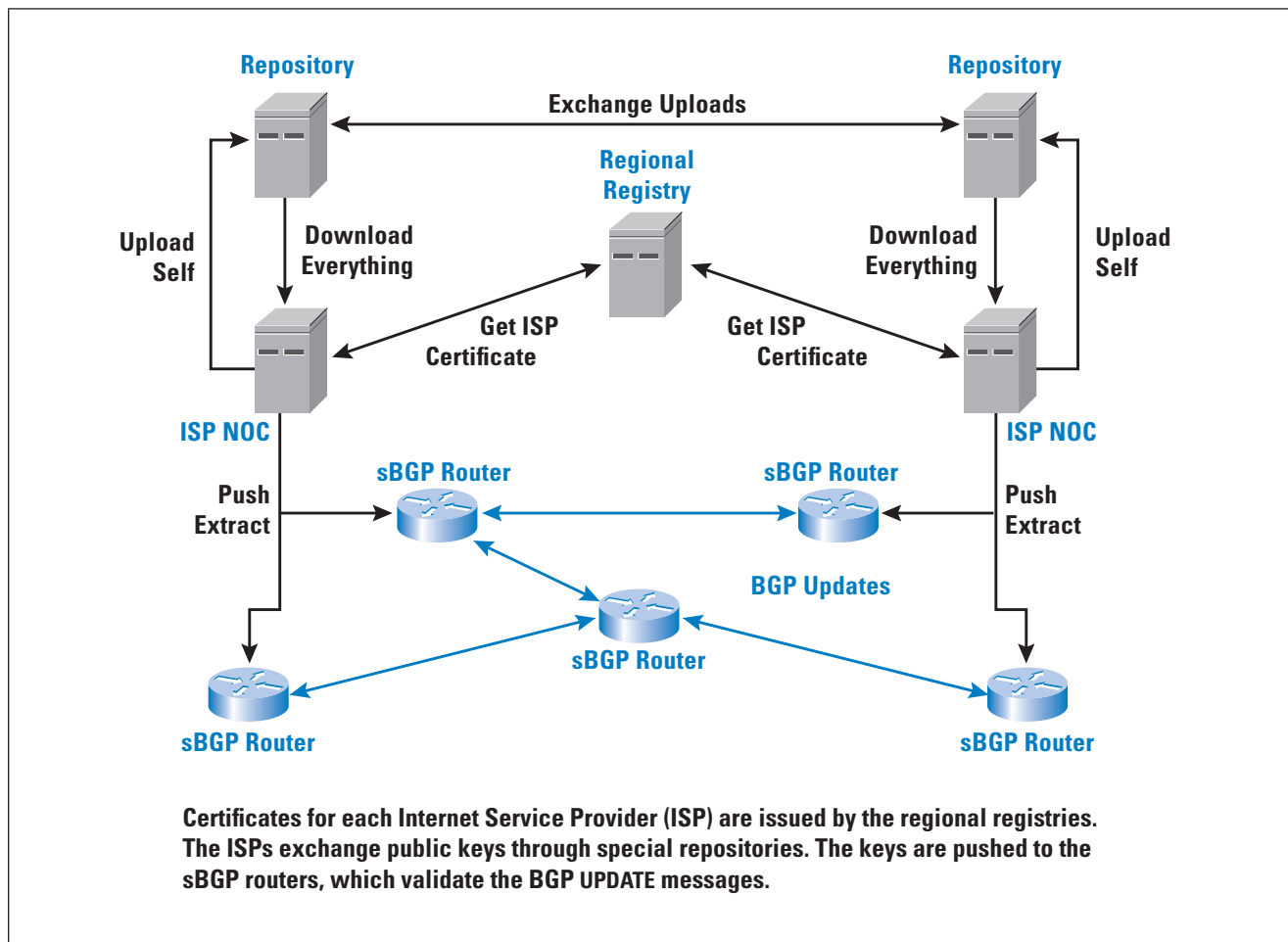
The basic security framework proposed in sBGP is that of digital signatures, X.509 certificates, and PKIs to enable BGP speakers to verify the identities and authorisation of other BGP speakers, AS administrators, and address prefix owners.

The verification framework for sBGP requires a PKI for address allocations, where every address assignment is reflected in an issued certificate^[17]. This PKI provides a means of verification of a “right-of-use” of an address. A second PKI maps the assignment of ASes, where an AS number assignment is reflected in an issued certificate, and the association between an AS number and a BGP speaking router is reflected in a subordinate certificate. In addition, sBGP proposes the use of IPsec to secure the inter-router communication paths.

sBGP also proposes the use of *attestations*. Produced by an address holder, an *address attestation* authorises a nominated AS to advertise itself as the origin AS for a particular address prefix. A *route attestation* is produced by an AS holder; it attests that a BGP speaker is an authorised member of that AS and that it has received a specified route. The address and AS PKIs, together with these attestations, allow a BGP speaker to verify the origination of a route advertisement and verify that the AS Path as specified in the BGP UPDATE is the path taken by the routing UPDATE message via the sequence of nested route attestations.

Figure 1 shows inter-operation and information exchange between sBGP elements.

Figure 1: sBGP



sBGP proposed to distribute the address attestations and the set of certificates that compose the two PKIs via conventional distribution mechanisms outside of BGP messages. For route attestations, it is necessary to pass these attestations via path attributes of the BGP UPDATE message, as an additional attribute of the UPDATE message.

Numerous significant issues have been identified with sBGP, including the computation burden for signature generation and validation, the increased load in BGP session restart, the issue of piecemeal deployment and the completeness of route attestations, and the requirement that the BGP UPDATE message has to traverse the same AS sequence as that contained in the UPDATE message^[18,19].

Secure Origin BGP

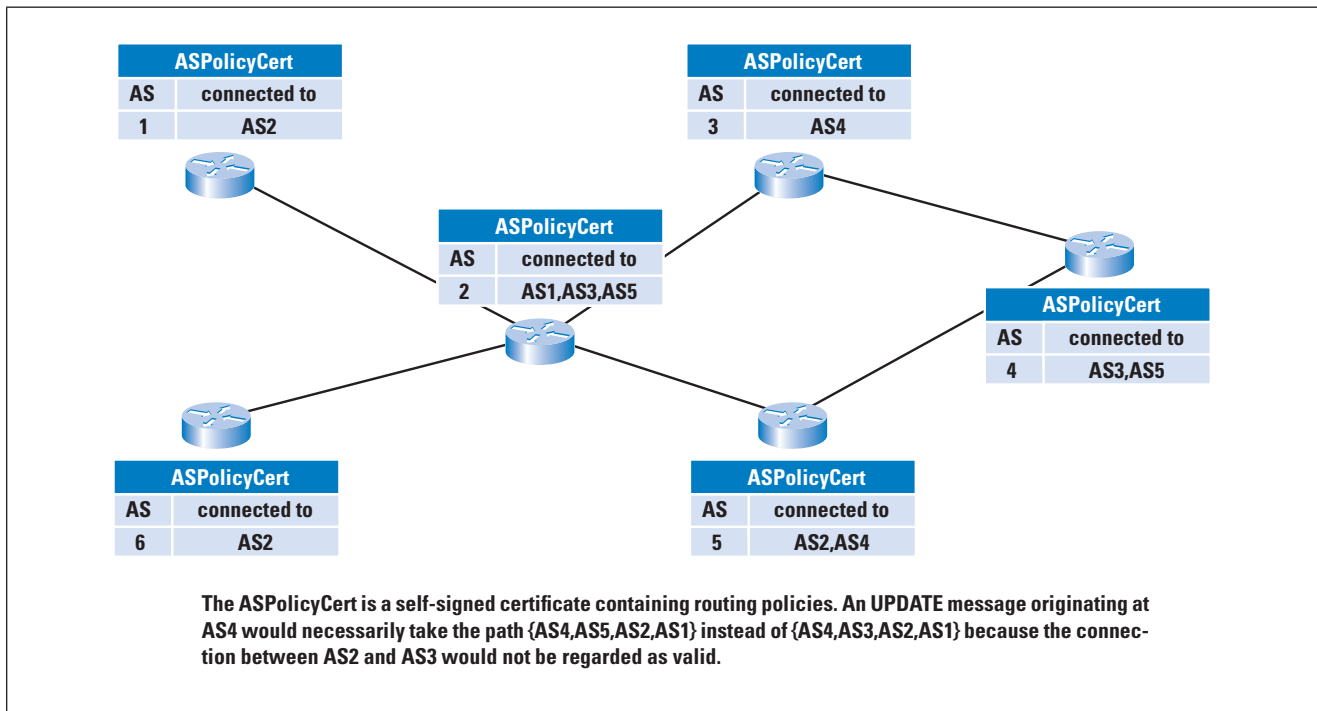
Secure Origin BGP (soBGP)^[20,21] was a response to some of the significant issues that were raised with the sBGP approach, particularly relating to the update processing load when validating the chain of router attestations and the potential overhead of signing every advertised UPDATE with a locally generated router attestation.

The validation questions that soBGP posed also included the notion of an explicit authorisation from the address holder to the originating AS to advertise the prefix into the routing system. soBGP AS Path validation is quite different from that of sBGP, in that soBGP attempted to validate that the AS Path, as presented in the UPDATE message, represents a feasible inter-AS Path from the BGP speaker to the destination AS. This feasibility test is a weaker validation condition than validating that the UPDATE message actually traversed the AS Path described in the message.

soBGP avoids the use of a hierarchical PKI that mirrors the AS number distribution framework and nominates the use of a web of trust, or a reputation mechanism, as the means of validation of these certificates. At the time of its development, no Address or AS PKI had been devised or deployed, so this web-of-trust approach was a pragmatic response to this critical omission. soBGP uses the concept of an *AuthCert* to bind an address prefix to an originating AS. This *AuthCert* is not signed by the address holder, but by a private key that is bound to an AS via an *EntityCert*. soBGP deliberately avoided the use of a PKI that was derived from the established AS and address distribution framework. This consideration appears to have been pragmatic at the time, because no such PKI existed then, and it was unclear if the various address registries were in a position to undertake this type of role of administering such a specialised PKI in any case. This situation left open the problem of how to establish trust anchors for validation of these signed objects, a rather significant deficiency in the validation framework of soBGP.

Instead of sBGP route attestations, soBGP used the concept of an *ASPolicyCert* as the foundation for constructing the data for testing the feasibility of a given AS Path. An *ASPolicyCert* contained a list of the AS local peer ASes, signed by the AS private key. An AS peering was considered valid only if both ASes list each other in their respective *ASPolicyCerts*. Figure 2 depicts a possible soBGP peering network.

Figure 2: soBGP Peering Network



The overall approach proposed in soBGP represented a different set of design trade-offs to sBGP, where the amount of validated material in a BGP UPDATE message is reduced. This approach was intended to reduce the processing overhead for validation of UPDATE messages. In soBGP each local BGP speaker assembles a validated inter-AS topology map as it collects ASPolicyCerts, and each AS Path in UPDATE messages is then checked to see if the AS sequence matches a feasible inter-AS Path in this map. soBGP proposed to use BGP itself to flood ASPolicyCerts through the network, using a new BGP message type (a *Security Message*) for this function.

The use of *Web of Trust* and the avoidance of a hierarchical PKI for the validation of AuthCerts and EntityCerts could be considered a weakness in this approach, because the derivation of authority to speak about addresses is very unclear in this model, but this absence occurred because the protocol was developed prior to the completion of the work on the *Resource PKI* (RPKI).

It is clear that soBGP could be readily adapted to use the RPKI as its trust and authority framework.

The fact that soBGP used BGP itself to flood the security credentials through the network represented an interesting approach to the problem of distributing such credentials, but it also at the time raised some unanswered questions relating to partial deployment scenarios. Interest in continuing work on soBGP waned in the early 2000s, most likely because the level of operator demand was inadequate to sustain the development effort.

Pretty Secure BGP

Pretty Secure BGP (psBGP)^[22] put forward the proposition that the proposals relating to the authentication of an address in a routing context must either rely on the use of signed attestations that need to be validated in the context of a PKI or on the authenticity of information contained in Internet Routing Registries.

The weakness of routing registries is that the commonly used access controls to the registry are insufficient to validate the accuracy or the current authenticity of the information that is represented as being contained in a route registry object. The information may have been accurate at the time the information was entered into the registry, but it may no longer be accurate at the time the relying party accesses it.

The psBGP approach was also motivated by the proponents' opinion that a PKI could not be constructed in a deterministic manner because of the indeterminate nature of some forms of address allocations. This opinion led to the assertion that any approach that relies on trusted sources of comprehensive information about prefix assignments and the identity of current right-of-use holders of address space is not a feasible proposition. Accordingly, psBGP rejected the notion of a hierarchical PKI that could be used to validate assertions about addresses and their use.

Interestingly, although psBGP rejected the notion of a hierarchical address PKI, psBGP assumed the existence of a centralised trust model for AS numbers and the existence of a hierarchical PKI that allowed public keys to be associated with AS numbers in a manner that could be validated in the context of this PKI. This notion exposed a basic inconsistency in the assumptions that lie behind psBGP, namely that a hierarchical PKI for ASes aligned to the AS distribution framework was assumed to be feasible, but a comparable PKI for addresses was not. Given that the same distribution framework has been used for both resources in the context of the Internet, it is unclear why this distinction between ASes and addresses was necessary or even appropriate.

psBGP used a rating mechanism similar to that used by PGP^[23], but in this case the rating was used for prefix origination. An AS asserted the prefixes it originated and also could list the prefixes originated by its AS peers in signed attestation.

The ability of an AS to sign an attestation about prefixes originated by a neighbour AS allowed a psBGP speaker to infer AS neighbour relationship from such assertions, allowing the local BGP speaker to construct a local model of inter-AS topology in a fashion analogous to soBGP. One of the critical differences between psBGP and soBGP was the explicit inclusion of the *strict* AS Path validation test, namely that it was a goal of psBGP to allow a BGP speaker to verify that the BGP UPDATE message traversed the same sequence of ASes as is asserted in the AS Path of the UPDATE message. The AS Path validation function relies on a sequence of nested digital signatures of each of the ASes in the AS Path for trusted validation, using a similar approach to sBGP.

However, psBGP allowed for partial path signatures to exist, mapping the validation outcome to a confidence level rather than a more basic sBGP model of accepting an AS Path only if the AS Path in the BGP UPDATE message was completely verifiable.

The essential approach of psBGP was the use of a reputation scheme in place of a hierarchical address PKI, but the value of this contribution was based on accepting the underlying premise that a hierarchical PKI for addresses was infeasible. It is also noted that the basis of accepting inter-AS ratings in order to construct a local trust value was based on accepting the validity of an AS trust rating, which, in turn, was predicated upon the integrity of the AS hierarchical PKI. psBGP appeared to be needlessly complex and bears many of the characteristics of making a particular solution fit the problem, rather than attempting to craft a solution within the bounds of the problem space.

The use of inter-AS cross certification with prefix assertion lists introduces considerable complexity in both the treatment of confidence in the assertions and the resulting assessment of the reliability of the verification of the outcome. psBGP does not consider the alternate case where the trust model relating to addresses is based on a hierarchical PKI that mirrors the address distribution framework. In such a case, the calculation of confidence levels would be largely unnecessary. The major contribution of psBGP relates to the case of partial deployment of a security solution in relation to AS Path validation, with the calculation of a confidence rating in the face of partial security information.

Inter-domain Route Validation

All of the approaches to securing the semantics of BGP described in this section so far entail changes to the operation of BGP itself and operate most effectively in an environment of universal deployment. In practical terms this scenario is unlikely, and the experience with the uptake of modifications to BGP that supported 32-bit AS number values suggests that the public Internet has considerable inertia and is very resistant to adopting changes to BGP^[24]. In a system as large as the public Internet, long-term piecemeal deployment is a far more likely scenario.

The approach proposed with *Inter-domain Route Validation (IRV)*^[25] is not to modify the BGP protocol in any way, but to define a companion information-distribution protocol.

The intent here was to attempt to provide legacy compatibility and incremental deployment capability. The IRV approach replaced the concept of simultaneously feeding both routing information and associated credentials in BGP with the concept of moving the provision of credentials into a query response framework where the receiver of a route object can query the originating AS about the authenticity of a received route object, or request additional information relating to the object in a similar fashion to the information contained in an *Internet Routing Registry (IRR)*^[26].

In IRV, each AS is responsible for providing an IRV server capable of providing authoritative responses relating to prefixes originated by this AS. IRV is envisaged as being used to provide routing policy information, using the *Routing Policy Specification Language* (RPSL)^[27,28] structure that the IRRs already use, community configuration information, contact information, a local view of the routing system in terms of received route advertisements and withdrawals, and route updates that have been sent to neighbouring ASes.

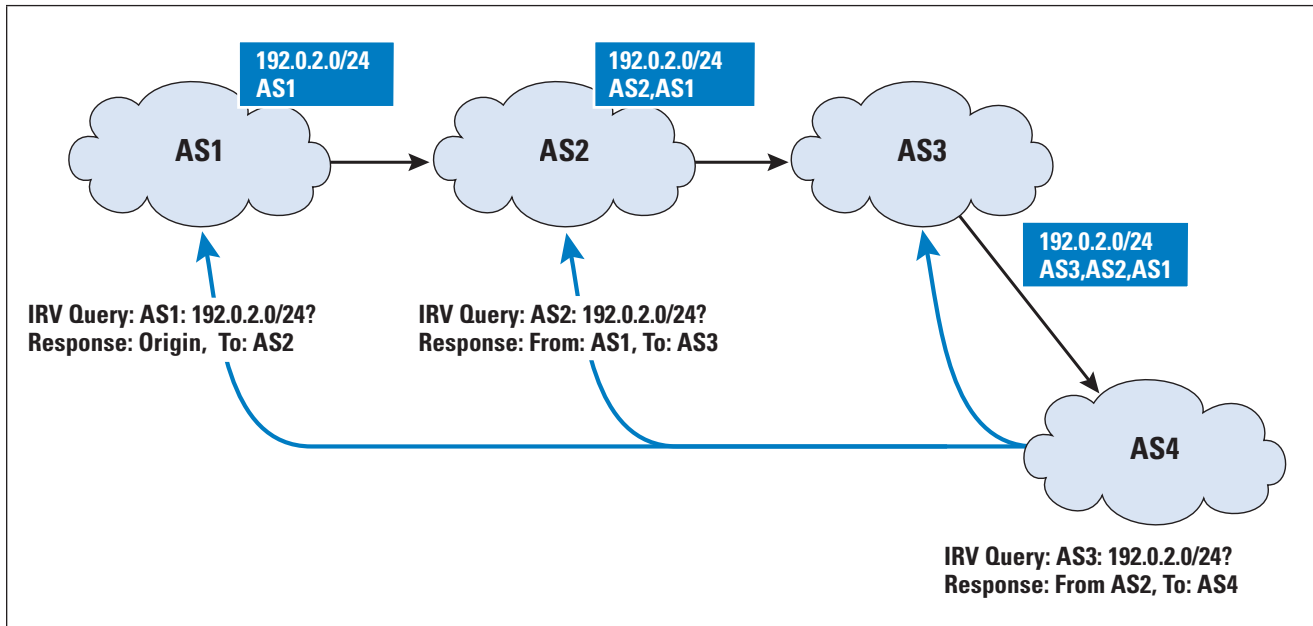
Assuming that there is a way to reliably query a per-AS IRV server and receive a response that can be validated, then AS origination validation in the IRV framework is a case of querying the originating AS IRV server with the origination query for the prefix in question and verifying the response. In a similar fashion, AS Path validation is a case of querying each IRV server of the AS in the AS Path, confirming that an advertisement was received from the previous AS in the AS Path, and that an advertisement has been sent to the next AS in the AS Path (Figure 3). This approach is midway between a strict AS Path test that validates that the UPDATE message was passed along the AS sequence described in the AS Path, and AS Path plausibility that validates that a set of AS peer connections that correspond to the AS sequence exists. Here the validation test is that each AS in the sequence is currently advertising this prefix to the next AS in sequence.

This IRV architecture has numerous issues that are not completely specified, including IRV discovery, IRV query redirection, authentication of queries and responses, selective responses, transport layer protection, and imposed overheads. It is unclear how an IRV response is to be validated, and how the relying party can verify that the received response originated from the IRV server of the AS in question, that the response has not been altered in any way, and that the response represents the actual held state in the queried AS. A similar concern lies in the estimation of additional overhead associated with performing a query to each AS in the AS Path for every received BGP UPDATE. Whether or not the query and response are preconditions to the local acceptance of a BGP route is also unspecified. While making validation of a route, a precondition for acceptance of a route would appear to offer a more robust form of security. It is also true that the IRV associated with the originating AS may be reachable only via the prefix being advertised, in which case the IRV would be unreachable until the route is accepted. It is also unclear to what extent the additional information that the IRV could provide would be useful within strict real-time constraints.

The IRV approach is essentially an extension of the IRR concept that further decentralises the publication point of routing information to individual ASes. It extends the IRR in a manner that is intended to provide adequate assurance that received responses are responses to the original query, that the response was formed by the authoritative IRV for an AS, that the response is complete and has not been altered in any way, and that the response is an accurate representation of the state of the remote AS, using DNS-style chained look-ups.

What is unclear here is whether this decentralisation has superior performance and security properties compared to an alternative approach of further augmentation to the existing IRR framework.

Figure 3: AS Path Verification Using IRV



A similar approach within the IRR framework that integrates the concept of an address and AS PKI could make provision for signed responses in a way that allows the IRR client to authenticate that the response is accurate, current, and contains information that has been digitally signed by the AS or prefix holder. In such a model of publication, the relying party is able to validate the authenticity of the IRR object independently of the manner in which the object was published or the manner in which it has been retrieved^[29].

Secure Path Vector Routing for Securing BGP

Secure Path Vector Routing for Securing BGP (SPV) is another proposal that explores the feasibility of using symmetric cryptographic operations to secure the AS Path in BGP UPDATE messages^[30] using hash chains and trees. The SPV study identified the following classes of path attacks:

- *Forgery*, where false paths are associated with routes in order to influence local route selection decisions
- *Modification*, where the path is altered in order to hide the UPDATE from a target AS or influence local route-selection decisions
- *Denial of Service*, where the attack attempts to overwhelm the intended victim's resources
- *Worm-holing*, where colluding adversaries assert false AS-to-AS links

The first two classes are attacks via BGP, whereas the second two could be more accurately classified as attacks on the routing system itself through multiparty collusion. SPV takes the approach of tree-authenticated hash values and applies it specifically to AS Path validation as an alternative to the nested digital signature structure proposed as the AS Path validation mechanism of sBGP. The SPV study paper claims significantly improved processor performance using this technique, based on the difference in computational complexity for asymmetric cryptography from symmetric cryptography as used in hash functions.

This proposal falls into the category of proposals that call for changes to the operation of the BGP protocol. In this case, the significant change is the requirement that all routes must be re-advertised to peers within a fixed time interval. This requirement is the weakest part of the approach in terms of performance evaluation, because much of the leverage in terms of scaling BGP is based on the use of a reliable transport protocol for BGP messages which, in turn, obviates any need for a BGP re-advertisement function. The need to regularly re-advertise the entire routing table to all peers has some adverse implications in terms of the performance of the protocol and its scaling capabilities.

SPV also assumes that the originating AS has knowledge of the private key associated with an address, as distinct from the more logical approach that an originating AS need only be able to produce an authority from the address allowing the AS to originate the advertisement. This approach, while efficient on processing speed, requires more storage; a higher level of time synchronisation; higher update rates within the BGP protocol, coupled with some form of loose time synchronisation; and complex key pair distribution. It has also been observed^[31] that SPV does not sufficiently protect against route forgery and eavesdropping or collusion attacks.

Signature Amortisation and Aggregate Signatures

If the signature load of sBGP is the problem, then how can this load be reduced? Numerous papers have addressed this question.

It may be possible to amortise the cost of signature validation over many messages^[32]. The technique signs a subset of the connected topology over which an UPDATE flows and places a topology description as a vector in an equivalent of an AS connectivity attestation that is flooded to all relying parties. The AS Path signing can then be generalized such that the same vector is reproduced in the signed data, with the AS neighbours who were passed the UPDATE messages marked in the bit vector. All AS neighbours can now receive the same UPDATE.

Related work^[33] combines the time-efficient approach of signature amortisation with space-efficient techniques of aggregate signatures to propose a set of constructions for aggregated path authentication that improve on the sBGP requirements for processing throughput and memory space.

Aggregate signatures apply to a collection of UPDATE messages that are to be sent to a peer. Instead of signing each UPDATE separately, the UPDATE messages are hashed into a *Merkle* hash tree^[34] and the root of the tree is signed, and the UPDATE and the root of the hash tree are sent as the signed UPDATE to each peer. This technique improves upon [35], which uses bilinear maps instead of Merkle hash trees.

Exploiting Path Stability

You also can mitigate the validation overhead by caching validation outcomes and reapplying the outcome if the same update information is received within the cache lifetime. A study by Butler, McDaniel, and Aiello^[36] noted that across a 1-month period less than 2% of advertised prefixes were advertised using more than 10 paths and less than 0.06% of prefixes were advertised with more than 20 paths.

Their paper proposed combining numerous approaches to reduce the AS Path validation workload. The first was the use of hash chains and signature aggregation, where a BGP speaker sends all local viable paths to its peers along with the tokens that represent hash chain anchors, allowing route change to be represented by an authentication token that can be validated by hash operations. The second part of the approach was to use Merkle hash trees to sign across a set of UPDATE messages that are queued awaiting the *Minimum Route Advertisement Interval* (MRAI) timer. The third part of the approach was to exploit the stability of path advertisements to amortise cryptographic operations over many validations by caching the cryptographic proofs. The paper asserted that simulations point to a reduction of the computational costs by as much as 97% over existing approaches using this approach.

Another approach, termed *pretty good BGP* (pgBGP)^[37], analyses path stability over a longer period of time and builds a local database that is then consulted in order to detect anomalous routes. The idea is that origin ASes usually do not suddenly change over time for certain prefixes, and such a sudden change might indicate an attack to the routing system. pgBGP does not provide completely automated security, because it does not eliminate any route advertisements, but rather puts them into quarantine for 24 hours (similar to route flap damping), giving operators the time to decide how to classify the event. You can deploy this proposal incrementally, and it imposes little overhead on the routing system. It is a method to mitigate effects of an attack to the routing system, not an effective mechanism for prevention of such attacks.

Detecting Prefix Hijacking

One special case of routing attacks that is considered a major threat and evokes high interest in the research community is *prefix hijacking*. A considerable amount of research has been undertaken to provide security against this single form of attack. The approaches describe possible methods of detecting prefix hijacking^[38,39,40,41], as well as complete systems and implementations of prefix hijacking detection in order to possibly react to the attack.

These systems^[42,43,44,45] rely on existing external route-monitoring databases like *Route Views*^[46] or need special routing registries to be deployed to detect prefix hijacking. The quality of such prefix hijack detection systems is strongly dependent on the quality of the route databases, all of which have some level of perspective bias given that all views of the BGP routing system are relative to the location of the collector.

Another method to detect prefix hijacking is to look for *Multiple Origin AS* (MOAS)^[47,48], which can be either a sign of multi-homing an AS or of bogus route announcements, thus prefix hijacking.

A different approach is presented for *iSPY*^[49], which tries to detect prefix hijacking by continuously probing known transit ASes in order to detect whether the prefix owned by the probing AS has been hijacked through a path change in the routing fabric to reach the address prefix.

Secure BGP and BGP Dynamics

If securing BGP is a case of applying cryptographic operations to BGP UPDATE messages, then the other approach to reducing the security overhead is to exploit the dynamic behaviour of these messages.

The BGP update pattern is addressed in [50], where a study of BGP update dynamics showed that a cache of 10,000 prefix and AS Path validation outcomes, or less than 5% of the total number of distinct routed entries, would achieve a cache rate of between 30% to 50% using a simple *Least Recently Used* (LRU) cache-replacement algorithm.

When distance-vector algorithms react to a change in prefix reachability, many UPDATE messages are generally observed before the routing system reaches a stable state. A study of BGP convergence across the global Internet concluded that the severity of path exploration and the convergence speed depend on the relative positions of the event origin and the observer^[51].

This study aligned the originator and the observer in terms of the “tiering” of *Internet Service Providers* (ISPs) and noted that these extended convergence times and larger path exploration events occurred at lower levels of the tiering hierarchy. It hypothesised that the richer inter-connectivity that was typically prevalent at such lower levels in the tiering hierarchy was a major contributing factor here. Fail-over and new route announcements converge in similar times, while route withdrawals have far longer convergence times.

A similar study on BGP path exploration characteristics proposed modifications to the BGP UPDATE message intended to identify and limit the path exploration behaviour of BGP^[52].

If a significant level of update load is related to path exploration and a significant level of AS Path security overhead is related to validation of short-term transient routing states associated with path exploration, then another direction in terms of reducing security overheads is to limit path exploration behaviour. An approach to do so by selective damping of BGP updates that are characteristic of BGP path exploration following a withdrawal at source is described in *Path Exploration Damping*^[53,54].

Further study of BGP update behaviour has explored the level of determinism that exists in the BGP route-selection process and noted that in the absence of the *Multiple Exit Discriminator* (MED) and *Route Reflectors*, the process can be considered to be a deterministic one^[55]. The paper suggests some refinements to BGP that could achieve a similar outcome to MEDs and Route Reflectors while preserving the deterministic route-selection property. The question this paper raises is that most security proposals view AS Path validation as an “after-the-event” activity because of the assumed lack of predictability in BGP. This paper questions this basic assumption and raises the possibility of path security as a provisioning activity, which, in turn raises some interesting performance optimisations for BGP path security as a provisioning exercise rather than a reactive task.

Securing the Data Plane

Securing BGP is not only a matter of securing the control plane, but also of securing the data plane^[56] and ensuring that the status of the forwarding table is consistent with the advertised BGP routing information.

A study by Mao et al.^[57] showed that up to 8% of the paths advertised through the control plane do not match the actual paths in the data plane. The data plane is subject to not only attacks that try to subvert the routing system, but also to synthetic BGP announcements from network operators that could enable the theft of carriage capacity. It is, therefore, necessary to provide security for the whole data path, not only on a next-hop basis as *Stealth Probing*^[58] intends to, because carriers might span over multiple ASes and synthesise false routing information that spans multiple AS hops.

Proposed approaches focus mainly on probing the full data path through packet injection, trying to detect and isolate malicious routers. In “Secure Traceroute”^[59], a modified *traceroute* is used to control which path data packets actually take and compares it to the actual AS Path of the routing table, effectively detecting malicious ASes. Secure Traceroute comes with the overhead of a PKI and related key exchange and no chance for piecemeal deployment.

The Fatih approach^[60] instead focuses on using traffic summary functions, and comparing their results with those of other routers, allowing detection of ASes that provide anomalous values. These traffic summary functions seem to be prone to inaccuracy because of a variety of applications running on routers that might alter the packet flow, and their application appears infeasible in routers with very high packet volumes.

The solution proposed as *Listen and Whisper*^[61] tries to detect inaccuracies in the data plane (the Listen part) but focuses also on control-plane security (the Whisper part) and aims to provide an almost complete BGP security solution, combining both parts. Compared to sBGP, Listen and Whisper should be classified as a “just-too-late” solution for BGP security, like many solutions that try to ensure data plane/control plane consistency. Like other data-plane security solutions, this approach seems infeasible, because it tries to detect data-plane anomalies by analysing individual TCP flows, and scaling this approach to the high-speed core of the Internet presents some practical challenges.

An approach that aims towards high performance and possible partial deployment is described in [62]. Its focus is to ensure that the data path always conforms to the announced AS Path, and is achieved by probing data paths by injecting tagged IP packets, or by using IP options. Similar to pgBGP, it leaves the decision of which action to take towards a malicious router to the network operator and builds up a small database to detect possible malicious routers. It deploys the roles of *verifiers* and *provers* on certain ASes, with the verifier being an AS that wants to verify a certain route, and the prover being an AS that helps the verifier in the process by replying on probe data.

Even though all these approaches intend to provide a certain level of data-plane security, and also a certain level of control-plane security, none provides comprehensive data-plane security. Authenticity of a data path from start to end could easily be forged by two ASes deploying tunnels between them, and thus disabling the possibility to effectively verify the data path by a third party.

IETF Activity – RPKI, ROV, BGPsec, and ASPA

Following numerous efforts to make progress in this area, the IETF charted a *Routing Protocol Security Requirements Working Group* (RPSEC) in 2002 to develop a common set of security requirements for routing protocols. The activity concluded in 2009. In terms of the study of inter-domain security requirements, the work stalled on some fundamental and evidently irreconcilable disagreements over the issue of the requirements for AS Path security^[63,87], and the BGP-related working drafts from the RPSEC Working Group were never published as RFCs.

Based on the initial RPSEC work on security of route origination, the *Internet Engineering Steering Group* (IESG) chartered the *Secure Inter-Domain Routing Working Group* (SIDR) in 2006^[64]. The charter for this effort presented some problems, in that it was stalled in assuming security requirements for AS Path validation and had to await results from the RPSEC activity. Given that RPSEC was unable to agree on a requirement for AS Path security, the initial work in SIDR was concentrated on securing the origination of routing information rather than its propagation through the inter-domain space.

Notably in retrospect, SIDR was also constrained from making any changes to the BGP protocol, implying that any security framework applied to the operation of BGP was to be positioned as an overlay rather than a basic change to the BGP protocol itself. This decision turned out to be very important because it precluded some design decisions that would turn out to be critical for the SIDR design work.

The initial SIDR products were a collection of specifications that described a profile for a PKI for IP addresses and AS numbers (the RPKI), as well as a model for publication and maintenance of local cache, discussed earlier in Part 1 of this survey. From this foundation, the SIDR Working Group moved on to Route Origination Validation.

Route Origination Validation

Route Origination Validation (ROV) builds upon the earlier work in the Routing Registry effort, where a prefix holder is able to publish information as to how an address prefix is to be announced into the routing system by nominating the AS number(s) that are permitted to originate a routing announcement for the prefix. In the RPKI framework this information is published as a signed *Route Origin Authorization* (ROA)^[65,66].

A ROA, which is signed by a prefix holder, denotes a permission given by the address prefix holder for an AS to originate a route.

Many additional implications are associated with publishing a ROA. The first is that no other AS has permission to announce that prefix when a cryptographically valid ROA exists in the RPKI system. If the prefix holder wishes to authorize multiple ASes to originate a route for this prefix, the prefix holder must generate multiple ROAs, meaning that an address holder can declare that a prefix should not be routed at all by issuing a ROA that provides a permission to AS0. Secondly, the ROA denies permission for any AS to originate a prefix that is more specific than the prefix listed in the ROA. You can use a *MaxLength* attribute of a ROA to define a range of more specific prefix lengths than a ROA permits. Thirdly, there is no acknowledgement of the ROA on the part of the AS. A prefix holder may publish a ROA providing a permission to an AS that is unaware of the permission.

The RPKI framework has no symmetric instrument relating to the AS holder. An AS holder does not have the ability to issue a signed attestation that lists all the prefixes that it intends to originate in the routing system.

One more important component of the ROV framework is the *RPKI to Router Protocol* (RTR)^[67]. This protocol allows you to remove a crypto engine from a router and operate on a dedicated platform. The result of this local processing of ROA data is expressed in the form of a filter list, which is implemented as a shared state between a RTR server and one or more RTR client routers. This mechanism offloads most of the RPKI overheads from the router and leaves just a residual filtering function on the router.

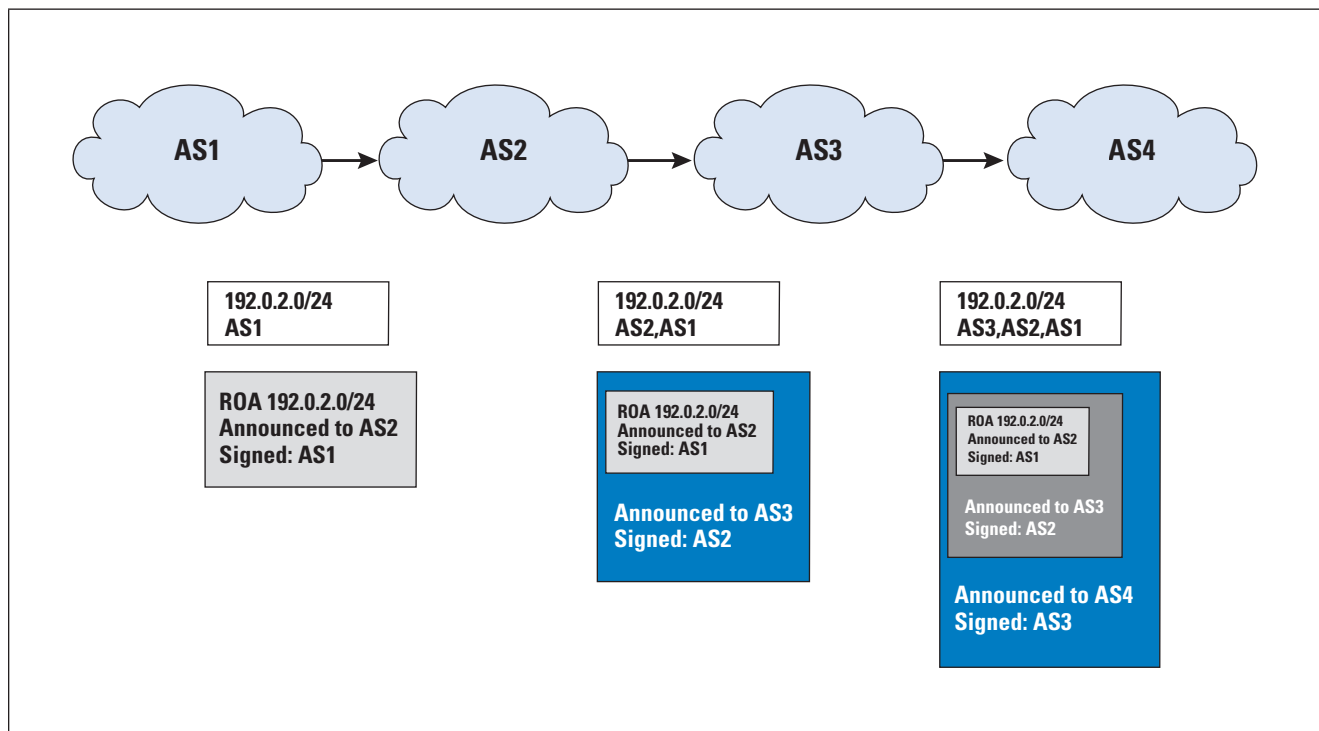
BGPsec

The SIDR working group commenced work on an extension to BGP what would allow validation of the AS Path attribute in 2011, and the standard track specification of BGPsec was published in 2017^[68].

Unlike ROV, BGPsec is not implemented in an off-router mode but is implemented through the definition of nontransitive BGP AS Path attributes. These attributes carry the digital signatures produced by the AS that propagates a BGP UPDATE message. These signatures, signed by the AS, provide confidence that every AS listed in the AS Path attribute has handled the propagation of this prefix, that the order in the AS Path is the exact order of propagation of the UPDATE message through the inter-domain routing space, and that each AS listed has explicitly authorised the propagation of an UPDATE message to its eBGP peer.

BGPsec appears to be solidly based on the concepts described in the earlier sBGP work^[8]. In essence, each eBGP speaker generates a digital signature that covers the information it received (including that digital signature) and the AS number to whom this UPDATE is to be sent (Figure 4). There is a wealth of detail behind this simple overview, but it can be summarised by the observation that this mechanism ties the AS Path in the UPDATE message to the sequence of ASes that handled the propagation of the route object. A detailed exposition of BGPsec design decisions is available in [69].

Figure 4: BGPsec Handling of AS Path Signature Structure



Stepwise AS Path validation cannot tolerate AS Sets in this approach, nor AS Confederation Sets, nor sets that are in the process of being deprecated in response to this limitation^[88]. In a similar vein, BGP Route Reflectors require special processing, as do private AS numbers.

This design approach has numerous consequences.

The first, and perhaps the most important consequence, is that piece-meal incremental deployment is simply not possible in BGPsec. When an UPDATE is passed from a BGPsec BGP speaker to a non-BGPsec BGP speaker, all BGPsec attributes are lost, meaning that if the UPDATE is further propagated to a BGPsec BGP speaker, the initial BGPsec information is unavailable. In today's Internet, the consequences of this highly constrained deployment scenario are prohibitive factors for adoption.

This approach also places a high crypto processing load on BGPsec-aware BGP speakers. There is some scepticism that this load is a feasible impost on the routing infrastructure of the Internet, and this scepticism guided the design of the ROV RTR approach. However, for BGPsec, not only are routers expected to process the BGPsec messages, but they also hold secure private keys to perform signing in real time for outgoing UPDATE messages.

Thirdly, while this approach can provide some assurance regarding the "correct" operation of the BGP protocol and can detect efforts to tamper with update messages, there is no protection against spurious WITHDRAW messages, no ability to ascertain the alignment of the route object with the forwarding state of the network, and no protection of alignment of the UPDATE with the policy state. In other words, route leaks can still occur in BGPsec.

In summary, BGPsec represents a relatively high overhead to pay for a limited set of assurances and a limited protective capability. Furthermore, a more extreme view says that BGPsec cannot achieve any of the security properties because of the fundamental design principles of BGP and BGPsec. In one research paper^[70], it is asserted that routes can still be hijacked in BGPsec, and routing loops can still appear. The authors of the paper hope to stimulate further dialog to rethink the fundamental tenets of BGP and BGPsec designs by publishing their analysis of the observed shortcomings of BGPsec.

Autonomous System Provider Authorization

The issue with the overall SIDR approach to BGP security is that if BGPsec is impractical, we cannot rely on ROV alone. All a determined routing attacker would need to do is tack on the originating AS to a synthesised AS Path and then could place any AS sequence in the AS Path attribute of a synthetic route.

ROV represents a substantial effort to get the infrastructure deployed, but without any form of AS Path protection the level of protection ROV offers is minimal at best.

The conclusion is that ROV needs to be accompanied by some form of AS Path validation if it is to be useful.

Many proposals to address this shortfall have been made. An interesting approach is *Peer Lock*^[71], which is based on the observation that the core of the routed Internet is a small set of Tier 1 ASes, and no customer of an AS should be announcing a route where the AS Path includes any of these Tier 1 networks. Secondly, no more than two of these Tier 1 ASes should appear in any AS Path, and if there are two such ASes in the AS Path they should be adjacent. This approach does not necessarily catch much in the way of deliberate efforts to generate a synthetic AS Path, but it can be effective in catching many common forms of route leaks, and its implementation is quite simple and very lightweight.

Can we do better?

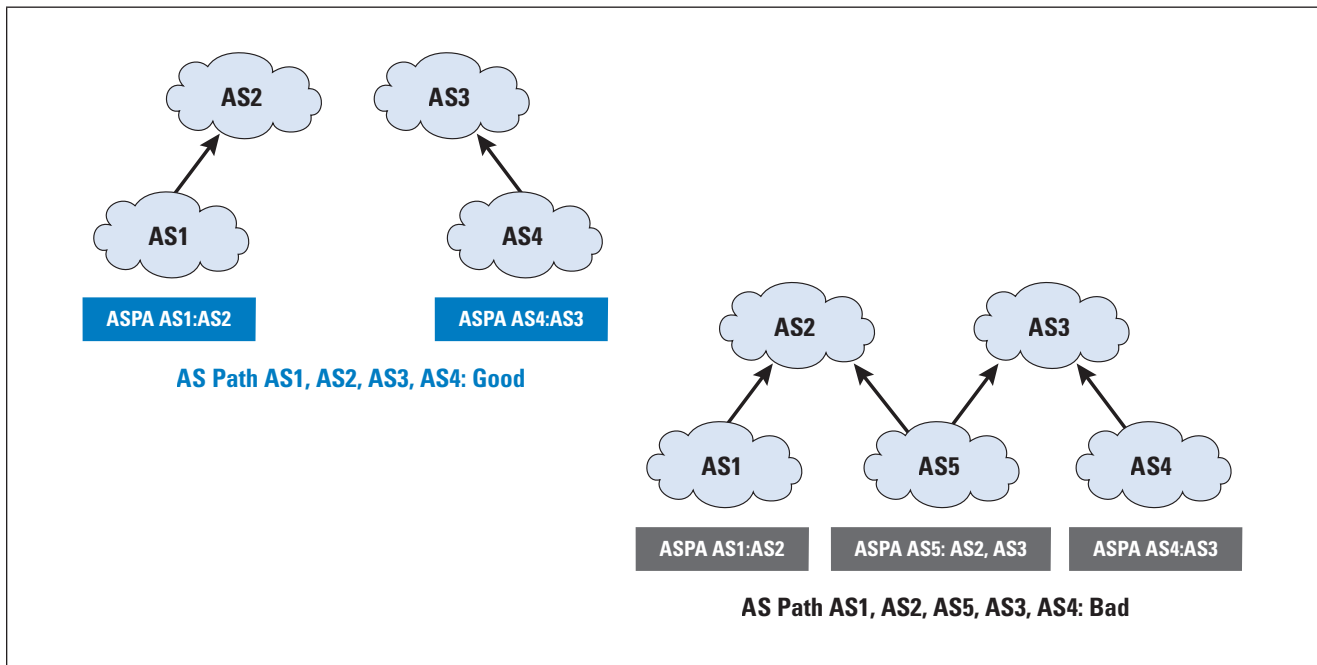
In what appears to be a replay of the situation from around 2000 when soBGP was proposed as a lighter weight response to the crypto load associated with sBGP in the area of AS Path validation, a proposal to use RPKI-signed AS adjacency attestations as a response to the issues with BGPsec has been made.

The proposal has a slight twist, however, which is different from soBGP in that an element of routing policy is also used in the *Autonomous System Provider Authorization* (ASPA) proposal^[72]. Instead of an AS listing its adjacent ASes in the inter-domain routing space and requiring both ASes to list each other as BGP neighbours before accepting the AS adjacency as valid, the ASPA framework requires an AS to list only its adjacent ASes that act in a transit provider role to the issuing AS. Given that a common criticism of BGPsec, sBGP, and soBGP was that these proposals were incapable of identifying route leaks (because route leaks represent a violation of route policy as distinct from a violation of the BGP protocol itself). ASPA provides a means of identifying such route leaks.

The ASPA relationship is a graph fragment in the directed graph that describes the inter-AS topology^[73]. The property that the ASPA proposal uses is described as “valley-free” AS Paths. All AS Paths can be characterised by zero or more paired relationships from customer-to-provider (up), zero or one peer-to-peer relationships (flat), and zero or more provider-to-customer relationships (down). In other words, all viable AS Paths are a sequence of customer-to-provider (up) AS pairs, then a peer AS pair, and then a set of provider-to-customer (down) AS pairs. Any AS sequence that contains a down and then an up (or a “valley”) represents a customer AS leaking routes learned from one provider to another (Figure 5).

ASPA requires any AS that issues an ASPA object to comply with the constraint that the providers listed in an AS ASPA are the complete set of providers for that AS.

Figure 5: ASPA and Route Leaks



ASPA still provides some benefit, even in scenarios of partial deployment. After an AS issues an ASPA, a routing attacker can include this AS in a synthetic AS Path attribute only if it also includes an adjacent provider AS, and the synthetic AS pair can be inserted in the “front part” of the AS Path (customer-to-provider) only if the order is preserved, and in the “back part” of the AS Path (provider-to-customer) in reverse order. Like soBGP, the use of ASPAs does not necessarily prevent the synthesis of AS Paths by a routing attacker, but it limits what you can use to make such synthetic paths, and the greater the use of ASPAs the more it becomes the case that the only AS Paths that can be synthesised are viable BGP AS Paths in any case. soBGP termed this constraint *AS Path Plausibility*, and the same condition applies to ASPA.

It’s evidently still early days for ASPA, and after 3 years the work remains a study item in the *SIDR Operations* (SIDROPS) Working Group of the IETF. Part of the issue here is that the SIDROPS Working Group has had its collective attention diverted away from the issues of BGP security mechanisms and AS Path validation and has taken on the role of the RPKI operational maintenance working group. In addition, in the area of RPKI operations the topic that presently takes up the working group’s attention is not the PKI itself, but the ongoing ramifications of the original design decision to use an out-of-band client-pull credential distribution mechanism for RPKI distribution. The emerging observation is that this original design choice is sufficiently flawed that the efforts in the working group to adjust the parameters of this distribution system will in all likelihood be unable to adequately address the operational issues that accompany scaling up the use of the RPKI credential system.

It may be productive at this point in time to reopen the question of how to use BGP itself to perform a *just-in-time* push-based distribution of BGP security credentials, but within the structure of the IETF it is difficult for an operationally focussed working group to perform protocol development work. However, it's an equally difficult ask for the IETF to reopen a protocol design effort on BGP security so soon after the closure of the original SIDR effort. The protracted and painful saga of the DNSSEC development effort in the IETF is one that many participants in the IETF are unwilling to repeat for BGP security.

Open Questions on Securing BGP

It appears to some observers that no current solution to routing security has found an adequate balance between appropriate security and acceptable deployment overhead^[74,75], and that's an observation that I can agree with. We are just not there yet.

Current research on BGP performance is focused on topics related to scalability, convergence times, stability, and consistency, while the questions on security research have been focused on the integrity, authenticity, authority, and verifiability of routing information. These two fields of research are inherently connected, in that a more stable routing system that can provide clear indications when convergence to a stable routing state is achieved is believed to also provide clear indications of when verification of routing information is appropriate.

In exploring the threat model for BGP, it is noted that BGP was designed to support inter-domain routing between trusted networks, while today's networks operate in a looser confederation that does not exhibit the same mutual trust properties. Not only are the TCP sessions that BGP uses vulnerable to attack, and the messages that BGP uses vulnerable to alteration that would disrupt the network routing system, the integrity of the operation of BGP is also threatened by misconfiguration, where incorrect information is injected into the routing system unintentionally, and by router vulnerabilities where a compromised routing system can exploit its trusted role and intentionally inject false information into the routing system.

Some of these attacks are intended to overwhelm a BGP speaker and force its reset, because BGP is a method of directly accessing the processing unit of a router and a saturation attack can cause processor and memory overload. Other attacks are aimed at altering the forwarding state of a router, generating an incorrect or unintended forwarding state for one or more prefixes. Other forms of attack are aimed at causing a BGP speaker to become unstable and thereby disrupt the forwarding function and impact on applications. A BGP session that is being continually reset will cause large local traffic bursts as neighbouring BGP speakers continually resend their routing tables upon each reset, and the continued instability will trigger a flap damping response in other BGP speakers.

The factors that contribute to these vulnerabilities include a lack of BGP message integrity checks, an as yet partial ability to check the authority of an originating AS to actually originate an advertisement for a prefix, and an inability to verify the accuracy, completeness, and authenticity of AS Path attributes of a routing advertisement. The use of the RPKI to support address attestations, as in ROAs, provides a very robust means of detecting incorrect origin route objects, as long as the RPKI itself is accurately aligned to the address distribution framework and as long as the RPKI is generally, if not universally, used.

In contrast, robust solutions to the problem of AS Path authentication have been elusive so far. BGPsec provides a robust method of path validation but has been assessed to be significantly expensive in terms of processor and memory cost, and also detrimental to BGP convergence times, and it requires comprehensive adoption to be effective. Efforts to substitute AS Path plausibility in place of actual AS Path validity, as is the case with ASPA, offer a different level of robustness that appears to be more practically achievable.

The study of approaches to securing BGP has raised several questions about the behaviour of inter-domain routing and the most effective approach to securing BGP. These questions include consideration of security topics and raise the issue of whether it is possible to secure the routing information to the extent that the routing information being presented is tightly aligned to the associated forwarding state^[76]:

- Is it possible to secure this association of routing information to the chained forwarding state? Can a BGP speaker validate that not only the AS path as presented in a BGP route advertisement matches the BGP propagation path taken by the prefix advertisement, and also that the current forwarding state of the network to reach the address prefix is aligned to this AS Path and this alignment can be validated? To put it simply, can a router validate that a route matches the forwarding path? This question is not one that is directly addressed within any of the current set of inter-domain routing security measures.
- A related issue concerns the overheads of securing BGP and the scaling properties of BGP. Is BGP too monolithic a protocol even before adding security capabilities? BGP simultaneously performs the functions of exchanging reachable prefixes, maintaining an inter-domain network topology, binding prefixes to paths, and implementing routing policy. Would inter-domain routing be more scalable if these functions were performed by separate protocols? Adding security and authentication within BGP, as in the sBGP model, increases the complexity of the protocol and may diminish its long-term prospects for scalability across ever larger and denser inter-domain topologies. At the same time, using a separate mechanism to flood security credentials in a manner that is entirely distinct from BGP itself, as used in the Route Origination Validation framework, becomes a source of additional operational complexity and potential vulnerability, even though the BGP protocol itself is unaltered.

Following are several practical and some more fundamental questions relating to securing BGP:

- The first is a practical question relating to the inevitable design trade-off between the level of security and the performance overheads of processing security credentials. The question concerns what aspects of securing BGP should be considered essential and what is simply desirable, but not essential. Our level of understanding as to what aspects of BGP performance and load are critical for the robust operation of network applications and what are not so critical appears to be less than comprehensive. The impact of performance trade-offs in BGP in terms of time to converge, the size of the routing space, the router memory and processing load, and scaling capability are not well understood to the extent that there is a commonly accepted answer here.
- The next question is whether verification of the correct operation of the BGP protocol is sufficient, or whether the policy intent of the routing environ is equally critical. For example, if a stub network were to leak the routes it learned from one transit network to another transit network, this route leak would, in the normal situation, be regarded as contrary to routing policies, but there is no violation of the BGP protocol itself. If we want to also include alignment to routing policies, then the question arises as to how such policies are to be expressed, who has the authority to express them, and how BGP speakers reconcile local routing policies with external routing policies when the policies differ.
- The next question is whether securing the operation of the BGP protocol (securing the control plane) is sufficient in and of itself to adequately mitigate the vulnerabilities in the overall routing system, or whether it is also necessary to include mechanisms that extend the security model to validate that the routing information represents current forwarding state in each routing element in the network (securing the data plane). One answer to this question is that securing one element of a system with multiple components does not necessarily address the underlying vulnerabilities of the entire system. The more common outcome is that such work exposes the residual vulnerabilities in other components, and that an effective security system needs to address all components of the routing system. While it may be possible for a BGP speaker to be able to validate that the originating AS did indeed originate the prefix advertisement and that the AS Path accurately represents the propagation path of this advertisement through the network, that is not the basic question in terms of the properties of the overall system.
- The more basic question here is whether a BGP speaker can verify that if it decides to forward a packet on the next hop along a path indicated by the routing system as the optimal path to a destination, is this choice indeed the optimal local choice, and does this next-hop decision pass the packet “closer” to the destination address?

- If a comprehensive security framework is proving to be elusive in terms of deployment considerations, then could a less comprehensive approach offer acceptable outcomes? Many security frameworks demonstrate a profile of diminishing returns, where the incremental cost of deploying additional security capabilities increases, while the incremental benefit in terms of risk mitigation decreases. In the case of securing BGP, could an approach of reducing the security credential generation and validation workload, through reducing the amount or timeliness of validated information, represent an acceptable trade-off? We see a practical form of this question today, where the capabilities the Route Origination Validation offer can mitigate some forms of routing incidents but are ineffectual against other forms of route manipulation that preserve the origination data. Practically, is this mitigation enough? Or do we need to also deploy some mechanism that allows detection of various forms of AS Path manipulation? A similar question relates to the comparison of the earlier soBGP and sBGP models. Is Path Plausibility sufficient? Did the mechanisms of soBGP exercise sufficient levels of constraint such that any synthesised path is close enough to a viable network path that the difference is of little consequence from a security perspective? This question is being replayed today when we consider the relative merits of the ASPA approach against the heavier weight of the BGPsec fully signed AS Path attribute.
- A final question here concerns the practicalities of deployment. The Internet is now far too large to sustain the concept of a *Flag Day* for deployment of any technology, and it is not possible to assume that a technology would be universally adopted without a protracted period of piecemeal deployment as part of a transitional interval. Indeed, as the Internet continues to grow and the diversity within the Internet increases, the anticipated transitional periods become indefinite, and piecemeal deployment becomes a continuing factor rather than a temporary transitional factor. The questions to consider include whether it is even possible to deploy high-integrity security using partial deployment scenarios, or whether the BGP protocol is too incomplete in terms of its information-distribution properties to allow robust validation of the intended forwarding state? Does securing forwarding imply carrying additional information relating to the routing and forwarding state coupling in addition to routing that would be entirely impractical in a partial deployment scenario?

Conclusions

BGP has proven surprisingly resilient in terms of its longevity of useful operational life, despite early predictions of its imminent demise in favour of IDRP^[12]. BGP-4 has routed the inter-domain Internet since late 1993, and the number of routed elements for the IPv4 Internet “default-free zone” grew from under 20,000 distinct prefixes to some 1,000,000 distinct prefixes by mid-2021, with a further 130,000 prefixes in the IPv6 network^[10].

Despite the changes in the IPv4 address infrastructure due to exhaustion of the registry free pools, the growth in the number of routing IPv4 prefixes appears to continue unabated, and together with the continued deployment of IPv6, these numbers are expected to continue to rise in the coming years.

Because of its extensibility and large installed base, BGP-4 will likely remain the only inter-domain routing protocol in the foreseeable future for the Internet (although the term “foreseeable” is prudently measured in units of years and perhaps not in decades). So far, BGP has not changed in any substantive manner, including in its security properties.

There is ample evidence from reports of use of unregistered addresses^[77] or of “routing incidents”^[78] that BGP is the subject of various forms of accidental inattention and possibly deliberate forms of abuse. Current efforts at mitigation of these forms of abuse appear in the inter-domain routing space to be less than fully adequate, and the ease with which unauthorised or bogus route objects can be injected into the inter-domain routing system remains a continuing threat issue for the security, stability, and utility of the Internet. We appear to be getting very comfortable in operating a network that experiences a continuing stream of routing incidents, both intentional and unintentional, and the longer this situation persists the more we are resigned to just accept it as the status quo for the Internet and place the onus on applications and content-distribution systems to defend themselves from routing attack. Like many unintended outcomes, it's not the outcome we would prefer to have, nor is it necessarily the optimal outcome in terms of collective cost and benefit, but it's the outcome many of us have simply accepted. All change comes at a price, and the more we resign ourselves to operating networks in the face of a poorly secured routing system the greater the effort required to make the case that the cost of a change to improve this situation will be money and effort widely spent.

References

- [0] Geoff Huston, “A Survey on Securing Inter-Domain Routing, Part 1 – BGP: Design, Threats and Security Requirements,” *The Internet Protocol Journal*, Volume 24, No. 3, October 2021.
- [1] Steven Michael Bellovin, “Security problems in the TCP/IP Protocol Suite,” *ACM SIGCOMM Computer Communication Review*, Volume 19, Issue 2, April 1, 1989.
- [2] Vijay Gill, John Heasley, and David Meyer, “The Generalized TTL Security Mechanism (GTSM),” RFC 3682, February 2004.
- [3] Carlos Pignataro, Pekka Savola, David Meyer, Vijay Gill, and John Heasley, “The Generalized TTL Security Mechanism (GTSM),” RFC 5082, October 2007.
- [4] Andy Heffernan, “Protection of BGP Sessions via the TCP MD5 Signature Option,” RFC 2385, August 1998.

- [5] Ronald L. Rivest, “The MD5 Message-Digest Algorithm,” RFC 1321, April 1992.
- [6] Marcus Leech, “Key Management Considerations for the TCP MD5 Signature Option,” RFC 3562, July 2003.
- [7] Ronald P. Bonica, Allison Mankin, and Joe Touch, “The TCP Authentication Option,” RFC 5925, June 2010.
- [8] Karen Seo and Stephen Kent, “Security Architecture for the Internet Protocol,” RFC 4301, December 2005.
- [9] Radia Perlman, “Network Layer Protocols with Byzantine Robustness,” MIT Doctoral Thesis, August 1988.
<https://dspace.mit.edu/handle/1721.1/14403>
- [10] Bradley R. Smith and Jose Joaquin Garcia-Luna-Aceves, “Securing the border gateway routing protocol,” in *Proceedings of GLOBECOM '96, 1996 IEEE Global Telecommunications Conference*, November 1996.
- [11] Bradley R. Smith and Jose Joaquin Garcia-Luna-Aceves, “Efficient Security Mechanisms for the Border Gateway Routing Protocol,” *Computer Communications*, Volume 21, No. 3, March 1998.
- [12] “Protocol for Exchange of Inter-Domain Routing Information Among Intermediate Systems to Support Forwarding of ISO 8473 PDUs,” ISO/IEC 10747, October 1994.
<https://standards.globalspec.com/std/9960/iso-iec-10747>
- [13] Tony Bates, Randy Bush, Tony Li, and Yakov Rekhter, “DNS-based NLRI origin AS verification in BGP,” Internet Draft, work in progress, July 1998.
<https://datatracker.ietf.org/doc/html/draft-bates-bgp4-nlri-orig-verif-00>
- [14] Scott Rose, Matt Larson, Dan Massey, Rob Austein, and Roy Arends, “DNS Security Introduction and Requirements,” RFC 4033, March 2005.
- [15] Lutz Donnerhacke and Wouter Wijngaards, “DNSSEC protected routing announcements for BGP,” Internet Draft, work in progress, May 2008.
<https://datatracker.ietf.org/doc/html/draft-donnerhacke-sidr-bgp-verification-dnssec/>
- [16] Stephen Kent, Charles Lynn, and Karen Seo, “Secure Border Gateway Protocol (SBGP),” *IEEE Journal on Selected Areas in Communications*, Volume 18, Issue 4, April 2000.
- [17] Karen Seo, Charles Lynn, and Stephen Kent, “Public-key infrastructure for the Secure Border Gateway Protocol (S-BGP),” *Proceedings DARPA Information Survivability Conference and Exposition II. DISCEX'01*, August 2002.

- [18] Stephen Kent, Charlie Lynn, Joanne Mikkelsen, and Karen Seo, "Secure Border Gateway Protocol (S-BGP) – Real World Performance and Deployment Issues," in *Proceedings of the 7th Annual Network and Distributed System Security Symposium*, February 2000.
- [19] Meiyuan Zhao, Sean W. Smith, and David M. Nicol, "Evaluating the Performance Impact of PKI on BGP Security," in *4th Annual PKI R&D Workshop*, NIST, April 2005.
<https://nvlpubs.nist.gov/nistpubs/Legacy/IR/nistir7224.pdf>
- [20] Russ White, "Securing BGP Through Secure Origin BGP," *The Internet Protocol Journal*, Volume 6, No. 3, September 2003.
- [21] Russ White, "Architecture and Deployment Considerations for Secure Origin BGP (soBGP)," Internet Draft, work in progress, June 2006.
<https://datatracker.ietf.org/doc/html/draft-white-sobgp-architecture-02>
- [22] Paul C van Oorschot, Tao Wan, and Evangelos Kranakis, "On interdomain routing security and pretty secure BGP (psBGP)," *ACM Transactions on Information System Security*, Volume 10, Issue 3, July 2007.
- [23] Philip R. Zimmermann, *The Official PGP User's Guide*, ISBN-13: 978-0262740173, MIT Press, 1995.
- [24] Geoff Huston, "Exploring Autonomous System Numbers," *The Internet Protocol Journal*, Volume 9, No. 1, March 2006.
- [25] G. Goodell, W. Aiello, T. Griffin, J. Ioannidis, P. McDaniel, and A. Rubin, "Working Around BGP: An Incremental Approach to Improving Security and Accuracy of Interdomain Routing," in *Proceedings of Internet Society Symposium on Network and Distributed System Security (NDSS 03)*, February 2003.
- [26] Tony Bates, Elise Gerich, Laurent Joncheray, Jean-Michel Jouanigot, Daniel Karrenberg, Marten Terpstra, and Jessica Yu, "Representation of IP Routing Policies in a Routing Registry (ripe-81++)," RFC 1786, March 1995.
- [27] David Kessens, Tony Bates, Cengiz Alaettinoglu, David Meyer, Curtis Villamizar, Marten Terpstra, Daniel Karrenberg, and Elise Gerich, "Routing Policy Specification Language (RPSL)," RFC 2622, June 1999.
- [28] Joao Damas, Andrei Robachevsky, Larry Blunk, and Florent Parent, "Routing Policy Specification Language next generation (RPSLNg)," RFC 4012, March 2005.
- [29] Robert Kisteleki and Jos Boumans, "Securing RPSL Objects with RPKI Signatures," Internet Draft, work in progress, October 2008.
<https://datatracker.ietf.org/doc/html/draft-kisteleki-sidr-rpsl-sig>

- [30] Yih-Chun Hu, Adrian Perrig, and Marvin Sirbu, “SPV: Secure Path Vector Routing for Securing BGP,” *ACM SIGCOMM Computer Communication Review*, Volume 34, Issue 4, October 2004.
- [31] Barath Raghavan, Saurabh Panjwani, and Anton Mityagin, “Analysis of the SPV Secure Routing Protocol: Weaknesses and Lessons,” *ACM SIGCOMM Computer Communication Review*, Volume 37, Issue 2, April 2007.
- [32] David M. Nicol, Sean W. Smith, and Meiyuan Zhao, “Efficient Security for BGP Route Announcements,” TR-2003-440, Dartmouth College, Computer Science, 2003.
- [33] Meiyuan Zhao, Sean W. Smith, and David M. Nicol, “Aggregated Path Authentication for Efficient BGP security,” in *CCS '05: Proceedings of the 12th ACM Conference on Computer and Communications Security*, November 2005.
- [34] Ralph C. Merkle, “Protocols for Public Key Cryptosystems,” *IEEE Symposium on Security and Privacy*, April 1980.
- [35] Dan Boneh, Craig Gentry, Ben Lynn, and Hovav Shacham, “Aggregate and verifiably encrypted signatures from bilinear maps,” in *Advances in Cryptology - EUROCRYPT 2003*, Lecture Notes in Computer Science, Volume 2656, Springer Verlag, January 2003.
- [36] Kevin Butler, Patrick McDaniel, and William Aiello, “Optimizing BGP Security by Exploiting Path Stability,” in *CCS '06: Proceedings of the 13th ACM conference on Computer and Communications Security*, October 2006.
- [37] Josh Karlin, Stephanie Forrest, and Jennifer Rexford, “Pretty Good BGP: Improving BGP by Cautiously Adopting Routes,” in *ICNP '06: Proceedings of the 2006 IEEE International Conference on Network Protocols*, IEEE Computer Society, November 2006.
- [38] Jian Qiu, Lixin Gao, Supranamaya Ranjan, and Antonio Nucci, “Detecting Bogus BGP Route Information: Going Beyond Prefix Hijacking,” in *Third International Conference on Security and Privacy in Communications Networks and the Workshops — SecureComm 2007*, September 2007.
- [39] Changxi Zheng, Lusheng Ji, Dan Pei, Jia Wang, and Paul Francis, “A Light-Weight Distributed Scheme for Detecting IP Prefix Hijacks in Real-Time,” *ACM SIGCOMM Computer Communication Review*, Volume 37, Issue 4, October 2007.
- [40] Xin Hu and Z. Morley Mao, “Accurate Real-time Identification of IP Prefix Hijacking,” in *SP '07: Proceedings of the 2007 IEEE Symposium on Security and Privacy*, IEEE Computer Society, May 2007.

- [41] Noor Hadi Hammood and Bahaa Al-Musawi, "Using BGP Features Towards Identifying Type of BGP Anomaly," in *Proceedings of 2021 International Congress of Advanced Technology and Engineering (ICOTEN)*, July 2021.
- [42] Christopher Kruegel, Darren Mutz, William Robertson, and Fredrik Valeur, "Topology-Based Detection of Anomalous BGP Messages," in *Recent Advances in Intrusion Detection*, Lecture Notes in Computer Science, Volume 2820, Springer Verlag, February 2003.
- [43] Mohit Lad, Daniel Massey, Dan Pei, Yiguo Wu, Beichuan Zhang, and Lixia Zhang, "PHAS: A Prefix Hijack Alert System," in *USENIX-SS'06: Proceedings of the 15th conference on USENIX Security Symposium*, USENIX Association, July 2006.
- [44] E-yong Kim, Klara Nahrstedt, Li Xiao, and Kunsoo Park, "Identity-Based Registry for Secure Interdomain Routing," in *ASIACCS '06: Proceedings of the 2006 ACM Symposium on Information, Computer and Communications Security*, March 2006.
- [45] Zheng Zhang, Ying Zhang, Y. Charlie Hu, and Zhuoqing Morley Mao, "Practical defenses against BGP prefix hijacking," in *CoNEXT '07: Proceedings of the 2007 ACM CoNEXT Conference*, December 2007.
- [46] University of Oregon Route Views Project:
<http://www.routeviews.org>
- [47] Xiaoliang Zhao, Dan Pei, Lan Wang, Dan Massey, Allison Mankin, S. Felix Wu, and Lixia Zhang, "An analysis of BGP multiple origin AS (MOAS) conflicts," in *IMW '01: Proceedings of the 1st ACM SIGCOMM Workshop on Internet Measurement*, November 2001.
- [48] Xiaoliang Zhao, Dan Pei, Lan Wang, Dan Massey, Allison Mankin, S. Felix Wu, and Lixia Zhang, "Detection of invalid routing announcement in the Internet," in *Proceedings. International Conference on Dependable Systems and Networks*, IEEE, December 2002.
- [49] Zheng Zhang, Ying Zhang, Y. Charlie Hu, Zhuoqing Morley Mao, and Randy Bush, "ISPY: Detecting IP Prefix Hijacking on My Own," *ACM SIGCOMM Computer Communication Review*, Volume 38, Issue 4, October 2008.
- [50] Geoff Huston, "Measures of self-similarity of BGP updates and implications for securing BGP," in *Proceedings of the 8th International Conference on Passive and Active Network Measurement (PAM 2007)*, Volume 4427, Springer Verlag, April 2007.
- [51] Ricardo Oliveira, Beichuan Zhang, Dan Pei, Rafit Izhak-Ratzin, and Lixia Zhang, "Quantifying Path Exploration in the Internet," in *IMC '06: Proceedings of the 6th ACM SIGCOMM Conference on Internet measurement*, October 2006.

- [52] Jaideep Chandrashekar, Zhenhai Duan, Zhi-Li Zhang, and Jeff Krasky, "Limiting Path Exploration in BGP," in *Proceedings of the IEEE 24th Annual Joint Conference of the IEEE Computer and Communications Societies*, March 2005.
- [53] Tony Li and Geoff Huston, "BGP Stability Improvements," Internet Draft, work in progress, June 2007.
<https://datatracker.ietf.org/doc/html/draft-li-bgp-stability-01.txt>
- [54] Geoff Huston, Mattia Rossi, and Grenville Armitage, "A Technique for Reducing BGP Update Announcements through Path Exploration Damping," in *IEEE Journal on Selected Areas in Communications*, Volume 28, Issue 8, October 2010.
- [55] Nick Feamster and Jennifer Rexford, "Network-Wide Prediction of BGP Routes," *IEEE/ACM Transactions on Networking*, Volume 15, Issue 2, April 2007.
- [56] Dan Wendlandt, Ioannis C. Avramopoulos, David G. Andersen, and Jennifer Rexford, "Don't Secure Routing Protocols, Secure Data Delivery," in *Proceedings of the 5th ACM Workshop on Hot Topics in Networks (Hotnets-V)*, November 2006.
- [57] Zhuoqing Morley Mao, Jennifer Rexford, Jia Wang, and Randy H. Katz, "Towards an accurate AS-level traceroute tool," in *SIGCOMM '03: Proceedings of the 2003 Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications*, August 2003.
- [58] Ioannis C. Avramopoulos, and Jennifer Rexford, "Stealth probing: efficient data-plane security for IP routing," in *ATEC '06: Proceedings of the Annual Conference on USENIX '06 Annual Technical Conference*, USENIX Association, January 2006.
- [59] Venkata N. Padmanabhan and Daniel R. Simon, "Secure Traceroute to Detect Faulty or Malicious Routing," *ACM SIGCOMM Computer Communication Review*, Volume 33, Issue 1, January 2003.
- [60] Alper Tugay Mizrak, Yu-Chung Cheng, Keith Marzullo, and Stefan Savage, "Fatih: Detecting and isolating malicious routers," in *DSN '05: Proceedings of the 2005 International Conference on Dependable Systems and Networks*, IEEE Computer Society, July 2005.
- [61] Lakshminarayanan Subramanian, Volker Roth, Ion Stoica, Scott Shenker, and Randy H. Katz, "Listen and Whisper: Security Mechanisms for BGP," in *NSDI'04: Proceedings of the 1st Symposium on Networked Systems Design and Implementation*, USENIX Association, March 2004.
- [62] Edmund L. Wong, Praveen Balasubramanian, Lorenzo Alvisi, Mohamed G. Gouda, and Vitaly Shmatikov, "Truth in Advertising: Lightweight Verification of Route Integrity," in *PODC '07: Proceedings of the Twenty-sixth Annual ACM Symposium on Principles of Distributed Computing*, August 2007.

- [63] IETF Routing Protocol Security Requirements Working Group:
<https://datatracker.ietf.org/wg/rpsec/about/>
- [64] IETF Secure Inter-Domain Routing Working Group:
<https://datatracker.ietf.org/wg/sidr/about/>
- [65] Matt Lepinski, Derrick Kong, and Stephen Kent, “A Profile for Route Origin Authorizations (ROAs),” RFC 6482, February 2012.
- [66] Geoff Huston and George Michaelson, “Validation of Route Origination Using the Resource Certificate Public Key Infrastructure (PKI) and Route Origin Authorizations (ROAs),” RFC 6483, February 2012.
- [67] Rob Austein and Randy Bush, “The Resource Public Key Infrastructure (RPKI) to Router Protocol, Version 1,” RFC 8210, September 2017.
- [68] Matthew Lepinski, “BGPsec Protocol Specification,” RFC 8205, September 2017.
- [69] Kotikalapudi Sriram, “BGPsec Design Choices and Summary of Supporting Discussions,” RFC 8374, April 2018.
- [70] Qi Li, Yih-Chun Hu, and Xinwen Zhang, “Even Rockets Cannot Make Pigs Fly Sustainably: Can BGP be Secured with BGPsec?” in *Proceedings of the 2014 Network and Distributed System Security (NDSS) Symposium*, February 2014.
- [71] Job Snijders, “Practical everyday BGP filtering with AS_PATH filters: Peer Locking,” NANOG 56, June 2016.
https://archive.nanog.org/sites/default/files/Snijders_Everyday_Practical_Bgp.pdf
- [72] Alexander Azimov, Eugene Bogomazov, Randy Bush, Keyur Patel, and Job Snijders, “Verification of AS_PATH Using the Resource Certificate Public Key Infrastructure and Autonomous System Provider Authorization,” February 2021, Internet Draft, work in progress,
<https://datatracker.ietf.org/doc/html/draft-ietf-sidrops-aspa-verification>
- [73] Lixin Gao, “On Inferring Autonomous Relationships in the Internet,” *IEEE/ACM Transactions on Networking*, Volume 9, Issue 6, December 2001.
- [74] Robert Lychev, Sharon Goldberg, and Michael Shapira, “BGP Security in Partial Deployment: Is the Juice Worth the Squeeze?” in *SIGCOMM ’13: Proceedings of the ACM SIGCOMM 2013 Conference on SIGCOMM*, August 2013.
- [75] Cecilia Testart, “Reviewing a Historical Internet Vulnerability: Why Isn’t BGP More Secure and What Can We Do About It?” in *TPRC 46: The 46th Research Conference on Communication, Information and Internet Policy 2018*, September 2018.

- [76] Nick Feamster, Hari Balakrishnan, and Jennifer Rexford, “Some Foundational Problems in Interdomain Routing,” in *3rd ACM SIGCOMM Workshop on Hot Topics in Networks* (HotNets), San Diego, CA, November 2004.
- [77] Geoff Huston, the CIDR Report:
<https://www.cidr-report.org/as2.0/#Bogons>
- [78] MANRS Observatory,
<https://observatory.manrs.org/#/overview>
- [79] Brian Weis, “Why IPsec and BGP don’t play well together in real networks,” Security Area Working Group presentations, IETF 66, July 2006.
<https://www.ietf.org/proceedings/66/slides/saag-2.pdf>
- [80] Eric Rescorla, “The Transport Layer Security (TLS) Protocol Version 1.3,” RFC 8446, August 2018.
- [81] Jana Iyengar and Martin Thomson, “QUIC: A UDP-Based Multiplexed and Secure Transport,” RFC 9000, May 2021.
- [82] Internet Routing Registry: <http://www.irr.net>
- [83] Dave Mitchell, Larry J. Blunk, Danny McPherson, Shane Amante, and Eric Osterweil, “Considerations for Internet Routing Registries (IRRs) and Routing Policy Configuration,” RFC 7682, December 2015.
- [84] Yakov Rekhter, “Routing in a Multi-provider Internet,” RFC 1787, April 1995.
- [85] Sandy Murphy, Curtis Villamizar, Cengiz Alaettinoglu, and David M. Meyer, “Routing Policy System Security,” RFC 2725, December 1999.
- [86] Richard Steenbergen, “Examining the validity of IRR Data,” NANOG 44, October 2006.
https://archive.nanog.org/meetings/nanog44/presentations/Tuesday/RAS_irrdata_N44.pdf
- [87] Blaine Christian and Tony Tauber, “BGP Security Requirements,” Internet Draft, work in progress, November 2008.
<https://datatracker.ietf.org/doc/html/draft-ietf-rpsec-bgpsecrec-10>
- [88] Warren Kumari, “Recommendation for Not Using AS_SET and AS_CONFED_SET in BGP,” RFC 6472, December 2011.

GEOFF HUSTON, B.Sc., M.Sc. A.M., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

FCC Launches Inquiry To Reduce Cyber Risks

On February 25, 2022 the *Federal Communications Commission* (FCC) Chairwoman Jessica Rosenworcel shared with her colleagues a proposed action to help protect America's communications networks against cyberattacks. Earlier that week, the Department of Homeland Security warned U.S. organizations at all levels that they could face cyber threats stemming from the Russia-Ukraine conflict. The proposal would begin an inquiry into the vulnerabilities of the Internet's global routing system.

If adopted by a vote of the full Commission, this action, called a *Notice of Inquiry*, would begin a proceeding by seeking public comment on vulnerabilities threatening the security and integrity of the *Border Gateway Protocol* (BGP), which is central to the Internet's global routing system. The inquiry would also examine the impact of these vulnerabilities on the transmission of data through email, e-commerce, bank transactions, interconnected *Voice-over Internet Protocol* (VoIP), and 911 calls—and how best to address these challenges.

BGP is the routing protocol used to exchange reachability information among independently managed networks on the Internet. BGP's initial design, which remains widely deployed today, does not include explicit security features to ensure trust in this exchanged information. As a result, a bad network actor may deliberately falsify BGP reachability information to redirect traffic. Russian network operators have been suspected of exploiting BGP's vulnerability to hijacking in the past. "BGP hijacks" can expose Americans' personal information, enable theft, extortion, and state-level espionage, and disrupt otherwise-secure transactions.

Working with its federal partners, the Commission has urged the communications sector to defend against cyber threats, while also taking measures to reinforce the nation's readiness and to strengthen the cybersecurity of vital communications services and infrastructure, especially in light of Russia's actions inside of Ukraine. Chairwoman Rosenworcel also recently shared with her colleagues a Notice of Proposed Rulemaking that would begin the process of strengthening the Commission's rules for notifying customers and federal law enforcement of breaches of *Customer Proprietary Network Information* (CPNI). The inquiry under consideration would build on those efforts. For more information, visit: <https://www.fcc.gov>

APNIC Announces "hybrid" APNIC 54 Conference in Singapore

The *Asia Pacific Network Information Centre* (APNIC) is pleased to announce that APNIC 54 will include a face-to-face event in Singapore in September 8–15, 2022. The conference will provide full online participation support so all attendees—online or in-person—receive the best possible conference experience. The *Asia Pacific Regional Internet Governance Forum* (APrIGF) and the *Asia Pacific School on Internet Governance* (APSIG) intend to co-locate their 2022 meetings with APNIC 54.

This will be the first time an APNIC conference has included a face-to-face component since APNIC 49 in March 2020 held in Melbourne, Australia, in conjunction with the *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT).

As previously announced, the usual “rotation” of the location of APRICOT and APNIC conferences has been suspended since the start of the pandemic, but a decision to restart it will be considered by APNOG and APNIC in the coming months. When the conference rotation restarts, the first face-to-face APRICOT will be held in Manila, Philippines. More information about APNIC 54, including the venue, dates, online participation options, partner meetings and other details can be found here: <https://conference.apnic.net/54/>

Our Privacy Policy

The *General Data Protection Regulation* (GDPR) is a regulation for data protection and privacy for all individual citizens of the *European Union* (EU) and the *European Economic Area* (EEA). Its implementation in May 2018 led many organizations worldwide to post or update privacy statements regarding how they handle information collected in the course of business. Such statements tend to be long and include carefully crafted legal language. We realize that we may need to provide similar language on our website and in the printed edition, but until such a statement has been developed here is an explanation of how we use any information you have supplied relating to your subscription:

- The mailing list for *The Internet Protocol Journal* (IPJ) is entirely “opt in.” We never have and never will use mailing lists from other organizations for any purpose.
- You may unsubscribe at any time using our online subscription system or by contacting us via e-mail. We will honor any request to remove your name and contact information from our database.
- We will use your contact information only to communicate with you about your subscription; for example, to inform you that a new issue is available, that your subscription needs to be renewed, or that your printed copy has been returned to us as undeliverable by the postal authorities.
- We will never use your contact information for any other purpose or provide the subscription list to any third party other than for the purpose of distributing IPJ by post or by electronic means.
- If you make a donation in support of the journal, your name will be listed on our website and in print unless you tell us otherwise.

Check your Subscription Details!

If you have a print subscription to this journal, you will find an expiration date printed on the back cover. For the last couple of years, we have “auto-renewed” your subscription, but now we ask you to log in to our subscription system and perform this simple task yourself. The subscription portal is here: <https://www.ipjsubscription.org/>

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Kjetil Aas	Ron Buchalski	Dmitriy Dudko	Martijn Groenleer	David Jump
Fabrizio Accatino	Paul Buchanan	Andrew Dul	Geert Jan de Groot	Anders Marius
Michael Achola	Stefan Buckmann	Joan Marc Riera	Christopher Guemez	Jørgensen
Martin Adkins	Caner Budakoglu	Duocastella	Gulf Coast Shots	Merike Kaeo
Melchior Aelmans	Darrell Budic	Pedro Duque	Sheryll de Guzman	Andrew Kaiser
Christopher Affleck	BugWorks	Holger Durer	Rex Hale	Christos Karayiannis
Scott Aitken	Scott Burleigh	Mark Eanes	Jason Hall	Daniel Karrenberg
Jacobus Akkerhuis	Chad Burnham	Andrew Edwards	Darow Han	David Kekar
Antonio Cuñat Alario	Jon Harald Bøvre	Peter Robert Egli	Handy Networks LLC	Stuart Kendrick
Nicola Altan	Olivier Cahagne	George Ehlers	James Hamilton	Robert Kent
Shane Amante	Antoine Camerlo	Peter Eisses	Stephen Hanna	Jithin Kesavan
Marcelo do Amaral	Tracy Camp	Torbjörn Eklöv	Martin Hannigan	Jubal Kessler
Matteo D'Ambrosio	Ignacio Soto Campos	Y Ertur	John Hardin	Shan Ali Khan
Selva Anandavel	Fabio Caneparo	ERNW GmbH	David Harper	Nabeel Khatri
Jens Andersson	Roberto Canonico	ESdatCo	Edward Hauser	Dae Young Kim
Danish Ansari	David Cardwell	Steve Esquivel	David Hauweele	William W. H. Kimandu
Finn Arildsen	Richard Carrara	Jay Etchings	Marilyn Hay	John King
Tim Armstrong	John Cavanaugh	Mikhail Evstiounin	Headcrafts SRLS	Russell Kirk
Richard Artes	Lj Cemerar	Bill Fenner	Hidde van der Heide	Gary Klesk
Michael Aschwenden	Dave Chapman	Paul Ferguson	Johan Helsingius	Anthony Klopp
David Atkins	Stefanos Charchalakakis	Ricardo Ferreira	Robert Hinden	Henry Kluge
Jac Backus	Greg Chisholm	Kent Fichtner	Asbjørn Højmark	Michael Kluk
Jaime Badua	David Chosrova	Armin Fisslthaler	Damien Holloway	Andrew Koch
Bent Bagger	Marcin Cieslak	Michael Fiumano	Alain Van Hoof	Ia Kochiashvili
Eric Baker	Lauris Cikovskis	The Flirble Organisation	Edward Hotard	Carsten Koempe
Santosh Balagopalan	Guido Coenders	Gary Ford	Bill Huber	Richard Koene
William Baltas	Brad Clark	Jean-Pierre Forcioli	Hagen Hultzs	Alexader Kogan
David Bandinelli	Narelle Clark	Susan Forney	Kauto Huopio	Matthijs Koot
Benjamin Barkin-Wilkins	Horst Clausen	Christopher Forsyth	Kevin Iddles	Antonin Kral
Feras Batainah	Joseph Connolly	Andrew Fox	Mika Ilvesmaki	Robert Krejčí
Michael Bazarewsky	Steve Corbató	Craig Fox	Karsten Iwen	Mathias Körber
David Belson	Brian Courtney	Fausto Franceschini	David Jaffe	John Kristoff
Hidde Beumer	Beth and Steve Crocker	Valerie Fronczak	Ashford Jaggernaut	Terje Krogdahl
Pier Paolo Biagi	Dave Crocker	Tomislav Futivic	Thomas Jalkanen	Bobby Krupczak
Tyson Blanchard	Kevin Croes	Laurence Gagliani	Martijn Jansen	Murray Kuchera
John Bigrow	John Curran	Edward Gallagher	Jozef Janitor	Warren Kumari
Orvar Ari Bjarnason	André Danthine	Andrew Gallo	John Jarvis	George Kuo
Axel Boeger	Morgan Davis	Chris Gamboni	Dennis Jennings	Dirk Kurfuerst
Keith Bogart	Jeff Day	Xosé Bravo Garcia	Edward Jennings	Darrell Lack
Mirko Bonadei	Julien Dhallenne	Osvaldo Gazzaniga	Aart Jochem	Andrew Lamb
Roberto Bonalumi	Freek Dijkstra	Kevin Gee	Nils Johansson	Richard Lamb
Lolke Boonstra	Geert Van Dijk	Greg Giessow	Brian Johnson	Yan Landriault
Julie Bottorff	David Dillow	John Gilbert	Curtis Johnson	Edwin Lang
Photography	Richard Dodsworth	Serge Van Ginderachter	Richard Johnson	Sig Lange
Gerry Boudreaux	Ernesto Doelling	Greg Goddard	Jim Johnston	Markus Langenmair
L de Braal	Michael Dolan	Tiago Goncalves	Jonatan Jonasson	Fred Langham
Kevin Breit	Eugene Doroniuk	Ron Goodheart	Daniel Jones	Tracy LaQuey Parker
Thomas Bridge	Karlheinz Dölger	Octavio Alfageme	Gary Jones	Alex Latzko
Ilia Bromberg	Michael Dragone	Gorostiaga	Jerry Jones	Jose Antonio Lazaro
Václav Brožík	Joshua Dreier	Barry Greene	Michael Jones	Lazaro
Christophe Brun	Lutz Drink	Jeffrey Greene	Amar Joshi	Rick van Leeuwen
Gareth Bryan	Aaron Dudek	Richard Gregor	Javier Juan	Simon Leinen

Robert Lewis	Mohammad Moghaddas	Andrew Potter	SeenThere	Peter Tomsu Fine Art
Christian Libérale	Roberto Montoya	Eduard Llull Pou	Scott Seifel	Photography
Martin Lillepuu	Charles Monson	Tim Pozar	Yury Shefer	Joseph Toste
Roger Lindholm	Andrea Montefusco	David Raistrick	Yaron Sheffer	Rey Tucker
Link Light Networks	Fernando Montenegro	Priyan R Rajeevan	Doron Shikmoni	Sandro Tumini
Chris and Janet Lonvick	Joel Moore	Balaji Rajendran	Tj Shumway	Angelo Turetta
Sergio Loreti	John More	Paul Rathbone	Jeffrey Sicuranza	Michael Turzanski
Eric Louie	Maurizio Moroni	William Rawlings	Thorsten Sideboard	Phil Tweedie
Adam Loveless	Brian Mort	Mujtiba Raza Rizvi	Greipur Sigurdsson	Steve Ulrich
Josh Lowe	Soenke Mumm	Bill Reid	Fillipe Cajaiba da Silva	Unitek Engineering AG
Guillermo a Loyola	Tariq Mustafa	Petr Rejhon	Andrew Simmons	John Urbanek
Hannes Lubich	Stuart Nadin	Robert Remenyi	Pradeep Singh	Martin Urwaleck
Dan Lynch	Michel Nakhla	Rodrigo Ribeiro	Henry Sinnreich	Betsy Vanderpool
David MacDuffie	Mazdak Rajabi Nasab	Glenn Ricart	Geoff Sisson	Surendran Vangadasalam
Sanya Madan	Krishna Natarajan	Justin Richards	John Sisson	Ramnath Vasudha
Miroslav Madić	Naveen Nathan	Rafael Riera	Helge Skrivervik	Philip Venables
Alexis Madriz	Darryl Newman	Mark Risinger	Terry Slattery	Buddy Venne
Carl Malamud	Thomas Nikolajsen	Fernando Robayo	Darren Sleeth	Alejandro Vennera
Jonathan Maldonado	Paul Nikolich	Gregory Robinson	Richard Smit	Luca Ventura
Michael Malik	Travis Northrup	Ron Rockrohr	Bob Smith	Scott Vermillion
Tarmo Mammers	Marijana Novakovic	Carlos Rodrigues	Courtney Smith	Tom Vest
Yogesh Mangar	David Oates	Magnus Romedahl	Eric Smith	Peter Villemoes
Bill Manning	Ovidiu Obersterescu	Lex Van Roon	Mark Smith	Vista Global Coaching
Harold March	Tim O'Brien	Marshall Rose	Tim Sneddon	& Consulting
Vincent Marchand	Mike O'Connor	Alessandra Rosi	Craig Snell	Dario Vitali
Normando Marcolongo	Mike O'Dell	David Ross	Job Snijders	Jeffrey Wagner
Gabriel Marroquin	John O'Neill	William Ross	Ronald Solano	Don Wahl
David Martin	Jim Oplotnik	Boudhayan	Asit Som	Michael L Wahrman
Jim Martin	Packet Consulting	Roychowdhury	Ignacio Soto Campos	Laurence Walker
Ruben Tripiana Martin	Limited	Carlos Rubio	Evandro Sousa	Randy Watts
Timothy Martin	Carlos Astor Araujo	Rainer Rudigier	Peter Spekrijse	Andrew Webster
Carles Mateu	Palmeira	Timo Ruiters	Thayumanavan Sridhar	Tim Weil
Juan Jose Marin Martinez	Alexis Panagopoulos	RustedMusic	Paul Stancik	Jd Wegner
Ioan Maxim	Gaurav Panwar	Babak Saberi	Ralf Stempfner	Westmoreland
David Mazel	Manuel Uruena Pascual	George Sadowsky	Matthew Stenberg	Engineering Inc.
Miles McCredie	Ricardo Patara	Scott Sandefur	Martin Štěpánek	Rick Wesson
Brian McCullough	Dipesh Patel	Sachin Sapkal	Adrian Stevens	Peter Whimp
Joe McEachern	Alex Parkinson	Arturas Satkovskis	Clinton Stevens	Russ White
Alexander McKenzie	Craig Partridge	PS Saunders	John Streck	Jurrien Wijnhuizen
Jay McMaster	Dan Paynter	Richard Savoy	Martin Streule	Derick Winkworth
Mark Mc Nicholas	Leif Eric Pedersen	John Sayer	David Strom	Pindar Wong
Olaf Mehlberg	Rui Sao Pedro	Phil Scarr	Colin Strutt	Makarand Yerawadekar
Carsten Melberg	Juan Pena	Gianpaolo Scassellati	Viktor Sudakov	Phillip Yialeloglou
Kevin Menezes	Chris Perkins	Elizabeth Scheid	Edward-W. Suor	Janko Zavernik
Bart Jan Menkveld	Michael Petry	Jeroen Van Ingen	Vincent Surillo	Bernd Zeimetz
Sean Mentzer	Alexander Peuchert	Schenau	Terence Charles Sweetser	Muhammad Ziad
William Mills	David Phelan	Carsten Scherb	T2Group	Ziauddin
David Millsom	Derrell Piper	Ernest Schirmer	Roman Tarasov	Tom Zingale
Desiree Miloshevic	Rob Pirnie	Philip Schneck	David Theese	Jose Zumalave
Joost van der Minnen	Marc Vives Piza	Peter Schoo	Douglas Thompson	Romeo Zwart
Thomas Mino	Jorge Ivan Pincay	Dan Schrenk	Kerry Thompson	廖明沂.
Rob Minshall	Ponce	Richard Schultz	Lorin J Thompson	
Wijnand Modderman-	Victoria Poncini	Timothy Schwab	Fabrizio Tivano	
Lenstra	Blahoslav Popela	Roger Schwartz		



Follow us on Twitter and Facebook

@protocoljournal



<https://www.facebook.com/newipj>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Internet
Society



Diamond Sponsors

Your logo here!

Ruby Sponsors



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

John Crain, Senior Vice President and Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Principal – Internet Technology, Policy, and Advocacy
The Internet Society

Dr. Jun Murai, Founder, WIDE Project
Distinguished Professor, Keio University
Co-Director, Keio University Cyber Civilization Research Center, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

August 2022

Volume 25, Number 2

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
Parallel BGP Processing	2
Transport Versus Network ...	15
20 Years of SIP	25
Fragments	31
Thank You!	32
Call for Papers	34
Supporters and Sponsors	35

According to Wikipedia, the *Border Gateway Protocol* (BGP) “...is a standardized exterior gateway protocol designed to exchange routing and reachability information among *Autonomous Systems* (ASs) on the Internet. BGP is classified as a path-vector routing protocol, and it makes routing decisions based on paths, network policies, or rule-sets configured by a network administrator.” We’ve covered numerous aspects of BGP in this journal, most recently in our two-part article by Geoff Huston entitled “A Survey on Securing Inter-Domain Routing.” In this issue, a team of engineers from Juniper Networks describes a method for running BGP processing in parallel using a concept known as *sharding*.

In our second article, Geoff Huston takes a closer look at the transport and network functions in today’s ever-changing Internet. Many network elements such as firewalls and *Network Address Translators* (NATs) use the transport protocol *header* to make decisions on how to handle traffic, but concerns about pervasive monitoring and information leakage have led to various forms of encryption-based solutions and an ongoing debate within the Internet technical community and beyond.

Using the Internet for teleconferencing or telephony is not a particularly new idea. I fondly remember taking part in experiments between the *Norwegian Defence Research Establishment* (NDRE), MIT’s Lincoln Laboratories, *University of Southern California’s Information Sciences Institute* (USC-ISI), and *University College London* (UCL) as early as 1977 when I was doing my military service at NDRE. You can find out more about these early developments by searching for the article “Linear Predictive Coding and the Internet Protocol.” *Voice over IP* (VoIP) as we know it today became a reality in 2002 with the publication of RFC 3261, which describes the *Session Initiation Protocol* (SIP). In our final article, Jonathan Rosenberg gives a retrospective on 20 years of SIP.

Publication of *The Internet Protocol Journal* is made possible by the generous support of numerous individuals and organizations. Please consider making a donation or getting your company to sign up for a sponsorship. As always, we welcome your feedback and suggestions on anything you read in this journal. Letters to the Editor may be edited for clarity and length and can be sent to ipj@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

Parallel BGP Protocol Processing

by Sanjay Khanna, Jaihari Loganathan, and Ashutosh Grewal, Juniper Networks

Managing large inter- and intra-*Border Gateway Protocol* (BGP) domains places a large computational load on the CPU of a router, adversely affecting its performance and increasing the BGP convergence time. To address these problems, we have architected a solution that splits a *BGP Routing Information Base* (RIB) across concurrently running BGP threads. These parallel running threads run the same code on multiple CPU cores concurrently. Each of these threads maintains a RIB *shard*, a subset of the RIB. This parallel BGP processing improves the read-side performance of processing incoming UPDATE messages. A set of parallel running I/O threads generate outbound UPDATE messages and improve the write-side performance of BGP. This entire design uses a lockless mechanism to allow parallel processing on each CPU core independently. A testbed representing a Tier 1 service provider *Route Reflector* network was used to verify and quantify the performance of the implementation. BGP in this topology receives several copies of the global Internet routing table (~800,000 routes). Our results show that performance improves as parallelism and scale are increased. The speedup we can attain gets better as more CPU cores are available for RIB sharding. These gains are bounded by the extent to which the BGP update processing can run in parallel.

Terms and Definitions

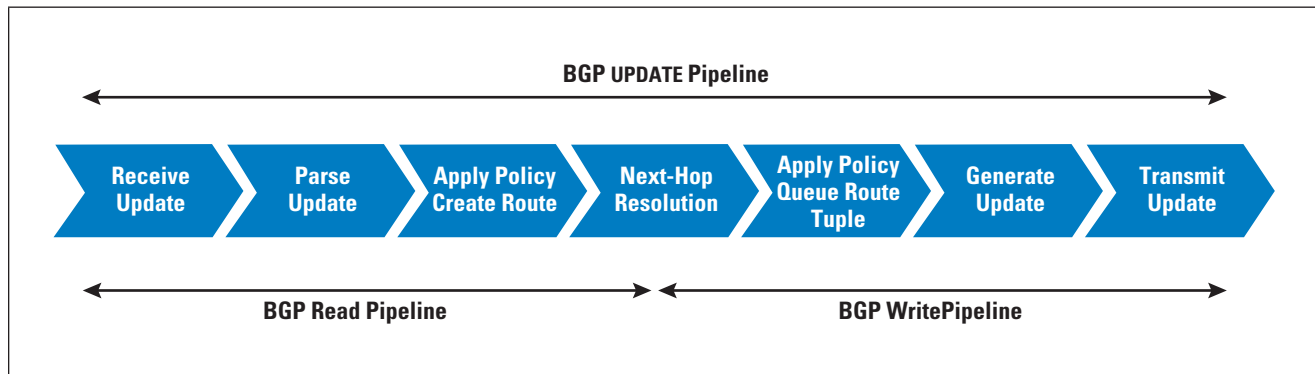
The BGP RIB conceptually consists of four parts:

- *Adj-RIBs-In*: Stores unprocessed routing information that has been learned from BGP updates received from peers. The routes contained in Adj-RIBs-In are considered feasible routes.
- *Loc-RIB*: Contains the routes that the BGP speaker has selected by applying the decision process (route selection, import policy) to the routes contained in Adj-RIBs-In. These routes populate the routing table (RIB) along with routes that other routing protocols discover.
- *Adj-RIBs-Out*: Contains the routes that the BGP speaker advertises to its peers in BGP UPDATE messages. Export routing policies determine what routes are placed in Adj-RIBs-Out.
- *Outbound Route Tuple* is a route in Adj-RIBs-Out. Every tuple consists of a prefix, associated BGP metrics, and information about peers in a peer group that will be sent to this tuple in an UPDATE message.

BGP Update Processing Pipeline

Figure 1 illustrates the BGP^[1] UPDATE message processing pipeline. This pipeline can be further subdivided into the *Read-Side Pipeline* to identify the inbound processing of an UPDATE message, and the *Write-Side Pipeline*, which concerns generation of outbound UPDATE messages to be sent to peers.

Figure 1: BGP Pipeline



The read-side pipeline consists of the following stages:

- *Receive Update*: Routes are advertised between BGP speakers in an UPDATE message. Multiple routes that have the same PATH attributes are advertised in a single UPDATE message. This message is received over an established *Transmission Control Protocol* (TCP) socket.
- *Parse Update*: A BGP UPDATE message is parsed for prefixes, and PATH attributes are canonicalized into a local state after validating the data in *Protocol Data Units* (PDUs).
- *Apply Import Policy and Create Route*: BGP stores routing information learned from the inbound UPDATE messages in a RIB called *Adj-RIB-In*. BGP applies the import policy on incoming BGP routes and—if permitted by policy—performs a best-route selection. The best routes are used to populate the local RIB, called *Loc-RIB*. These routes are then used to program a *Forwarding Information Base* (FIB) and generate outbound updates to BGP peers.
- *Next-hop Resolution*: BGP uses a local routing table to find the reachability information for a BGP next-hop, which may be several hops away. For example, *Interior Gateway Protocol* (IGP) metric, intermediate address, and outgoing interface are parts of resolving reachability for the BGP next-hop. BGP may have several routes to the same IP destination that have different degrees of preference. Although all accepted routes are installed in Loc-RIB, BGP may choose one route or multiple routes (BGP multipath case) as active routes.

The write-side pipeline consists of the following stages:

- *Apply Export Policy and Queue Outbound Route Tuple*: BGP applies export policy to routes in Loc-RIB and generates outbound RIBs called *Adj-RIBs-Out*. Adj-RIB-Out stores information that BGP uses to generate an outbound route tuple. Generally, several peers with the same policies are grouped into a *peer group*, and a single outbound route tuple is generated for each peer group.

- *Generate Update*: Route information in the outbound route tuple is converted into an UPDATE message. Destination prefixes that share the same PATH attributes are packed in a single message. This prefix packing reduces the number of UPDATE messages sent over TCP. Sending fewer messages improves local performance and does not impose extra work on the peer BGP routers.
- *Transmit Update*: UPDATE messages that were generated in a prior stage are sent over a TCP socket to a remote BGP peer. The same BGP UPDATE message is sent to multiple BGP peering sessions that share a common export policy, thereby amortizing the cost of UPDATE message generation.

BGP *convergence time* is the time taken by the router to process the incoming BGP UPDATE messages, passing them through the read pipeline in Figure 1, and distributing the results by generating UPDATE messages to its peers. As networks scale up in the number of peers, the number of parallel inbound feeds, and the size of the network in terms of the number of prefixes, this convergence time can become very high. The result can be slower convergence of the entire network when device and link failures occur. The slower convergence generally leads to traffic loss and a traffic black-hole in the network. Most operators of large networks want a faster convergence to reduce these downtimes. To help improve the convergence time of a router, and to exploit multi-core CPUs available on routing engines, we decided to run the previously mentioned pipeline functions in parallel. We encountered several hurdles to running BGP update processing in parallel:

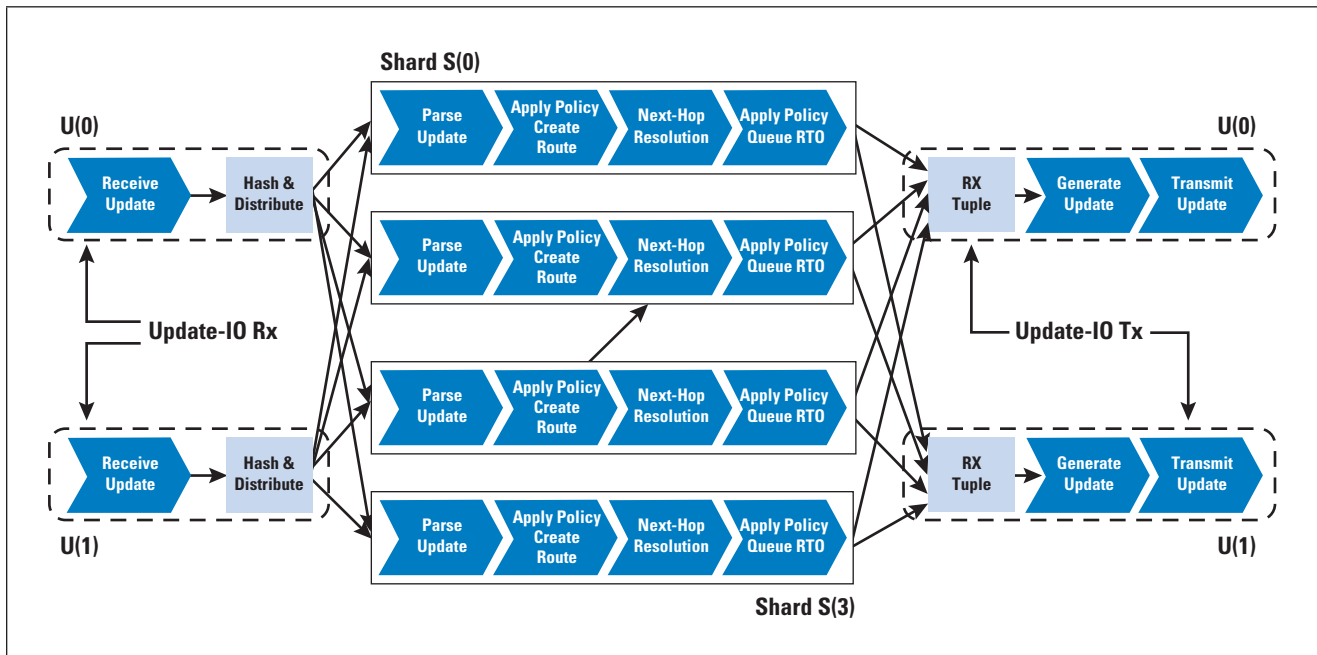
- RIBs are generally a shared collection of routes, and concurrent updates to RIBs need locks.
- Locking in general reduces concurrency and increases complexity.
- BGP next-hop resolution is a shared function and requires synchronization between parallel running threads.
- Prefix packing in outgoing UPDATE messages deteriorates because each concurrent pipeline produces more frequent and smaller UPDATE messages with fewer prefixes.
- Other protocols in an IP router need active BGP routes to implement other useful functions like *Layer 3 Virtual Private Network* (L3VPN) and programming the FIB.
- Interface state (like links and IP addresses) is a shared collection too, and concurrent threads meant more locking across concurrently running threads.

The BGP software architecture described in the following sections addresses these concurrency limitations and improves the BGP convergence time at very high scales. This solution performs better when input scale increases and the required number of concurrent pipelines increases.

Parallel BGP Protocol Processing

Figure 2 illustrates the high-level architecture of parallel BGP processing. The pipeline is broken into three logical parts running concurrently on two kinds of threads of execution: “Shard(S)” and “Update-IO(U)” threads. Shard threads process a subset of UPDATE message prefixes, execute most of the BGP read pipeline, and generate a tuple for an Update-IO thread. An Update-IO thread processes a tuple for a set of BGP peers in a peer group. Update-IO threads also receive the UPDATE message from peers (over TCP sockets) and distribute the messages to shard threads. These threads process tuples from shard threads and generate UPDATE messages to transmit to peers.

Figure 2: Parallel Processing of BGP UPDATE Messages



- *Update-IO Rx Processing:* Running in an update-IO thread, this stage receives update messages from peers, sanity checks the messages, and computes a hash on every prefix in the message to determine which shards will receive the message. In the best possible scenario a single shard gets the entire message, and in the worst-case scenario the same message is shared with every shard thread.
- *Shard Processing:* A shard thread receives a copy of the message and processes prefixes that it owns (ignoring the ones that it doesn't). The ownership of a prefix is decided by computing the hash on the prefix. Thereafter each shard will follow the entire read pipeline on its prefixes and generate tuples for the Update-IO thread matching the peer group.
- *Update-IO Tx Processing:* An Update-IO thread receives tuples and processes them into BGP UPDATE messages. These messages are sent towards the peer via TCP sockets. This stage is responsible for packing tuples from several shard threads into UPDATE messages to improve the packing of prefixes.

This architecture exploits the *Single Program Multiple Data* (SPMD) model of parallelism to do concurrent read and write pipeline processing. Concurrently executing threads do not share any state with each other, thereby eliminating any need for locks. Shards and Update-IO threads maintain a parallel ecosystem of collections of objects (like BGP peers, peer groups, RIBs, interfaces, configuration, etc.). Wherever needed, message passing is used to achieve eventual consistency of the routing state in the system. At any point of time, concurrently executing functions running at different rates may be in different states. However, all these functions executing on different cores will eventually reach the same final state as if there were a single executing pipeline. Reaching eventual consistency quickly is an important outcome of this architecture.

Sharding Several Kinds of BGP RIBs

This architecture requires all routes (BGP and otherwise) with the same IP address to always be assigned to the same shard so that the best active route calculation for all routes matching an IP address is handled in a single shard. This requirement guarantees the correctness of the active route selection algorithm. *Multiprotocol BGP* (MP-BGP) extensions^[2] allow BGP to carry routing information for multiple network layers and address families. BGP routes for each of these address families are saved in several RIBs, and each RIB has its own network prefix in a route. For example, the IPv4 Unicast RIB has the IPv4 address and prefix mask length as the route destination. The L3VPN RIB has the IP address, route distinguisher, and mask length as the route destination. For shard assignment, the hash is computed only on the address and prefix length part of the route destination, and the rest is ignored.

Next-Hop Resolver

The main job of the resolver is to translate protocol next-hops into forwarding next-hops using *helper routes*. A protocol next-hop is an IP address of a remote BGP peer, most commonly an *Interior BGP* (IBGP) peer. A protocol next-hop, by itself, is insufficient to make a forwarding decision. To forward a packet, a router needs to know the directly connected next-hop. This information is derived from helper routes that provide reachability information for the protocol next-hop. Since the resolver is a central function and the concurrent shard threads also need the services of this resolution, a mechanism is needed to distribute resolver information to shard threads. One way to achieve this distribution is to run the resolver as a service and all shard threads to register next-hop IP addresses for resolution. BGP in each shard gets reachability notifications for registered IP addresses. These notifications populate the local BGP neighbor reachability information of the shard and trigger routing updates local to a shard. These updates can activate a BGP route when a next-hop is reachable, change the BGP route if reachability changes, and inactivate the BGP route if a next-hop is unreachable.

VPN and Sharding

You can apply sharding to *Virtual Route Forwarding* (VRF) and *Virtual Private Network* (VPN)^[3] routing tables also. Each shard thread hosts a slice of the VPN and VRF RIB table as determined by the hash. The *Route Distinguisher* (RD) of VPN routes is excluded from hash calculations to allow routes with the same prefixes but different RDs to be correctly processed in the associated shard. VPN label allocation is a central service because a single pool of *Multiprotocol Label Switching* (MPLS) labels is generally available to send out. A shard thread that wants an MPLS label for a VPN route requests this centralized service for labels. Target routes—needed for VPN processing—are stored in a separate RIB. Since this RIB is smaller, it is duplicated in all shard threads.

Large Peer Groups and Update-IO Parallelism

A peer group is usually assigned to an Update-IO thread that manages packing and generation of updates for that group. Multiple groups get assigned to different update threads for parallelism. In certain use cases, one or more large peer groups can include a very large number of BGP peers, and as a result we would not be able to distribute the load of such peer groups over several Update-IO threads effectively. To handle such a special case, we split such configured peer groups into several logically split peer groups. Each split group is allocated a subset of peers from a large peer group and assigned to an Update-IO thread.

BGP Graceful Restart Handling with RIB Sharding

BGP *Graceful Restart* (GR)^[4] processing requires sending of an *End-of-Rib* (EOR) notification after initial download of routes to a peer that is coming up after a failure. Shard threads contribute to the initial download of routes to a peer independently and in parallel. But the EOR message is sent in a coordinated fashion from the main thread after all shard threads complete the initial download of their slices of RIB to the peer. Likewise, consumption of an inbound EOR requires coordination from shard threads. The inbound EOR message is sent to all shard threads.

BGP Optimal Route Reflection and Sharding

You can configure BGP *Optimal Route Reflection* (ORR)^[5] with *Intermediate System-to-Intermediate System* (IS-IS) and *Open Shortest Path First* (OSPF) on a Route Reflector to advertise the best path to the BGP ORR client groups. Configuration is done by using the IGP metric after calculating the *Shortest Path First* (SPF) from a client's perspective. For BGP ORR to work in a Route Reflector, BGP requires assistance from an IGP implementation to calculate an IGP metric for a prefix from a client's perspective. In a sharded BGP architecture, ORR must be built as a service like a resolver. Shard threads register ORR reachability to this service, and in turn receive notifications about reachability of the registered prefixes.

Configuration, CLI Show Commands, Telemetry, SNMP, and Sharding

The shard architecture requires that each shard thread processes configuration independently of other shard threads. This way each shard has its own configuration state to work with. To present a consistent view to the user, some of the *show* commands in a router must collate information about RIBs from all shard threads and present a system view of the RIBs to the user. The same is true about telemetry streaming and the *Simple Network Management Protocol* (SNMP) get/walk of tables in a shard. For debugging purposes, we also implemented a view into the RIBs of each shard thread.

Routing and BGP RIB Sharding

Figure 3a illustrates an approximate high-level view of how routing protocols were implemented in the JUNOS *Routing Protocols Daemon* (RPD) before BGP RIB sharding was implemented. It shows a single thread that runs all routing protocols (including BGP), and INFRA, which includes interfaces, routing tables, next-hop tables, route resolution, and FIB programming. When BGP RIB sharding and Update-IO are configured, then additional threads are spawned, as shown in Figure 3b.

Figure 3a: Single Threaded Routing Protocol Daemon

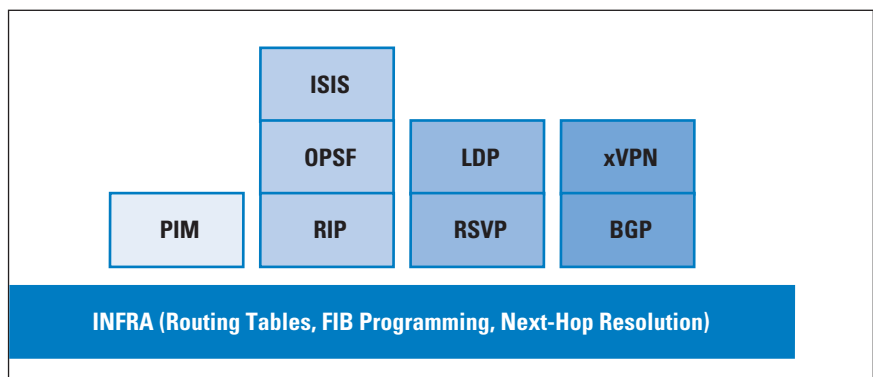
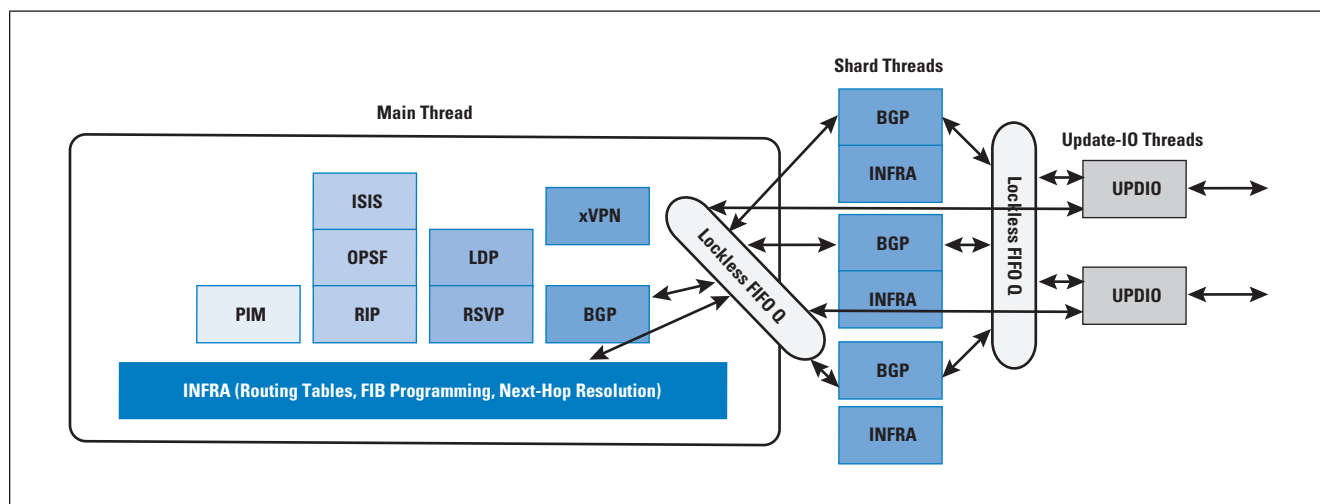


Figure 3b: Shard Routing Protocol Daemon



Main thread shards and Update-IO threads do not share any mutable state, and as a result no locks or critical sections are executing between these threads. As explained earlier, shard threads do the read-side processing of incoming messages in parallel and Update-IO threads do write-side processing of outgoing messages in parallel. Shard and Update-IO threads like the main thread process the routing configuration independently in a lockless manner. Coordination of peer state, interface state, and RIB routes and counters between main and other threads is done via *Interprocess Communication* (IPC) messages. These messages are sent over an IPC channel consisting of a pair of lockless *First-In, First-Out* (FIFO) queues. To allow for state compression, IPC channel queue depths are kept small. Socket read/write readiness is used for IPC channel back pressure between threads. Overall, all message passing between threads is very fast and efficient.

The following message types are sent from the main thread to the shard threads:

- Interface state, such as link information, address families configured, IP address, and state of the links is used by routing tables code and the BGP protocol. As a result, each shard has its own copy of interface state information.
- Route messages for non-BGP routes are distributed to a shard that is found via hash computation on the IP prefix of the route. Also, the entire route target RIB is sent as route messages to all shard threads.
- Configuration indications are sent as messages to signal availability of the updated configuration database.
- When *show* commands, SNMP *get* requests, and telemetry requests are sent, responses from shards are sent to the main thread to service the *Command-Line Interface* (CLI), SNMP requests, and telemetry streaming.
- BGP peer state transition messages are sent from the main thread to the shard threads. A handler in shard threads takes appropriate actions on BGP peer objects local to the shard.
- Next-hop resolution messages are sent from the resolver service in the main thread to shard threads. Shard threads register with this service for a BGP next-hop.
- VPN label messages are sent to shard threads whenever the main thread receives a request from a shard.

BGP in the main thread also shares state with Update-IO threads via IPC messages. These messages follow:

- The BGP peer *Finite State Machine* (FSM) in the main thread sends peer state messages to Update-IO threads. Such messages result in state changes for BGP peer objects in the Update-IO threads.
- *Route tuple messages*: Outbound BGP route messages are sent from the main thread to one of the Update-IO threads. Update-IO threads consume these messages to generate BGP UPDATE messages for a set of peers.

Shard threads send the following messages to the main and Update-IO threads:

- Route messages are sent from the shard to send active BGP-only routes to the main thread. The main thread adds these received routes to its RIB.
- Shards send *show*, SNMP, and telemetry responses to the main thread to assist in supporting these administrative functions.
- Shards send resolver registration requests for BGP next-hops, and they get resolver responses from the main thread.
- VPN label requests are sent to the main thread to get a pool of labels to support VPN functions in the shard threads.
- Outbound route tuples are sent to Update-IO threads per peer group.

When a shard thread receives UPDATE messages, it does the PDU processing on the prefixes that it owns and ignores those it does not. This processing uses the same hashing scheme that Update-IO threads use for distributing inbound prefixes. As a result, BGP routes are added to the RIBs of a shard thread. If the active route for a prefix in a shard is a BGP route, then it must be added to the FIB. Queueing points are used on the shard send side (before the IPC) to dampen the route churns due to link and peer flaps. The active route is added to the queue and then distributed (via IPC) to the main thread.

The main thread can program the FIB. Any subsequent state changes associated with BGP prefixes in shard threads may result in changes to the current active route in the shard. When this active route changes or is deleted in a shard, a new message is queued to be sent to the main thread.

The main thread centralizes answering BGP next-hop resolution, and it also programs the FIB. Any changes in the next-hop reachability are announced to shard threads. Shard threads react to such announcements and cause the necessary state changes in the Update-IO threads and the active route redistribution to the main threads. Update-IO threads generate outbound UPDATE messages and announce the prefix changes to their neighbors. In the end all threads, the FIB, and the network will have the consistent view of any prefix.

When a shard thread has run the BGP export policy and decided which BGP peers across all peer groups will receive a route, a new set of IBGP route announcement tuple messages (tuples) are queued for Update-IO threads. Shard threads multiplex several tuples back to an Update-IO thread. These tuples carry state about the BGP PATH attributes and prefixes. These PATH attributes are assumed immutable and are shared in a reference-counted manner between threads.

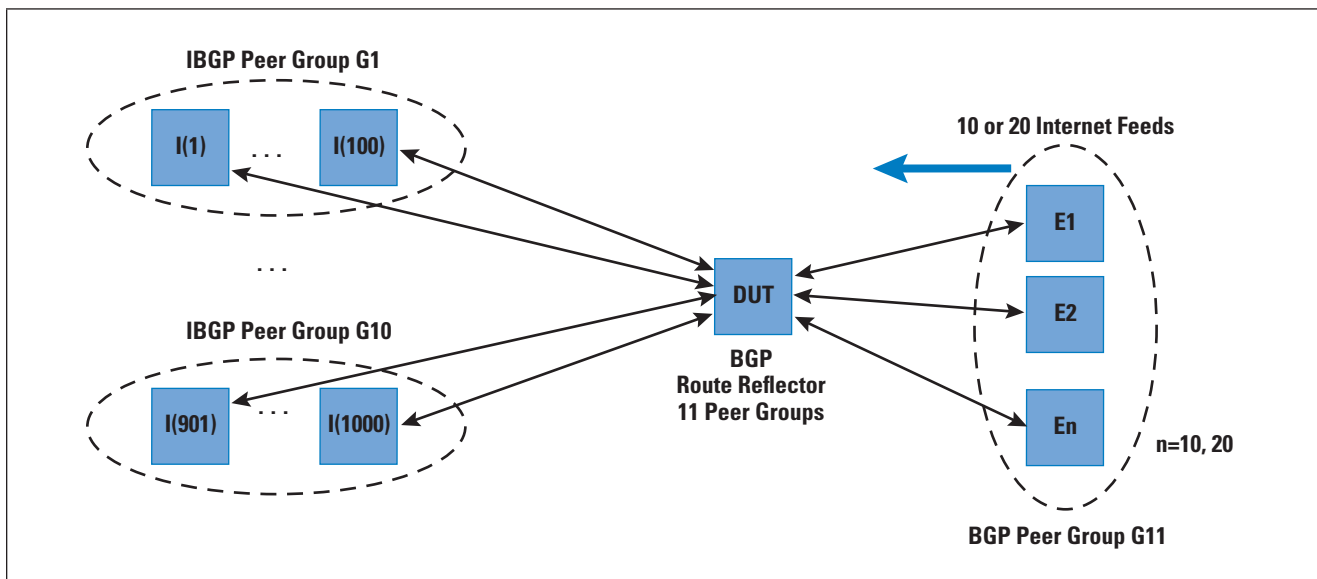
The Update-IO thread puts many tuples with the same BGP PATH attributes in a single BGP UPDATE message for transmission. Packing helps ensure that the number of BGP UPDATE messages is limited to prevent flooding of the local TCP/IP stack with too many send calls. Also, packing ensures that downstream BGP routers are not flooded with excessive messages. It ensures that gains from concurrent read processing in the shard threads to CPU overhead are not lost by handling more I/O messages.

Various queueing points are present in the main, shard, and Update-IO threads in our routing implementation. As these routes are learned and selected to distribute, the queueing points help us compress fast occurring state changes and reduce the churn from spreading from a source to the consumers. This churn reduction is very important when operating at scale and with several producers and consumers. Shard threads compress route changes for prefixes being redistributed to the main thread. For example, a route addition followed by a route deletion nullifies both. The main thread has a similar queueing logic to compress state churn when downloading routes to the FIB. Shard threads have a queue towards Update-IO threads to dampen the tuple churn in the BGP. Similarly, the Update-IO threads also use tuple queues to pack and suppress state before final state is disseminated to the peers.

Performance Measurements

To measure the gains of the approach described previously, we used a testbed where the *Device Under Test* (DUT) is a scaled BGP Route Reflector (Figure 4).

Figure 4: Route Reflector Topology



The RR receives several Internet feeds (800,000+ routes per feed) from an *External BGP* (EBGP) peer group of 10 BGP peers and reflects the best path for each destination towards 1,000 BGP client peers divided among 10 peer groups. Thus, a total of 11 BGP peer groups are in this testbed.

This setup represents a high scale of incoming UPDATE messages, and even higher scale of outbound UPDATE message generation. We used Internet feed instead of canned routes from a protocol tester like IXIA to mimic the real-world scenario where routes have variable PATH attributes. The DUT is a Linux Ubuntu with 18 servers with 8 cores (Xeon CPU E5-2640 v3 @ 2.6 GHz) and hyperthreading is turned off. The RR server on the DUT runs in a Linux Docker Container. The RR clients were also running our BGP implementation, where these clients drop all incoming PDUs after checking them for validity. This modification was done to ensure these 1,000 clients are never a bottleneck when receiving UPDATE messages from the DUT.

Table 1. Read-Side Performance Measurements

Number of Shard Threads	Convergence Time (Seconds)	Scale Up (S)
0	73	1.00
1	76	0.96
2	42	1.74
3	31	2.35
4	26	2.81
5	20	3.65
6	18	4.06

Table 2. BGP Write-Side Performance (4 Shards)

Number of Update-IO Threads	Convergence Time (Seconds)	Scale Up (S)
0	352	1.00
1	259	1.36
2	141	2.50
5	67	5.25
10	54	6.52

Table 3. BGP Full Pipeline Performance (4 Shards)

Number of Update-IO Threads	Convergence Time (Seconds)	Scale Up (S)
0	347	1.00
1	282	1.23
2	161	2.16
5	93	3.73
10	77	4.51
11	78	4.45

The testbed was run in three modes:

- Add a discard export policy to the DUT to ensure that no tuples to Update-IO threads are generated and we can get the measurements of the read side of the BGP pipeline. Table 1 shows the performance numbers for this mode of measurement. We noticed a maximum scale-up of 4x on six shards, and beyond six shards our gains were not beyond 4x.
- After the DUT has all the inbound Internet feeds and the RIBs have settled down, we delete the discard export policy, and this deletion triggers generation of tuples towards Update-IO threads and update messages start streaming toward the 1,000 peers. This step emulates the write side of the BGP pipeline and helps us measure that performance. Table 2 presents the measurements for this mode with 4 shard threads. The numbers of Update-IO threads were chosen to ensure 10 peer groups map to an even number of threads. We noticed gains linearly increasing as the number of threads increased.
- To measure the full pipeline performance, we repeated the previous 2 modes of testing without any discard export policy. We ran tests with 4 shard threads and up to 11 Update-IO threads. As expected, we noticed several fold improvements in the performance of BGP convergence time as we increased the number of I/O threads, as shown in Table 3.

Conclusions

Sharding in databases, web servers, and parallel computing concepts of SPMD have been used many times in the industry. This article is the first one where both concepts are used with BGP. We have designed and implemented a solution of splitting BGP into read and write pipelines and data (the RIBs). Our technique of sharing nothing among threads implementing BGP read and BGP write pipelines yields significant gains in scale and performance without impacting the convergence properties of the BGP protocol. This design also maintains BGP prefix packing and reduces the impact on local and remote routers. Lastly, we have implemented a design that keeps providing performance gains as parallelism increases, but the gains are limited by Amdahl's Law^[6].

References and Further Reading

- [1] Yakov Rekhter, Susan Hares, and Tony Li, "A Border Gateway Protocol 4 (BGP-4)," RFC 4271, January 2006.
- [2] Tony Bates, Ravi Chandra, David Katz, and Yakov Rekhter, "Multiprotocol Extensions for BGP-4," RFC 4760, January 2007.
- [3] Eric Rosen and Yakov Rekhter, "BGP/MPLS IP Virtual Private Networks (VPNs)," RFC 4364, February 2006.
- [4] Srihari R. Sangli, Enke Chen, Rex Fernando, John Scudder, and Yakov Rekhter, "Graceful Restart Mechanism for BGP," RFC 4724, January 2007.

- [5] Robert Raszuk, Bruno Decraene, Christian Cassar, Erik Aman, and Kevin Wang, “BGP Optimal Route Reflection (BGP ORR),” RFC 9107, August 2021.
- [6] Amdahl, Gene M., “Validity of the Single Processor Approach to Achieving Large-Scale Computing Capabilities,” *AFIPS Conference Proceedings* (30): 483–485, 1967.

SANJAY KHANNA holds a B.E. from Delhi University, a M.S., and a Ph.D. from Old Dominion University. Since 1993 he has worked in several IP networking-related jobs at IBM, Ericsson, Extreme Networks, and Juniper Networks. For the last several years he has been working on modernizing Juniper’s routing stack. He is a member of ACM. He can be reached at: skhanna@juniper.net

JAIHARI LOGANATHAN has been working in the networking industry since the early 1990s. He has a bachelor’s degree in computer science. He has worked on a variety of networking equipment and solutions from access modems, switches, security gateways, data center fabric, to core routers. His career spans several networking and cloud startups. He is a subject matter expert in many things networking. He has also helped start companies in search and networking space. He is currently working as a Distinguished Engineer at Juniper networks. He can be reached at: jlogan@juniper.net.

ASHUTOSH GREWAL holds a B.Tech. from the National Institute of Technology Durgapur and an M.S. from North Carolina State University. Since 2012, he has worked in the JUNOS routing protocols group at Juniper Networks as a Software Engineer on a variety of BGP, routing, and networking-related projects. He can be reached at: agrewal@juniper.net

Check your Subscription Details!

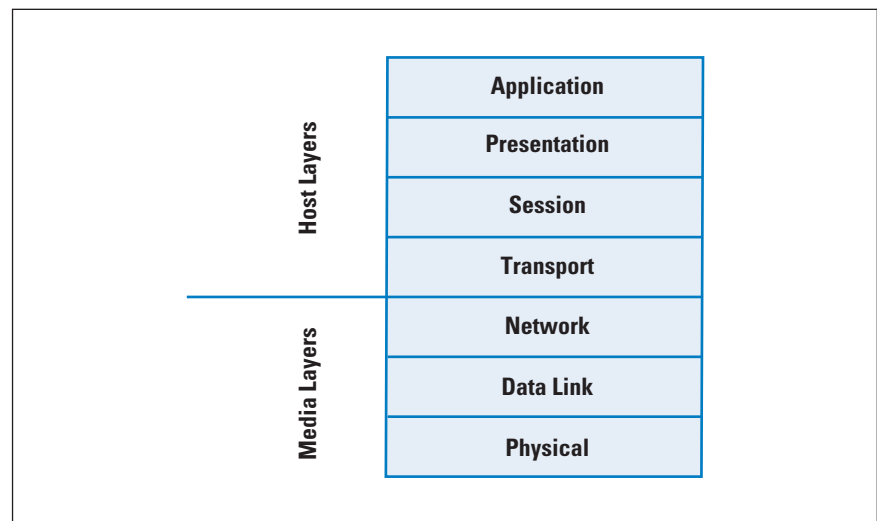
If you have a print subscription to this journal, you will find an expiration date printed on the back cover. For several years, we have “auto-renewed” your subscription, but now we ask you to log in to our subscription system and perform this simple task yourself. Make sure that both your postal and e-mail addresses are up-to-date since these are the only methods by which we can contact you. If you see the words “Invalid E-mail” on your copy this means that we have been unable to contact you through the e-mail address on file. If this is the case, please contact us at ipj@protocoljournal.org with your new information. The subscription portal is located here: <https://www.ipjsubscription.org/>

Transport Versus Network

by Geoff Huston, APNIC

One of the basic tools in network design is the so-called “stacked” protocol model. This model was developed in the late 1970s as part of a broader effort to develop general standards and methods of networking. In 1983, the efforts of the *Consultative Committee for International Telephony and Telegraphy* (CCITT) and *International Organization for Standardization* (ISO) merged to form *The Basic Reference Model for Open Systems Interconnection*, usually referred to as the *Open Systems Interconnection Reference Model*, or the “OSI Model”^[0]. This model included a seven-layer abstract model of networking that defined standard behaviours of both the overall network functions and the various components of the network (Figure 1).

Figure 1: OSI Reference Model



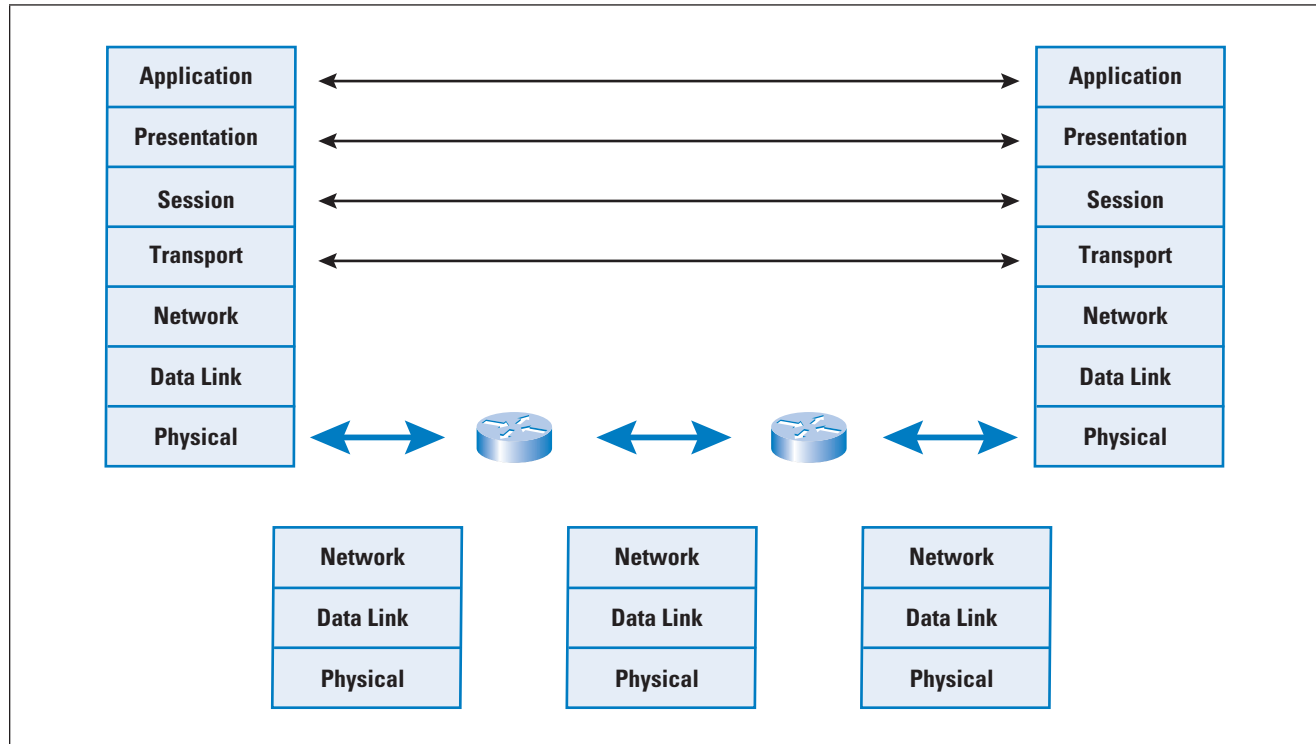
The model segmented functions into two parts:

- The *Media Layers* handle the encoding of binary data into the physical transmission media, the data link layer describes the data frames used between two inter-connected “nodes,” and the network layer manages a multi-node network, including addressing and routing behaviours that manage the transmission of data between attached hosts.
- The *Host Layers* concern functions on end hosts. These layers encompass the transport layer that performs data segmentation into packets, end-to-end flow control, packet loss recovery, and multiplexing. The layers above the transport layer in the OSI model are the session layer, the presentation layer, and the application itself.

In a model of a network as a collection of interior *nodes* and a set of attached *hosts*, the interior nodes use only the network layer to make forwarding decisions for each packet that it handles, while the hosts use the transport layer to manage the data flow between communicating hosts.

The implication of this model is that there is a delineation between node and host functions and a clear delineation of the data that they need to perform their functions. Nodes do not need to have any knowledge of the settings used at the transport level, and, similarly, hosts have no need to access the network layer (Figure 2).

Figure 2: Host and Node Functions



In the context of the Internet Protocol Suite, the network-layer function is encoded as the IP header of a data packet, and the transport-layer function is encoded as the transport header, conventionally as either a *Transmission Control Protocol* (TCP) or a *User Datagram Protocol* (UDP) header (although IP also defines other headers). In terms of the Internet architectural model, a packet should be deliverable through the Internet no matter what transport header it may have attached to it.

Therefore, it should not matter in the slightest what value you put in the IP Protocol field in IP packet headers. It's really none of the network's business!

The same almost applies to the *Extension Header* fields in IPv6, although this area is one where, inexplicably, IPv6 scrambled the egg and some extension headers are addressed to network elements, namely the *Hop-by-Hop* and *Routing* Extension Headers, while the remainder are ostensibly addressed to the destination host. If Extension Headers were defined exclusively as host (or destination) extensions, then the IPv6 networks should ignore them, while if they were intended to be network options, then hosts should ignore them. Perhaps it's another of those areas where theory and practice just don't align well.

In a strict sense the Protocol field in the IPv4 packet header need never have been placed in the IPv4 header in the first place. The particular transport protocol that the communicating hosts use is none of the network's business, in theory, meaning that if the two communicating hosts decide to deliberately obscure the transport protocol control settings from the network, then that should not matter in the slightest to the network.

In today's public Internet it appears to matter a lot that the transport protocol header is visible to the network. In fact, not only should the transport protocol be visible to the network, but the particular transport protocol that the hosts select also matters to the network. That is because many elements of today's network not only peek into the transport headers of the packets that they carry, but they also rely on the information in this transport header. Firewalls are a classic example of this reliance, but there are also *Network Address Translators*, *Equal-Cost Multi-Path Load Balancers*, and *Quality-of-Service Policy Engines*, to name a few. These network functions make assumptions about the visibility of transport headers in the IP packet in order to make consistent decisions about packet handling for all packets within a single transport flow. Often these network functions take it one step further and they process packets with the well-known transport headers (typically restricted to just TCP and UDP) and discard all else. It's even gone further than that, and we have reached the point that today's generally accepted rule is that unfragmented IP packets that contain a TCP transport header that includes one end using Port 443, and unfragmented IP packets that contain a UDP transport header where one end uses Port 53 stand the best chance of getting their data payload through to the intended destination. Every effort to augment this remarkably constrained set of packet profiles increases the probability of network-based disruption of communication.

Encrypted Transport Headers

If it is a self-limiting action to use a novel transport protocol in the public Internet, then why are we even considering the option of encrypting transport protocols to make all transport headers opaque to the network?

One answer is "Edward Snowden." When Snowden made his pervasive monitoring revelations^[1], the *Internet Engineering Task Force* (IETF) responded in what could be called a "like-for-like" reaction and came to a consensus position that "Pervasive Monitoring Is an Attack"^[2]. The general response to this form of insidious attack was to increase the level of encryption of Internet traffic to lift the degree of difficulty in carrying out network-based surveillance. Not only does this IETF response encompass the use of *Transport Layer Security* (TLS) to encrypt session payloads on the Internet wherever possible and shift the application behaviour profiles to make this the default action, but it also shifted our attention to other areas of Internet communication where compromise of the trust model was thought to be an issue.

The actions of the *Domain Name System* (DNS) protocol have been drawn into this IETF universal obfuscation effort, as has the transmission of transport protocol headers. We are long past the time when hosts were ill-equipped to perform encryption functions, and now robust encryption is not a luxury option with limited use, but rather something every user should reasonably expect to use as a minimum requirement. If the objective is to limit the information leakage in all aspects of the communications environment on the Internet, then the control meta-data is as important as the data itself. Applying confidentiality to transport header fields can certainly improve users' privacy and can help mitigate certain attacks or manipulation of packets by devices on the network path^[3].

However, I suspect that this privacy argument is only one part of the story, and while these measures to encrypt Internet traffic play to a popular concern of the surveillance state operating in a largely unchecked manner, they may not lie at the heart of why obscuring host functions from the network is a path that some parts of the Internet ecosystem vigorously pursue today with transport header encryption.

It's not clear that the objective is here, and as with all interdependent complex systems, deliberately obscuring one aspect of the system from another typically offers both benefits and downsides. In April 2019, the *Internet Architecture Board* (IAB) published RFC 8546, titled "The Wire Image of a Network Protocol"^[4]. It's a short document (9 pages) by today's RFC standards, but brevity does not necessarily imply clarity. This document appears to have cloaked its message in such a dense level of abstract terminology that it managed to say very little of practical use! The IAB document appears to have been prompted by the protracted debate in the QUIC Working Group over the use of the visible spin bit in the QUIC transport protocol^[5], and I suspect that it started as an effort to argue for some levels of transport behaviour visibility to the network, but the IAB's prognostications on the topic have offered little useful or informative comment apart from illustrating the level of collective angst that this issue has generated! The IAB is not the most prolific of commentators, and any matter that provokes an IAB response, no matter how cryptic that response may be, does illustrate that the topic is one of general concern rather than being just a rather esoteric tussle buried deep down in the design of a particular protocol.

This topic of encrypted transport headers is a transport topic, so it is natural to ask whether the IETF's *Transport Area* can do any better than the IAB in providing a clear and informed exposition of the issues here. The Transport Area Working Group of the IETF has completed an RFC on this topic, "Considerations around Transport Header Confidentiality, Network Operations, and the Evolution of Internet Transport Protocols"^[6]. To quote from its abstract: "This document discusses the possible impact when network traffic uses a protocol with an encrypted transport header. It suggests issues to consider when designing new transport protocols or features."

This document strikes me as an effort to produce a slightly more practically focused commentary on header encryption than the earlier IAB effort. At 49 pages it certainly cannot be considered a brief document, but does this extended commentary do any better in terms of clarity of the arguments being considered?

The document first looks at some rationales for the use of information on the network contained in headers. It cites the situation of link aggregation, and the problem of packet re-ordering in such scenarios. The common response to re-ordering is for the network to peer down into the transport header to gain a more granular view of a traffic flow than that which can be derived from source and destination IP address pairs. It is the IPv4 proxy for the IPv6 *Flow Label*. (Although the IPv6 Flow Label is so confused as to its intended role it's hard to understand how the IPv6 Flow Label field is useful in any case whatsoever!)

The document references differential service efforts that attempt to perform selective damage on traffic flows under the guise of “Quality of Service.” (That “Quality” label always seems to me to have an Orwellian connotation, and a more honest label would be “Selective Service Degradation,” or even just “Carriage Standover Services”). The document also enumerates the ways network operators can perform network analysis of using transport-level information, including traffic profile analysis, latency, and jitter and packet loss. However, the document strikes me as presenting a somewhat disingenuous set of rationales. For me, it is akin to a voice telephony operator justifying its eavesdropping on phone conversations on the basis of a baseless assertion that the information gathered by such wiretapping, or in other words knowledge of what people are saying to each other over a telephone connection, can be used to make the telephone network better! The document also uses the last recourse of the desperate, by invoking a nebulous concept of “security,” claiming that if network operators were no longer able to eavesdrop on the transport parameters of active sessions, then somehow the operator’s ability to run a secure network would be compromised in some unspecified way.

Obviously, none of the rationales presented in this document can withstand much in the way of close scrutiny.

It also appears to take a privacy-oriented stance in its analysis, and it seems to me that the privacy argument is largely an overt excuse for a more substantial difference of opinion between *content* and *carriage*. To a large extent, the issue from the perspective of the application is that the efforts of network operators to perform “traffic grooming” through transport header manipulation amounts to little more than inflicting damage on application data flows, and thereby pushes the network to a lower level of carriage efficiency. And this issue of the use of networks to selectively degrade transport performance in the name of network service quality is perhaps where we should look for the real tensions between networks and hosts in today’s Internet.

Transport Protocol Meddling

To look down this path we might want to start with the tensions between hosts and networks on the Internet.

In the telephone world, the network operator controlled all traffic. What you leased from the network was either a virtual circuit capable of passing a real-time voice conversation, or a fixed-capacity channel between two end points. If you used one of these channels, you couldn't go any faster than the contracted speed, and if you went slower, you did not release common capacity for anyone else to use. Obviously, the network charged more for leases of higher capacity. Packet networks changed all that. The network had no enforcement, and various applications (or traffic flows) competed with each other for the common transmission resource. Networks that wanted to control the allocation of shared common communications resources to clients had a problem.

This allocation control was the motive for a large body of work on the Internet during the 1990s and 2000s over what was called *Quality of Service* (QoS)^[7]. The network operator wanted to offer (no doubt for some premium) a “higher-quality” service to some clients and some traffic profiles. But if a network has a fixed-capacity offering a larger slice of the network resources to some clients, inevitably it will offer less capacity to the others. One common theme of much of this work was that while it was possible for the network to disrupt a communication session in various ways to make it go slower, it was a lot more challenging (or even impossible in many ways) to make a session go faster.

Thus, in order to offer preferential treatment to a class of traffic flows, a good way was to make all the other flows go slower! The intended effect was to clear some space for sessions that were intended to be favoured to expand their sending windows and occupy this cleared network space. So-called *Performance-Enhancing Proxies* were not really able to make the selected TCP sessions go faster per se, but they were able to make other concurrent TCP sessions go slower, and thereby make some space for the selected sessions to have a lower packet-loss probability and hence achieve a higher data-throughput rate. One way of using this form is session throttling to drop packets. A subtler way, but also very effective, is to alter the TCP control parameters. If the offered TCP window size parameter is reduced, then senders will conveniently throttle their sending rate accordingly.

Pretty obviously, this selective behaviour of throttling active TCP sessions by networks was not something that applications viewed as a sympathetic act, and there have been two major responses from the application side. One is the use of a different congestion-control algorithm that is a lot less sensitive to packet loss and more sensitive to changes in the end-to-end bandwidth delay product across network paths. This method is called the *Bottleneck Bandwidth and Round-trip* (BBR) TCP control protocol, which is a relatively new TCP sender-side control algorithm.

But BBR is still susceptible to on-path manipulation of the TCP window size, and protecting the session from this form of network interference is where encrypted transport headers emerged and became an important objective. This response is the second one, executed by obscuring where the TCP control information is actually carried in the packet.

As we've already noted, you just can't remove a visible transport header from IP packets in the Public Internet, and even encrypting the TCP header would probably incur the same drop response from the network. But hosts have the option to ignore these transport header settings. So, while the host can't remove a visible transport header, they can make the headers meaningless.

One option is to use a "dummy" outer TCP wrapper as fodder for networks that want to peek at the transport layer and manipulate the session settings while hiding the real TCP control header inside an encrypted payload. There would be little in the way of a visible network signature that this manipulation is happening, apart from the observation that the TCP end hosts would appear to be unresponsive to manipulation of their window parameters.

However, the problem with this approach is that these days the application is actually trying not only to take control over its transport session parameters from a meddling network, but also to assert the same control over the platform in which the application is hosted. In theory, the application could use "raw IP" interfaces into the platform I/O routines, but in practice in deployed systems it is close to impossible. Platforms used in production systems tend to treat applications with suspicion. (Given the proliferation of malware, this level of paranoia on the part of the platform is probably warranted.) It is quite a challenge to disable all forms of how the platform handles the transport protocols and pass control of the transport protocol from the kernel into the applications space.

For this reason, it is logical to take the approach QUIC uses, where the shim wrapper of QUIC uses UDP as a visible transport header and pushes the TCP header into the encrypted payload part of the IP packet. UDP is close to ideal in this case as there are no transport controls in the protocol, just the local port numbers. QUIC looks to the network a lot like a UDP session that uses a TLS-like session encryption because in so many ways it is a UDP session that uses TLS. The change is that the end-to-end TCP flow control is now truly an end-to-end flow because only the two applications at the "ends" of the QUIC transport can see end-to-end transport-control parameters that are embedded in the end-to-end encrypted UDP payload. The host platform control over UDP packets is perfunctory, and the application is then allowed to assume complete control over the transport behaviour of the session.

Content Versus Carriage

Perhaps this shift to opaque transport headers goes a little further than just a desire for greater levels of protected autonomous control by applications. The shift that QUIC represents could be seen as the counter move by content providers to another round of a somewhat tired old game play by network operators to extract a tax from content providers by holding their content traffic to ransom, or, as it came to be known, a tussle over *Network Neutrality*.

There have been times when network operators have implemented measures to throttle certain forms of traffic that they asserted was using their network in some vaguely unspecified manner that was “unfair” in some way. The vagueness of all this discussion is probably attributable to a baser desire on the part of the carriage operator, which was to extort a carriage toll from content providers in a crude form of basic blackmail: “My network, my rules. You customer, you pay!”

I suspect that many carriage providers in this industry, who are witnessing the content providers take all the money off the table, believe that they are the victims here. Their efforts to restore some of their lost revenue base has meant that they are looking to restore a “fair share” of revenue in forcing the giants of the content space to pay for their share of carriage costs. However, if the enforcement mechanism of this extortion pressure is through playing with the transport-control parameters of the traffic that transits the carriage network (or, in other words, holding the traffic to ransom), then the obvious response is to push the transport controls under the same encryption veil as the content itself to prevent such real-time manipulation of the traffic profile. And this explanation of why QUIC is so important is perhaps a more compelling one.

If this situation is a tussle for primacy in the tensions between carriage and content, then it looks like the content folks are gaining the upper hand. Through encryption at every level in the host part of the protocol stack, including at the transport layer, the content folks are withholding information from the carriage providers that would allow the carriage providers to selectively discriminate and play content providers off against each other. If all that the network can do is limited to fully encrypted UDP packet streams, then one stream looks much like another, and selective discrimination is just not feasible. And if that’s not enough, then padding and deliberate packet variation can blur most efforts at traffic profiling.

But when I say “content” I really mean “apps,” and when I say “apps” I actually mean “browsers,” so in reality I am really talking about Chrome, and when I say Chrome, I mean Google.

The massive dominance of mobile traffic in the industry and the massive dominance of Android in the mobile device environment tilts this space to an extraordinary degree.

Given this inherent level of control of all mobile devices, coupled with control of the majority browser platform in this space, it is hard to conceive how Google could possibly lose in this tussle. However, it is likely that if Google wins this particular battle with the carriage providers, there will be further battles to come. It is highly likely that the carriage industry will follow the lead from traditional print media and head to politicians with the case that Google's destruction of the business model for the provision of national communications infrastructure is counter to national interests, and political intervention is necessary to restore some balance into the market and allow the market for carriage to be a viable investment vehicle. Or, to put in more crudely, if Google has destroyed the residual value of the contained carriage market, then Google should now pay carriage operators to restore its viability.

At this point all technical considerations of encryption and information leakage, and even all market considerations of the viability of various business models, just walk out the door, and in their place comes a bevy of lawyers and politicians. Strategic national interest is always a strong argument to make, and when we get over the various nebulous threats by actors to quit national markets, we then get down to the real question of: "What is a tenable business relationship between carriage and content?"

In such a politically charged space the choices at that point are either that the various market players will compromise and reach some outcome that they can all live with, or the politicians will attempt to impose an outcome that will in all likelihood be far more disagreeable for all!

Whatever the outcome in the next few years, it should be fun to watch this drama play out. Don't forget to bring popcorn!

References and Further Reading

- [0] Wikipedia, "OSI model,"
https://en.wikipedia.org/wiki/OSI_model
- [1] Cullen Jennings, Brian Trammell, Christian Huitema, Bruce Schneier, Ted Hardie, Richard Barnes, and Daniel Borkmann, "Confidentiality in the Face of Pervasive Surveillance: A Threat Model and Problem Statement," RFC 7624, August 2015.
- [2] Stephen Farrell and Hannes Tschofenig, "Pervasive Monitoring Is an Attack," RFC 7258, May 2014.
- [3] Al Morton and Kathleen Moriarty, "Effects of Pervasive Encryption on Operators," RFC 8404, July 2018.
- [4] Brian Trammell, "The Wire Image of a Network Protocol," RFC 8546, April 2019.
- [5] Geoff Huston, "Just One Bit," *The ISP Column*, March 2018.
<https://www.potaroo.net/ispcol/2018-03/onebit.html>

- [6] Godred Fairhurst and Colin Perkins, “Considerations around Transport Header Confidentiality, Network Operations, and the Evolution of Internet Transport Protocols,” RFC 9065, July 2021.
- [7] Geoff Huston, “QoS — Fact or Fiction?” *The Internet Protocol Journal*, Volume 3, No. 1, March 2000.
- [8] Geoff Huston, “A Quick Look at QUIC,” *The Internet Protocol Journal*, Volume 22, No. 1, March 2019.
- [9] Geoff Huston, “Anatomy: Inside Network Address Translators,” *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [10] Dave Oran, “Considerations in the Development of a QoS Architecture for CCNx-Like Information-Centric Networking Protocols,” RFC 9064, June 2021.

GEOFF HUSTON, B.Sc., M.Sc. A.M., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net



SIP phones have replaced traditional telephones even in the offices of The Internet Protocol Journal.

20 Years of SIP — A Retrospective

by Jonathan Rosenberg, *Five9*

June of 2022 marked the twentieth anniversary of the publication of the *Session Initiation Protocol* (SIP), documented in RFC 3261^[1,2,3,4]. When it was published in June of 2002, it set records for the longest specification produced by the *Internet Engineering Task Force* (IETF), at 269 pages. The IETF produces the technology standards that make the Internet work. Protocols like *Internet Protocol* (IP), *Transmission Control Protocol* (TCP), and *Hypertext Transfer Protocol* (HTTP)—all now part of mainstream vernacular—came out of the IETF. It was a monumental effort to produce, involving a dedicated author team that worked full time for months to ensure that the specifications were correct, consistent, and complete. I had the great fortune to be the lead author of this document, an accomplishment that was the defining moment for my career.

The RFC 3261 author team included Robert Sparks, Jon Peterson, Alan Johnston, Allison Mankin, Jonathan Rosenberg, Gonzalo Camarillo, and Henning Schulzrinne.

In the 20 years (and almost countless extensions to SIP) that followed, it is hard to dispute that SIP has been a major success. At its core, SIP enabled the transformation of the telecommunications industry from one based on hardware to one based on software—colloquially known as *Voice over IP* (VoIP). A 20th anniversary is the ideal time for a retrospective, to consider both its positives and negatives. On the plus side, this transformation resulted in the re-engineering of the phone network, the creation of new markets and market categories, and the creation of jobs and livelihoods. On the negative side, it has exacerbated the scourge of robocalling.

The Phone Network Re-Engineered

Prior to SIP, the telephone network was built using telephone switches based on custom hardware. These switches were made by a small set of vendors and were a completely verticalized solution from the physical networking layer to the application layer software.

With the mainstream adoption of IP networks, it became possible to replace that hardware with general-purpose computers running SIP-based software applications. This replacement resulted in a dramatic reduction of cost compared to the prior generation. SIP further reduced this cost by enabling these software applications to run on machines in a small number of data centers that might be far away from the people talking to each other, while still keeping the audio delays to a minimum. This centralization of the software was a dramatic shift in how the phone network worked.

This new paradigm had the most immediate impacts on the way corporate phone systems were built. Before SIP, businesses needed to put hardware-based phone systems called *Private Branch Exchanges* (PBXs) in each building and wire them up on a separate network from the IP network used for everything else.

SIP allowed enterprise IT departments to ditch this separate network and reuse the IP network for voice. It also allowed them to put the software in their data centers and eliminate the hardware in each building. This development represented a huge improvement in costs and reduction in complexity. The final “icing on the cake” was that SIP enabled video, instant messaging, and presence too, spawning the creation of desktop applications—called *softphones*—that allowed users to place calls, have video meetings, and chat. Businesses far and wide adopted these phones. Today, almost all business phone systems are based on SIP.

With its success in corporations, pressure grew for the phone companies (that is, the telcos), to provide a way for businesses to connect their phone systems to the rest of the phone network using SIP. Prior to this time, businesses could use software within their corporate campus, but they needed to switch to hardware to connect to the rest of the world. And so, “SIP Trunking” was born, providing a way to send and receive calls into the traditional phone network using SIP-based software applications. Its adoption was rapid, and it was the first step in transforming the edge of the telco networks from hardware to software.

Around the same time, mobile phone operators were seeing an explosion in usage due to smartphones. These mobile phones had two distinct wireless connections—an old one just for voice, and a new one for data. To expand capacity, they needed to reclaim the voice channel and use it for data. They could do it by switching the voice to VoIP, which would require them to replace their own voice hardware with SIP-based software. The wireless industry produced an expansive set of specifications on how to build a SIP-based replacement for mobile phone networks, called the *Internet Multimedia Subsystem* (IMS). IMS was finally deployed in the late 2010s. Today, most mobile phone calls use a SIP client built into the phone and traverse a SIP network deployed and operated by the mobile carriers. This change is largely invisible to mobile phone users, but not entirely. SIP also enabled the usage of higher quality wideband voice for phone calls, creating an audio experience that is more like listening to music, and you may have noticed this difference in more and more calls you make.

In a similar fashion, wireline telco providers saw a surge in demand for data. To make the jump to next-generation data access technologies like fiber, they needed to get rid of their separate voice networks and move to voice over IP too. Today, if you have one of these higher speed data networks and still have an analog phone in your home, the analog signal is converted to VoIP using a SIP client in the modem at your house, and then processed by a SIP network that the carrier operates.

The final piece of the puzzle is how carriers themselves connect to each other. This process has gradually migrated to SIP too, using carrier versions of SIP trunking. This change is now accelerating, since the conversion is needed to enable the deployment of *Secure Telephone Identity Revisited* and *Signature-based Handling of Asserted information using toKENs* (STIR/SHAKEN), a SIP-based technology to combat robocalling^[5,6,7].

Without a doubt, this transformation of the telecommunications technology stack—that SIP enabled—has massively impacted the world, enabling lower costs, more bandwidth for data, better quality for voice, and added video.

Market Category Creation

This transformation of the telecommunications industry also created entirely new markets and market categories that didn't exist before SIP. To enumerate just a few of them:

- *IP PBX*: The *IP Private Branch Exchange* (PBX) provides phone services for businesses. This market was created as a direct replacement for the legacy hardware-based PBX products that preceded it. Cisco Systems led this market, which never had a product in the PBX market, along with incumbents like Avaya, Siemens, and Nortel, many of which had legacy products along with the newer IP-based ones. This market is now itself shrinking, being replaced by *Unified Communications as a Service* (UCaaS).
- *SIP Trunking*: This market is estimated to be around \$13B in 2021^[8] and is a replacement for legacy hardware phone network access technologies.
- *SIP Hardphones*: Before SIP, the PBX vendors made their own phones, and a given phone could only work with their own hardware. With SIP, it became possible for vendors to produce phones that could work with many different IP PBXs. These phones were often produced at low cost. Vendors include Yealink, Cisco, Grandstream, and Avaya.
- *Session Border Controller* (SBC): The usage of SIP trunking drove demand for a new category of product that could serve as a SIP firewall of sorts, managing the boundary between an enterprise and a carrier, or between carriers. Ribbon and Oracle are the market leaders, with a market size estimated at USD \$709M in 2022^[9].
- *Internet Multimedia Subsystem* (IMS): Market leaders include Ericsson, Siemens, and Nokia. The market size was USD \$1.8B in 2019^[10].
- *Communications Platform as a Service* (CPaaS): This market category is an entirely new one, enabled by the transformation of telecommunications to software. CPaaS vendors offer *Application Programming Interfaces* (APIs) that allow developers to build telecom applications easily. These APIs allow for sending of SMSs, placing and receiving of phone calls, and so on. Twilio created this market and is still the market leader. SIP enabled the CPaaS vendors to gain low-cost and global access to telephone services, and without SIP, the market could not have existed. The market is huge and growing—estimated at USD \$5.2B in 2021^[11], (though most of it is for sending *Short Message Service* (SMS), where SIP has been less impactful).
- *Unified Communications as a Service* (UCaaS) puts the IP PBX in the cloud so businesses can consume voice and video communications services from the cloud.

Market leaders include RingCentral, 8x8, Cisco Systems, Microsoft, and Zoom. All of these vendors depend on SIP-based interconnection to the telephone network. This market is really big—estimated to be USD \$28.9B in 2021!^[12]

- *Contact Center as a Service* (CCaaS) enables delivery of contact center software from the cloud, including voice response systems, agent desktop applications, and call distribution software. Vendors include Five9, Gensys, and NICE/InContact. Like UCaaS, these vendors depend on SIP to interconnect to the telephone network. This market was valued at USD \$4.8B in 2021.^[13]

When put together, SIP created or enabled these (no less than eight) distinct markets, representing approximately USD \$50B in market value!

Job Creation

For me, the greatest source of satisfaction from the success of SIP is when I hear from someone that they have built their careers and their livelihood around this technology. SIP is complex, and like any complex technology that many vendors use in many ways in many markets, expertise in it becomes a marketable skill.

Many LinkedIn profiles list “SIP” as a skill. Many are software developers, but many other jobs require SIP expertise. SIP network engineers and technicians build, deploy, and operate SIP networks. Sales and marketing engineers configure and demonstrate SIP-based products. IT workers who manage business communications for their companies need to understand SIP too. A search on LinkedIn for people matching “SIP” yields approximately 239,000 results.

Many companies now exist that provide SIP certifications and training—for example the SIP School.^[14] SIP is taught in many graduate classes that cover computer networking, and some even have dedicated courses just on VoIP.

It’s hard to know how many jobs SIP has created, but it would not be unreasonable to guess it is somewhere in the ballpark of 100,000 jobs. If you add the folks working in technical roles across the companies in the markets that SIP created, along with those working in telcos or in IT departments providing VOIP, it is easy to see how the number could be that large.

The Downside: Robocalling

Almost all technologies that have brought great benefits have come with some drawbacks. There is no better example than the automobile, which has brought countless benefits, but also caused 42,915 deaths in 2021 due to automobile accidents. The Internet too, has brought countless benefits, but has also brought with it problems that are becoming more apparent. SIP has had far less impact as other technologies, so its drawbacks are fewer, but they do exist.

Without a doubt, the biggest drawback has been the rise of robocalling and the fake caller IDs that come with it. Telemarketing calls predate SIP for sure. However, as SIP reduced the costs of placing calls and made it possible to make calls using off-the-shelf software, it caused a sharp increase in the volume of these unwanted calls.

The problem is exacerbated by a design flaw in SIP—the lack of an authenticated caller ID. Without that, it is easy for callers to insert any phone number they want. This design defect was inherited from email, as SIP copied this aspect of its design from how email worked. After many years of failed attempts to resolve the problem, there is finally “light at the end of the tunnel” using a SIP-based technology called STIR/SHAKEN^[5,6,7].

In Conclusion

It’s been the highlight of my career to have had the fortune to be the lead author for a technology that, 20 years later, has had a profound impact on the world. By enabling the transformation of telecommunications from hardware to software, SIP drove a re-engineering of both mobile and wired phone networks that resulted in lower cost communications services and more bandwidth available for data. It brought video to the enterprise, created entirely new markets and some new market categories, and created at least 100,000 jobs. I try and remind myself of that fact every time I get one of those annoying robocalls.

References and Further Reading

- [0] This article was adopted from Jonathan Rosenberg’s blog:
<https://www.jdrosen.net/blog/20-years-of-sip-a-retrospective>
- [1] Jonathan Rosenberg, Henning Schulzrinne, Gonzalo Camarillo, Alan Johnston, Jon Peterson, Robert Sparks, Mark Handley, and Eve Schooler, “SIP: Session Initiation Protocol,” RFC 3261, June 2002.
- [2] Jonathan Rosenberg and Henning Schulzrinne, “An Offer/Answer Model with the Session Description Protocol,” RFC 3264, June 2002.
- [3] Henning Schulzrinne and Jonathan Rosenberg, “The Session Initiation Protocol: Providing Advanced Telephony Access Across the Internet,” *Bell Labs Technical Journal*, October–December 1998.
- [4] William Stallings, Session Initiation Protocol, *The Internet Protocol Journal*, Volume 6, No. 1, March 2003.
- [5] Numeracle, “STIR/SHAKEN: Everything you need to know about the FCC’s Call Authentication Framework,”
<https://www.numeracle.com/resources/stir-shaken-center>
- [6] IETF Datatracker, “Secure Telephone Identity Revisited (stir),”
<https://datatracker.ietf.org/wg/stir/documents/>
- [7] Metaswitch, “What are the STIR/SHAKEN Standards?”
<https://www.metaswitch.com/knowledge-center/reference/what-are-the-stir/shaken-standards>
- [8] The Business Research Company, “COVID-19 Impact On The Global Session Initiation Protocol (SIP) Trunking Services Market Outlook,” March 8, 2022.
<https://www.prnewswire.co.uk/news-releases/covid-19-impact-on-the-global-session-initiation-protocol-sip-trunking-services-market-outlook-878809709.html>

- [9] Future Market Insights, “Session Border Controller (SBC) Market Overview (2022–2032),”
<https://www.futuremarketinsights.com/reports/session-border-controller-market>
- [10] Fior Markets, “Global IP Multimedia Subsystem (IMS) Market Size to Expand Significantly of USD 8.26 Billion by 2027,”
<https://www.globenewswire.com/news-release/2022/05/19/2447271/0/en/Global-IP-Multimedia-Subsystem-IMS-Market-Size-to-Expand-Significantly-of-USD-8-26-billion-by-2027-Fior-Markets.html>
- [11] Future Market Insights, “Communications Platform as a Service (CPaaS) Market Outlook (2022–2032),”
<https://www.futuremarketinsights.com/reports/communications-platform-as-a-service-cpaas-market>
- [12] Fortune Business Insights, “Unified Communication as a Service (UCaaS) Market Size, Share & COVID-19 Impact Analysis, By Component (Telephony, Unified Messaging, Collaboration Platforms), By Delivery Model (Managed Services, and Hosted/Cloud Services), By Organization Size (Large Enterprises, SME’s), By Vertical (BFSI, IT and Telecommunications, IT-enabled Services (ITeS), Education, Retail and Consumer Goods), and Regional Forecast, 2021–2028,”
<https://www.fortunebusinessinsights.com/industry-reports/toc/unified-communication-as-a-service-ucaas-market-101934>
- [13] Fortune Business Insights, “Contact Center as a Service (CCaaS) Market Size, Share & COVID-19 Impact Analysis, By Function (Interactive Voice Response (IVR), Multichannel, Automatic Call Distribution, Computer Telephony Integration (CTI), Reporting and Analytics, Workforce Optimization, Customer Collaboration, and Others), By Enterprise Size (Small & Medium Enterprises and Large Enterprises), By Industry (BFSI, IT & Telecommunications, Government, Healthcare, Consumer Goods & Retail, Travel & Hospitality, Media & Entertainment, and Others), and Regional Forecast, 2022–2029,”
<https://www.fortunebusinessinsights.com/toc/contact-center-as-a-service-ccaas-market-104160>
- [14] The SIP School: <https://www.thesipschool.com/>
- [15] The SIP Forum: <https://www.sipforum.org/>

JONATHAN ROSENBERG is the Chief Technology Officer and Head of AI for Five9. He was previously CTO for Cisco Webex and Skype. He has been a frequent contributor to the IETF, with 72 RFCs. He’s the lead author of the Session Initiation Protocol (SIP) and related standards, such as ICE, STUN, TURN and SIMPLE. E-mail: jdrosen@jdrosen.net

IAB Comments on FCC Notice on Secure Internet Routing

The *Internet Architecture Board* (IAB), which provides oversight for the protocols and procedures used by the Internet and also handles the liaison management for the *Internet Engineering Task Force* (IETF), appreciates the opportunity to submit comments in response to the *Federal Communication Commission's* (FCC) Notice of Inquiry, "Secure Internet Routing"^[1]. The IETF is the main organization that works on standards relating to Internet technology. The mission of the IETF is to produce relevant technical documents that influence the way people design, use, and manage the Internet. The IETF is an open, diverse, global community of developers consisting of network operators, vendors, researchers and many other stakeholders.

The IETF originally developed the Internet protocol stack, including the routing system based on the *Border Gateway Protocol* (BGP), and continues to be responsible for maintaining and evolving the technical specifications that define the Internet and its protocols. The Internet's success has resulted from its flexible, modular architecture. BGP is the central protocol for providing global end-to-end connectivity across the world's heterogeneous network domains. It is fundamental to the operation of the Internet.

As in any protocol development, the adoption within the industry of new capabilities will vary. In recent decades, occurrences of BGP-related operational issues have increased. The existing BGP protocol stack is based on a design which can be extended, building on existing network investments. The IETF has two working groups dedicated to improving BGP interdomain routing, called *Inter-Domain Routing* (IDR) and *Global Routing Operations* (GROW). IDR is concerned with the correctness, robustness, and scalability of BGP. GROW is concerned with the operational problems associated with global routing systems, including measurement, policy, and security. The IETF will continue to evolve BGP to meet the needs of new network structures and applications, with a strong focus on security.

We believe in a continuous, modular, flexible evolution of the Internet and its protocols based on operational experience and requirements, where each service provider can determine their security needs based on their diverse requirements and in partnership with other providers. The success of future standardization efforts intended to increase routing security, will be highly dependent on educating BGP users about BGP operational issues and how well real-world deployment experience can be fed back into the multi-stakeholder standards development process, as opposed to a mandated top-down approach, which would fail to meet the diverse needs of the global community.

The FCC can support these efforts by supporting research and other work that help these communities to understand issues, develop solutions where needed, and deploy security technology more widely. The IAB believes that the IETF is an important partner in these efforts.

[1] "FCC Launches Inquiry To Reduce Cyber Risks," *The Internet Protocol Journal*, Volume 25, No. 1, April 2022, page 38.

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Kjetil Aas	Gareth Bryan	Dmitriy Dudko	Geert Jan de Groot	Anders Marius Jørgensen
Fabrizio Accatino	Ron Buchalski	Andrew Dul	Ólafur Guðmundsson	Merike Kaao
Michael Achola	Paul Buchanan	Joan Marc Riera	Christopher Guemez	Andrew Kaiser
Martin Adkins	Stefan Buckmann	Duocastella	Gulf Coast Shots	Christos Karayiannis
Melchior Aelmans	Caner Budakoglu	Pedro Duque	Sheryll de Guzman	Daniel Karrenberg
Christopher Affleck	Darrell Budic	Holger Durer	Rex Hale	David Kekar
Scott Aitken	BugWorks	Mark Eanes	Jason Hall	Stuart Kendrick
Jacobus Akkerhuis	Scott Burleigh	Andrew Edwards	Darow Han	Robert Kent
Antonio Cuñat Alario	Chad Burnham	Peter Robert Egli	Handy Networks LLC	Jithin Kesavan
William Allaire	Jon Harald Bøvre	George Ehlers	James Hamilton	Jubal Kessler
Nicola Altan	Olivier Cahagne	Peter Eisses	Stephen Hanna	Shan Ali Khan
Shane Amante	Antoine Camerlo	Torbjörn Eklöv	Martin Hannigan	Nabeel Khatri
Marcelo do Amaral	Tracy Camp	Y Ertur	John Hardin	Dae Young Kim
Matteo D'Ambrosio	Ignacio Soto Campos	ERNW GmbH	David Harper	William W. H. Kimandu
Selva Anandavel	Fabio Caneparo	ESdatCo	Edward Hauser	John King
Jens Andersson	Roberto Canonico	Steve Esquivel	David Hauweele	Russell Kirk
Danish Ansari	David Cardwell	Jay Etchings	Marilyn Hay	Gary Klesk
Finn Arildsen	Richard Carrara	Mikhail Evstiounin	Headcrafts SRLS	Anthony Klopp
Tim Armstrong	John Cavanaugh	Bill Fenner	Hidde van der Heide	Henry Kluge
Richard Artes	Lj Cemerias	Paul Ferguson	Johan Helsingius	Michael Kluk
Michael Aschwanden	Dave Chapman	Ricardo Ferreira	Robert Hinden	Andrew Koch
David Atkins	Stefanos Charchalakos	Kent Fichtner	Asbjørn Højmark	Ia Kochiashvili
Jac Backus	Greg Chisholm	Armin Fisslthaler	Damien Holloway	Carsten Koempe
Jaime Badua	David Chosrova	Michael Fiumano	Alain Van Hoof	Richard Koene
Bent Bagger	Marcin Cieslak	The Flirble Organisation	Edward Hotard	Alexander Kogan
Eric Baker	Lauris Cikovskis	Gary Ford	Bill Huber	Matthijs Koot
Santosh Balagopalan	Guido Coenders	Jean-Pierre Forcioli	Hagen Hultzs	Antonin Kral
William Baltas	Brad Clark	Susan Forney	Kauto Huopio	Robert Krejčí
David Bandinelli	Narelle Clark	Christopher Forsyth	Kevin Iddles	Mathias Körber
Benjamin Barkin-Wilkins	Horst Clausen	Andrew Fox	Mika Ilvesmaki	John Kristoff
Feras Batainah	Joseph Connolly	Craig Fox	Karsten Iwen	Terje Krogdahl
Michael Bazarewsky	Steve Corbató	Fausto Franceschini	David Jaffe	Bobby Krupczak
David Belson	Brian Courtney	Valerie Fronczak	Ashford Jaggernaut	Murray Kuchera
Richard Bennett	Beth and Steve Crocker	Tomislav Futivic	Thomas Jalkanen	Warren Kumari
Hidde Beumer	Dave Crocker	Laurence Gagliani	Martijn Jansen	George Kuo
Pier Paolo Biagi	Kevin Croes	Edward Gallagher	Jozef Janitor	Dirk Kurfuerst
Tyson Blanchard	John Curran	Andrew Gallo	John Jarvis	Darrell Lack
John Bigrow	André Danthine	Chris Gamboni	Dennis Jennings	Andrew Lamb
Orvar Ari Bjarnason	Morgan Davis	Xosé Bravo Garcia	Edward Jennings	Richard Lamb
Axel Boeger	Jeff Day	Osvaldo Gazzaniga	Aart Jochem	Yan Landriault
Keith Bogart	Julien Dhallenne	Kevin Gee	Nils Johansson	Edwin Lang
Mirko Bonadei	Freek Dijkstra	Greg Giessow	Brian Johnson	Sig Lange
Roberto Bonalumi	Geert Van Dijk	John Gilbert	Curtis Johnson	Markus Langenmair
Lolke Boonstra	David Dillow	Serge Van Ginderachter	Richard Johnson	Fred Langham
Julie Bottorff	Richard Dodsworth	Greg Goddard	Jim Johnston	Tracy LaQuey Parker
Photography	Ernesto Doelling	Tiago Goncalves	Jonatan Jonasson	Alex Latzko
Gerry Boudreaux	Michael Dolan	Ron Goodheart	Daniel Jones	Jose Antonio Lazaro
Leen de Braal	Eugene Doroniuk	Octavio Alfageme	Gary Jones	Lazaro
Kevin Breit	Karlheinz Dölger	Gorostiaga	Jerry Jones	Rick van Leeuwen
Thomas Bridge	Michael Dragone	Barry Greene	Michael Jones	Simon Leinen
Ilia Bromberg	Joshua Dreier	Jeffrey Greene	Amar Joshi	Robert Lewis
Václav Brožík	Lutz Drink	Richard Gregor	Javier Juan	Christian Liberale
Christophe Brun	Aaron Dudek	Martijn Groenleer	David Jump	Martin Lillepui

Roger Lindholm	Andrea Montefusco	David Raistrick	Scott Seifel	Peter Tomsu Fine Art
Link Light Networks	Fernando Montenegro	Priyan R Rajeevan	Paul Selkirk	Photography
Chris and Janet Lonvick	Joel Moore	Balaji Rajendran	Yury Shefer	Joseph Toste
Sergio Loreti	John More	Paul Rathbone	Yaron Sheffer	Rey Tucker
Eric Louie	Maurizio Moroni	William Rawlings	Doron Shikmoni	Sandro Tumini
Adam Loveless	Brian Mort	Mujtiba Raza Rizvi	Tj Shumway	Angelo Turetta
Josh Lowe	Soenke Mumm	Bill Reid	Jeffrey Sicuranza	Michael Turzanski
Guillermo a Loyola	Tariq Mustafa	Petr Rejhon	Thorsten Sideboard	Phil Tweedie
Hannes Lubich	Stuart Nadin	Robert Remenyi	Greipur Sigurdsson	Steve Ulrich
Dan Lynch	Michel Nakhla	Rodrigo Ribeiro	Fillipe Cajaiba da Silva	Unitek Engineering AG
David MacDuffie	Mazdak Rajabi Nasab	Glenn Ricart	Andrew Simmons	John Urbanek
Sanya Madan	Krishna Natarajan	Justin Richards	Pradeep Singh	Martin Urwaleck
Miroslav Madić	Naveen Nathan	Rafael Riera	Henry Sinnreich	Betsy Vanderpool
Alexis Madriz	Darryl Newman	Mark Risinger	Geoff Sisson	Surendran Vangadasalam
Carl Malamud	Thomas Nikolajsen	Fernando Robayo	John Sisson	Ramnath Vasudha
Jonathan Maldonado	Paul Nikolich	Michael Roberts	Helge Skrivervik	Philip Venables
Michael Malik	Travis Northrup	Gregory Robinson	Terry Slattery	Buddy Venne
Tarmo Mammers	Marijana Novakovic	Ron Rockrohr	Darren Sleeth	Alejandro Vennera
Yogesh Mangar	David Oates	Carlos Rodrigues	Richard Smit	Luca Ventura
John Mann	Ovidiu Obersterescu	Magnus Romedahl	Bob Smith	Scott Vermillion
Bill Manning	Tim O'Brien	Lex Van Roon	Courtney Smith	Tom Vest
Harold March	Mike O'Connor	Marshall Rose	Eric Smith	Peter Villemoes
Vincent Marchand	Mike O'Dell	Alessandra Rosi	Mark Smith	Vista Global Coaching
Normando Marcolongo	John O'Neill	David Ross	Tim Sneddon	& Consulting
Gabriel Marroquin	Jim Oplotnik	William Ross	Craig Snell	Dario Vitali
David Martin	Packet Consulting	Boudhayan	Job Snijders	Rüdiger Volk
Jim Martin	Limited	Roychowdhury	Ronald Solano	Jeffrey Wagner
Ruben Tripiana Martin	Carlos Astor Araujo	Carlos Rubio	Asit Som	Don Wahl
Timothy Martin	Palmeira	Rainer Rudigier	Ignacio Soto Campos	Michael L Wahrman
Carles Mateu	Alexis Panagopoulos	Timo Ruiters	Evandro Sousa	Laurence Walker
Juan Jose Marin Martinez	Gaurav Panwar	RustedMusic	Peter Spekrijse	Randy Watts
Ioan Maxim	Chris Parker	Babak Saberi	Thayumanavan Sridhar	Andrew Webster
David Mazel	Manuel Uruena Pascual	George Sadowsky	Paul Stancik	Tim Weil
Miles McCredie	Ricardo Patara	Scott Sandefur	Ralf Stempfner	Jd Wegner
Brian McCullough	Dipesh Patel	Sachin Sapkal	Matthew Stenberg	Westmoreland
Joe McEachern	Alex Parkinson	Arturas Satkovskis	Martin Štěpánek	Engineering Inc.
Alexander McKenzie	Craig Partridge	PS Saunders	Adrian Stevens	Rick Wesson
Jay McMaster	Dan Paynter	Richard Savoy	Clinton Stevens	Peter Whimp
Mark Mc Nicholas	Leif Eric Pedersen	John Sayer	John Streck	Russ White
Olaf Mehlberg	Rui Sao Pedro	Phil Scarr	Martin Streule	Jurrien Wijnhuizen
Carsten Melberg	Juan Pena	Gianpaolo Scassellati	David Strom	Derick Winkworth
Kevin Menezes	Chris Perkins	Elizabeth Scheid	Colin Strutt	Pindar Wong
Bart Jan Menkveld	Michael Petry	Jeroen Van Ingen	Viktor Sudakov	Makarand Yerawadekar
Sean Mentzer	Alexander Peuchert	Schenau	Edward-W. Suor	Phillip Yialeloglou
William Mills	David Phelan	Carsten Scherb	Vincent Surillo	Janko Zavernik
David Millsom	Harald Pilz	Ernest Schirmer	Terence Charles	Bernd Zeimet
Desiree Miloshevic	Derrell Piper	Benson Schliesser	Sweetser	Muhammad Ziad
Joost van der Minnen	Rob Pirnie	Philip Schneck	T2Group	Ziayuddin
Thomas Mino	Marc Vives Piza	James Schneider	Roman Tarasov	Tom Zingale
Rob Minshall	Jorge Ivan Pincay Ponce	Peter Schoo	David Theese	Jose Zumalave
Wijnand	Victoria Poncini	Dan Schrenk	Douglas Thompson	Romeo Zwart
Modderman-Lenstra	Blahoslav Popela	Richard Schultz	Kerry Thompson	廖明沂.
Mohammad Moghaddas	Andrew Potter	Timothy Schwab	Lorin J Thompson	
Roberto Montoya	Eduard Llull Pou	Roger Schwartz	Fabrizio Tivano	
Charles Monson	Tim Pozar	SeenThere		



Follow us on Twitter and Facebook

@protocoljournal



<https://www.facebook.com/newipj>

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Internet
Society



Diamond Sponsors

Your logo here!

Ruby Sponsors



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

John Crain, Senior Vice President and Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Principal – Internet Technology, Policy, and Advocacy
The Internet Society

Dr. Jun Murai, Founder, WIDE Project
Distinguished Professor, Keio University
Co-Director, Keio University Cyber Civilization Research Center, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

December 2022

Volume 25, Number 3

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
TCP and QUIC.....	2
Minimized DNS Resolution...	16
KINDNS.....	41
Supporters and Sponsors	45
Thank You!	46

Although most of the core protocols of the Internet have remained unchanged for many decades, there is still active work in the *Internet Engineering Task Force* (IETF) and elsewhere to enhance, improve, and even replace some functions provided by the protocols with respect to design and with focus on operational best practice. In this issue we will explore two areas where such work is in progress, namely transport and the *Domain Name System* (DNS).

In our first article, Geoff Huston compares the *Transmission Control Protocol* (TCP) to QUIC. QUIC has seen rapid deployment in the Internet largely due to its improved performance and extensibility, as well as privacy and security. Geoff predicts that QUIC may some day replace TCP as the major transport protocol in the Internet. The IETF has a working group dedicated to QUIC.

Another very active work area in the IETF focuses on the DNS. One topic of interest in the *DNS Operations* (DNSOP) working group is *Minimized DNS Resolution*. In our second article, Burton Kaliski, Jr. describes four different approaches to minimization: *Qname Minimization*, *NXDOMAIN Cut Processing*, *Aggressive DNSSEC Caching*, and *Local Root*. All of these approaches have been documented in RFCs and are in various stages of deployment across the Internet.

Our final article, by Adiel Akplogan, is an overview of *Knowledge-Sharing and Instantiating Norms for Domain Name System and Naming Security* (KINDNS) [pronounced “kindness”], an initiative launched by the *Internet Corporation for Assigned Names and Numbers* (ICANN) to promote DNS security and best practices.

I was very pleased to learn that the 2022 *Jonathan B. Postel Service Award* was awarded to my friend and mentor George Sadowsky for his work on the Internet Society’s *Developing Countries Workshops* in the 1990s. The tradition of training network engineers on all aspects of Internet technology continues at numerous conferences around the world to this day, most notably at events hosted by various *Network Operator Groups* (NOGs) and regional events such as APRICOT. For more details on the award, visit: <https://tinyurl.com/Postel2022>

As always, we welcome your feedback and suggestions on anything you read in this journal. Letters to the Editor may be edited for clarity and length and can be sent to ipj@protocoljournal.org

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

Comparing TCP and QUIC

by Geoff Huston, APNIC

A common view out there is that the QUIC transport protocol^[0, 4] is just another refinement to the original *Transmission Control Protocol* (TCP) transport protocol^[1, 2]. I find it hard to agree with this sentiment, because for me QUIC represents a significant shift in the set of transport capabilities available to applications in terms of communication privacy, session-control integrity, and flexibility. QUIC embodies a different communications model that makes it intrinsically useful to many more forms of application behaviours. Oh, yes. It's also faster than TCP! In my opinion it's likely that over time QUIC will replace TCP in the public Internet. So, for me QUIC is a lot more than just a few tweaks to TCP. Here we will describe both TCP and QUIC and look at the changes that QUIC has brought to the transport table.

However, we should first do a brief recap of TCP.

What Is TCP?

TCP is the embodiment of the end-to-end principle in the overall Internet architecture. All the functionality required to take a simple base of datagram delivery and impose upon this model an end-to-end signalling regime that implements reliability, sequencing, adaptive flow control, and streaming is embedded within the TCP protocol.

TCP is a *bilateral full-duplex* protocol. That means that TCP is a two-party communications protocol that supports both parties simultaneously, sending and receiving data within the context of a single TCP connection. Rather than impose a state within the network to support the connection, TCP uses synchronized state between the two end points, and much of the protocol design ensures that each local state transition is communicated to, and acknowledged by, the remote party without any mediation by the network whatsoever.

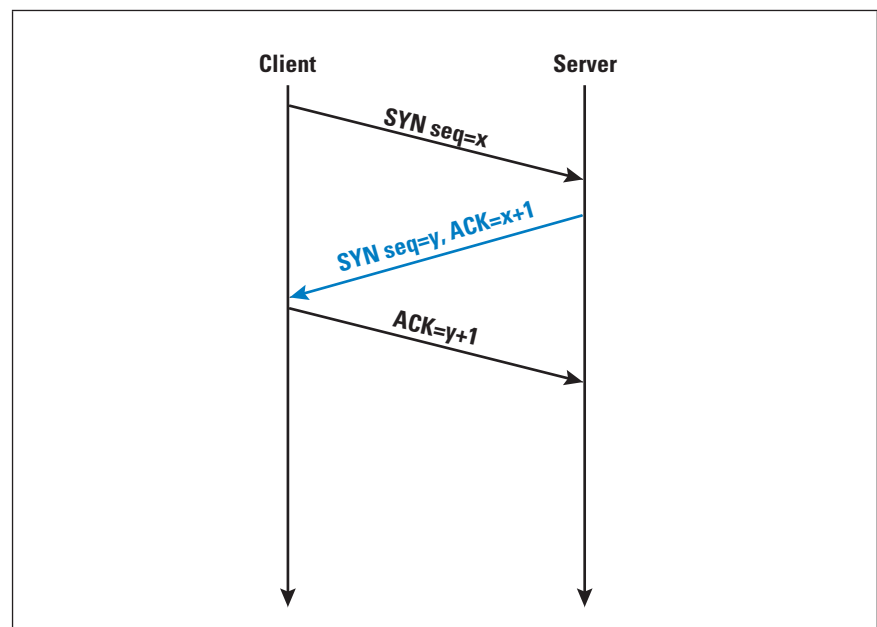
TCP is a *stream* protocol. The receiver sees the stream of data that the sender generates in precisely the same order as the sender generated. TCP is a true streaming protocol, and application-level network operations are not transparent. Other transport protocols have explicitly encapsulated each application transaction; for every sender's *write*, there must be a matching receiver's *read*. In this manner, the application-derived segmentation of the data stream into a logical record structure is preserved across the communication. TCP does not explicitly preserve such an implicit structure of the data, so there is no explicit pairing between *write* and *read* operations within the network protocol. For example, a TCP application may write three data blocks in sequence into the network connection, which the remote reader may collect in a single read operation. It is left to the application to mark the stream with its own record boundaries, if such boundaries exist in the data.

A rudimentary level of stream formatting is permitted within TCP through the concept of *urgent data* in which the sender can mark the end of a data segment that the application wants to bring to the attention of the receiver. The TCP data segment that carries the final byte of the urgent data segment can mark this data point, and the TCP receiving process has the responsibility to pass this mark to the receiving application.

The hosts at both ends identify the TCP connection by a 5-tuple of protocol identifier, source IP address, source port, destination IP address, and destination port.

The setup of a TCP connection requires a *three-way handshake*, ensuring that both sides of the connection have an unambiguous understanding of the byte-sequence values of the remote side. The operation of the connection setup is as follows: The local system sends an initial sequence number to the remote-end port using a SYN packet. The remote system responds with an *acknowledgement* (ACK) of the initial sequence number and the remote end's initial sequence number in a response SYN packet. The local end responds with an ACK of this remote sequence number. These handshake packets are conventionally TCP packets without any data payload. At this point, TCP shifts into a reliable data flow-control mode of operation. (Figure 1)

Figure 1: TCP 3-way Handshake



TCP is a *sliding window* protocol. The data stream is a sequence of numbered bytes. The sender retains a copy of all sent but as yet unacknowledged data in a local send buffer. When a receiver receives a data segment whose starting sequence number is the next expected data segment, it will send an ACK back to the sender with the sequence number of the end of the received data segment. This process allows the sender to discard all data whose sequence number is less than this received ACK sequence in the local send buffer and advance the send window.

When the received data is out of order, it will send an ACK back to the sender with the sequence number of the last in-order received data. In addition, the ACK message includes the size of the receiver's available buffer size (*receive window*). The volume of unacknowledged data must be no larger than this receiver window size. The overall constraint is that at all points in time the sender should ensure that the volume of unacknowledged data in flight in the network is the smaller of the advertised receive window size and the total capacity of the local send buffer.

TCP is an *ACK-clocked* flow-control protocol, in that within a static lossless mode of operation each received ACK packet indicates that a certain volume of data has been received at the receiver end (and hence has been removed from the network), and this clocking is accompanied by an advertised receive window that then permits the sender to inject the same volume of data into the network as the receiver has removed. Hence, the sending rate is governed by the received ACK rate.

However, TCP is not necessarily aware of the available path capacity of the network, and it must implement a control algorithm at the sending end that attempts to establish a dynamic equilibrium between the flow volume of the TCP session and all other concurrent TCP sessions that have path segments in common with this session. The mode of operation of this flow control is not fixed in the TCP specification, and numerous flow-control algorithms are in use. Many of these control algorithms use an induced instability in TCP through an approach of slow inflation of the sending window for each received ACK, and a rapid drop of the sending window in response to an indication of packet drop (3 duplicate ACKs). This process of sending rate inflation will stop when the send buffer is full, indicating that the sender cannot store any more sent data and must await ACKs before sending more data (send buffer rate limited). It will also stop sending rate inflation when the network cannot accept any further data in flight as the buffers of the network are already full, so further sent data will cause packet loss, which will be signalled back to the sender by duplicate ACKs.

This process has many outcomes relevant to service quality. First, TCP behaves adaptively rather than predictively. The flow-control algorithms are intended to increase the data-flow rate to fill all available network path capacity but also quickly back off if the available capacity changes because of network congestion or if a dynamic change occurs in the end-to-end network path that reduces this available capacity. Second, a single TCP flow across an otherwise idle network attempts to fill the network path with data, optimizing the flow rate (as long as the send buffer is larger than the network flow capacity). If a second TCP flow opens up across the same path, the two flow-control algorithms will interact so that both flows will stabilize to use approximately half of the available capacity per flow. More generally, TCP attempts to behave fairly, in that when multiple TCP flows are present the TCP algorithm is intended to share the network resource evenly across all active flows.

A design tension always exists between the efficiency of network use and enforcing predictable session performance. With TCP, you do not necessarily have predictable throughput but gain a highly utilized and efficient network.

TCP and TLS

Transport Layer Security (TLS)^[3] is handled as a further layer of indirection. When the TCP 3-way handshake is complete, the parties enter a TLS negotiation phase to allow authentication of the remote end of the connection, and to establish a session key that is used to manage the encryption of the session data.

TLS commences with an exchange of credentials. In version 1.3 of TLS (the latest version of this protocol), the client sends a *client hello* message that includes the TLS version the client supports, the cipher suites supported, the name of the service, and a string of random bytes known as the *client random*. The server responds with a *server hello* message that contains the public key certificate of the server, the *server random*, the chosen cipher suite, and a digital signature of the hello messages. Both ends now know each other's random values and the chosen cipher suite, so both can generate a master secret for session encryption. The client sends a *finished* message to indicate that the secure symmetric session key is ready for use (Figure 2).

Earlier versions of TLS used additional packets in the hello exchange that increased the time to complete the TLS handshake.

QUIC

We can now move on to QUIC. The QUIC transport protocol^[4] was apparently designed to address several issues with TCP and TLS, and in particular to improve the transport performance for encrypted traffic with faster session setup, and to allow for further evolution of transport mechanisms and explicitly avoid the emerging TCP ossification within the network.

It is a grossly inaccurate simplification, but at its simplest level QUIC is simply TCP encapsulated and encrypted in a *User Datagram Protocol* (UDP) payload. To the external network QUIC has the appearance of a bidirectional UDP packet sequence where the UDP payload is concealed. To the endpoints you can use QUIC as a reliable full-duplex data-flow protocol. Even at this level, QUIC has numerous advantages over TCP. The first lies in the deliberate concealing of the transport flow-control parameters from the network. The practice of deploying network middleware that rewrites TCP flow-control values to impair the behaviour of the application has not enjoyed widespread support from the application layer, and hiding these flow-control parameters from the network certainly prevents this practice. Second, it can allow the shift of responsibility for providing the transport protocol from the platform to the application. A tension between the application and the platform is longstanding.

Changes to kernel-level TCP are performed via updates to the platform software, and often applications have to wait for the platform to make changes before they can take advantage of the change. For example, if an application wanted to use the TCP *Bottleneck Bandwidth and Round-trip Propagation Time* (BBR) flow-control algorithm, then it would need to wait for a platform to integrate an implementation of the algorithm. By using a basic UDP interface to the transport services of the platform, you can lift the entire flow-control and encryption service into the application itself, if so desired. You may experience some performance penalty of shifting the transport code from the kernel to user space, but in return the application regains complete control of the transport service and allows it to operate in a mode that is not only independent of the platform, but also hidden from the platform. This shift gives the application environment greater levels of control and agility.

However, QUIC does a lot more than just wrapping up TCP in UDP, so let's look at the QUIC protocol in a little more detail.

QUIC Connections

A QUIC *connection* is a shared state between a single client and a single server. QUIC uses the combination of two numbers, one selected by each end, to form a pair of connection IDs. This pair of IDs acts as a persistent identity for the QUIC session, which is used to ensure that changes in addressing at lower protocol layers (addresses or ports) will not cause delivery of packets to a wrong recipient on the end host.

The primary function of a connection ID is to ensure that changes in addressing at lower protocol layers (IP source address and UDP source port numbers) do not cause packets for a QUIC connection to be dropped when the external IP address of an endpoint changes. Each endpoint selects a connection ID using an implementation-specific (and perhaps deployment-specific) method that allows identification of received packets with that connection ID by the endpoint upon receipt to the appropriate QUIC connection instance.

After an endpoint receives a packet with the same connection ID and a different IP address or UDP port, it will verify the peer's ownership of the new address by sending a *challenge frame* containing random data to this new address and waiting for an echoed response with the same data. This challenge and response exchange is performed within the established crypto state, so it is intentionally challenging for an eavesdropper to hijack a session in this way. The two endpoints can continue to exchange data after the verification of the new address.

This verification is particularly useful in terms of negotiating various forms of *Network Address Translation* (NAT) behaviour. NATs are intentionally transport-aware and for TCP, NATs will attempt to maintain a translation state until it observes the closing FIN protocol exchange. UDP offers no such externally visible clues as to the ending of a session, and NATs are prone to interpreting a silent period as a signal to tear down the NAT state.

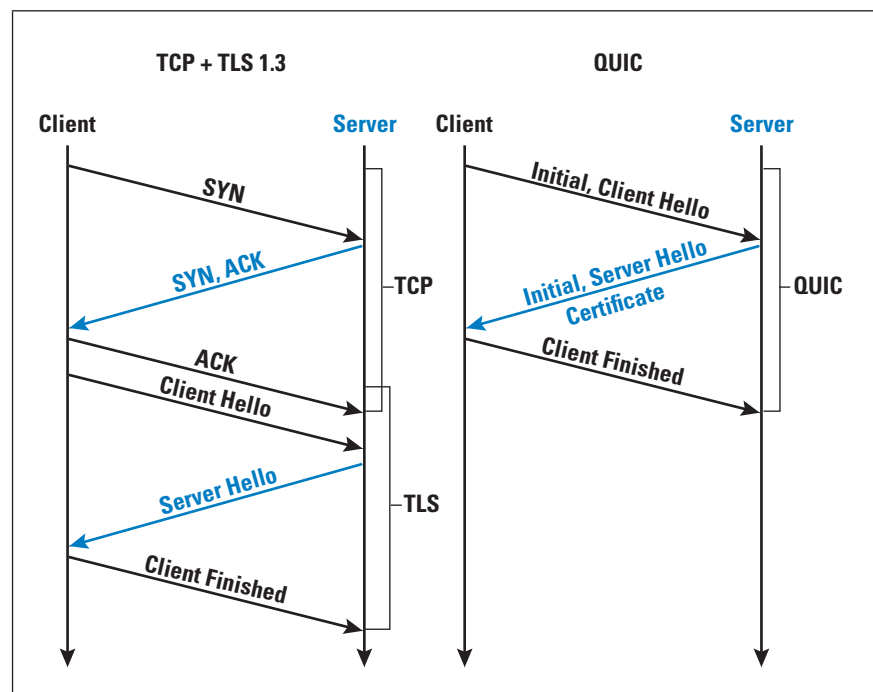
In such a case the next outbound packet might be assigned a new source address and/or UDP source port number by the NAT. It is also useful in terms of session resumption where the connection may have been idle for an extended period, and the NAT binding may have timed out. With TCP, any change in any of the four address and port fields of the connection 5-tuple will cause rejection of the packet as part of the TCP session. QUIC's use of a persistent connection ID permits the receiver to associate the new sender's address details with an existing connection.

You also can use this QUIC functionality of address agility in the context of network-level changes, such as a device switching between WiFi and cellular services while maintaining an active QUIC transport session.

QUIC Connection Handshake

A QUIC connection starts with a handshake that establishes a shared communications state and a shared secret using the QUIC-TLS protocol cryptographic handshake protocol in a single exchange. This protocol merges the TCP 3-way handshake and the TLS 1.3 3-way handshake into a single 3-packet exchange (Figure 2). This merge eliminates a full *Round Trip Time* (RTT) in the QUIC startup phase, which for short sessions is a very significant improvement in session performance.

Figure 2: TCP/TLS and QUIC Handshakes



QUIC also allows a client to send 0-RTT encrypted application data in its first packet to the server by reusing the negotiated parameters from a previous connection and a TLS 1.3 *Pre-Shared Key* (PSK) identity issued by the server, although these 0-RTT data exchanges are not protected against replay attack.

Packets and Frames

The QUIC protocol sends *packets* along the connection. Packets are individually numbered in a 62-bit number space. There is no allowance for retransmission of a numbered packet. If data is to be retransmitted, it is done in a new packet with the next packet number in sequence. That way there is a clear distinction between the reception of an original packet and a retransmission of the data payload.

You can load multiple QUIC packets into a single UDP datagram. QUIC UDP datagrams must not be fragmented, and unless the end performs *Path Maximum Transmission Unit* (PMTU) discovery, QUIC assumes that the path can support a 1,200-byte UDP payload.

A QUIC client expands the payload of all UDP datagrams carrying Initial packets to at least the smallest allowed maximum datagram size of 1,200 bytes by adding padding frames to the Initial packet or by coalescing a set of Initial packets. The payload of all UDP datagrams carrying ACK-eliciting Initial packets is padded to at least the smallest allowed maximum datagram size of 1,200 bytes. Sending UDP datagrams of this size ensures that the network path supports a reasonable PMTU in both directions. Additionally, a client that expands Initial packets helps reduce the order of amplitude gain of amplification attacks caused by server responses toward an unverified client address.

QUIC packets are encrypted individually so that the decryption process does not result in data decryption waiting for partially delivered packets. This encryption is not generally possible under TCP, where the encryption records are in a byte stream and the protocol stack is unaware of higher-layer boundaries within this stream. The additional inference from this per-packet encryption is that it's a requirement that QUIC IP packets are not fragmented. QUIC implementations typically use a conservative choice in the maximum packet size so that IP packet fragmentation does not occur.

A QUIC receiver ACKs the highest packet number received so far, together with a listing of all received contiguous packet number blocks of lower-numbered packets if there are gaps in the received packet sequence. Because QUIC uses purpose-defined ACK frames, QUIC can code up to 256 such number ranges in a single frame, whereas TCP *Selective Acknowledgment* (SACK) has a limit of 3 such sequence number ranges. This limit allows QUIC to provide a more detailed view of packet loss and reordering, leading to higher resiliency against packet losses and more efficient recovery. Lost packets are not retransmitted. Data recovery is performed in the context of each QUIC stream.

QUIC Streams

A QUIC connection is further broken into *streams*. Each QUIC stream provides an ordered byte-stream abstraction to an application similar in nature to a TCP byte stream. QUIC allows for an arbitrary number of concurrent streams to operate over a connection. Applications may indicate the relative priority of streams.

Because the connection has already performed the end-to-end association and established the encryption context, you can establish streams with minimal overhead. A single stream frame can open, pass data, and close down within a single packet, or it can exist for the entire lifetime of the connection.

By comparison, it is possible to multiplex a TCP session into streams, but all such multiplexed TCP streams share a single flow-control state. If the TCP receiver advertises a zero-sized window to the sender, then all multiplexed streams are blocked in a TCP scenario.

Each QUIC stream is identified by a unique *stream ID*, where its two least significant bits are used to identify which endpoint initiated the stream and whether the stream is bidirectional or unidirectional. The byte stream is segmented to data frames, and the stream frame offset is equivalent to the TCP sequence number, used for data-frame delivery ordering and loss detection and retransmission for reliable data delivery.

QUIC endpoints can decide how to allocate bandwidth between different streams, and how to prioritize transmission of different stream frames based on information from the application. This feature ensures effective loss recovery, congestion control, and flow-control operations, which can significantly impact application performance.

QUIC Datagrams

In addition to reliable streams, QUIC also supports an unreliable but secured data-delivery service with DATAGRAM frames, which will not be retransmitted upon loss detection^[5]. When an application sends a datagram over a QUIC connection, QUIC will generate a DATAGRAM frame and send it in the first available packet. When a QUIC endpoint receives a valid DATAGRAM frame, it is expected that it would deliver the data to the application immediately. These DATAGRAM frames are not associated with any stream.

If a received packet contains only DATAGRAM frames, then the ACK frame can be delayed, as the sender will not retransmit a frame when there is an ACK failure in any case. This service is not a reliable datagram service. If a sender detects that a packet containing a specific DATAGRAM frame might have been lost, the implementation may notify the application that it believes the datagram was lost. Similarly, if a packet containing a DATAGRAM frame is acknowledged, the implementation may notify the sender application that the datagram was successfully transmitted and received.

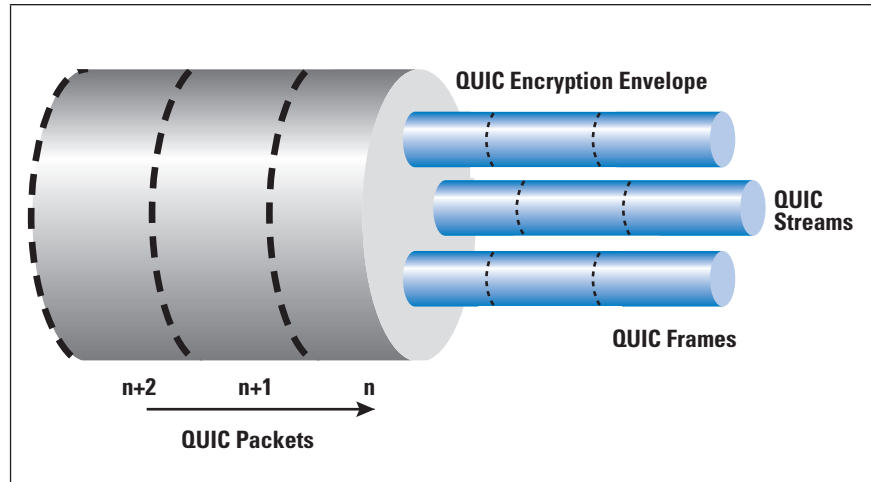
QUIC Frames

Each packet contains a sequence of *frames*. Frames have a frame-type field and type-dependant data. The QUIC standard^[4] defines 20 different frame types. They serve an analogous purpose to the TCP *flags*, carrying a control signal about the state of streams and the connection itself.

Frame types include padding, ping (or keepalive), ACK frames for received packet numbers, which themselves contain *Explicit Congestion Notification* (ECN) counts as well as ACK ranges, as well as stream data frames and datagram frames.

Figure 3 shows the larger organisation of QUIC connections, streams, and frames.

Figure 3: QUIC Logical Organisation



QUIC Recovery and Flow Control

QUIC packets contain one or more frames. QUIC performs loss detection based on these packets, not on individual frames. For each packet that the receiver acknowledges, all frames carried in that packet are considered received. The packet is considered lost if that packet is unacknowledged when a later sent packet has been acknowledged, and when a loss threshold is met.

QUIC uses two thresholds to determine loss. The first is a *Packet Reordering Threshold* t . When packet x is acknowledged, then all unacknowledged packets with a number less than $x - t$ are considered lost. The second is related to the QUIC-measured RTT interval, the *waiting time* w which is determined as a weight factor applied to the current estimated RTT interval. If the time of the most recent acknowledgement is t , then all unacknowledged packets sent before time $t - w$ will be considered lost.

For recovery, all frames in lost packets where the associated stream requires retransmission will be placed into new packets for retransmission. The lost packet itself is not retransmitted.

As with TCP's advertised receiver window, QUIC contains a mechanism to enable a QUIC receiver to control the maximum amount of data that a sender can send on an individual stream, and the maximum amount on all streams at any time. Also, as with TCP, QUIC does not specify the flow-control algorithm to be used by reliable streams, although one such sender-side congestion controller is defined in [6]. This algorithm is similar to TCP's *New Reno*^[7].

We now examine some problems with QUIC.

RPC Support

IP hosts commonly support just two transport services, UDP and TCP. UDP is a simple datagram-delivery service. Data encapsulated using UDP have no assured delivery. TCP, as we have seen, is a reliable streaming service. The TCP protocol repairs any packet loss or changes to the delivered packet sequence.

Another model, namely the *Remote Procedure Call* (RPC) model, emulates the functionality of procedure calls, and rather than the byte-stream model of TCP or the datagram model of UDP, the RPC model is a reliable request/reply model, where the reply is uniquely associated with the request. Perhaps the most well-known example today of an RPC framework is *gRPC*^[8]. *gRPC* is based on an HTTP/2 platform, implying that the framework is susceptible to head-of-line blocking as with any other TCP-based substrate.

The issue here is that a reliable byte stream is not the right abstraction for RPC, as the core of RPC is a request/reply paradigm, which is more aligned to a reliable messaging paradigm, with all that such a paradigm entails. A capable RPC framework needs to handle lost, mis-ordered, and duplicated messages, with an identifier space that can match requests and responses. The underlying message transport needs to handle messages of arbitrary size, which entails packetization adaptation within the transport.

The bidirectional stream framework is a reasonable match to the RPC communications model where each RPC can be matched against an individual stream. The stream is reliable and sequenced. The data framing is not contained in QUIC, and it is still an application task to add a record structure to an RPC stream, if that is what is required. The invocation overhead is low in that the encrypted end-to-end connection is already established.

It certainly appears that HTTPS behaves much more like RPC than a reliable byte stream. That can benefit applications that run over HTTP(S), such as *gRPC*, and a set of *Representational State Transfer* (REST)ful APIs.

Load-Balancing QUIC

In today's world of managing scale, it is very common to place a front-end load balancer across many servers. The load balancer in the TCP world typically categorizes packets as being in the same TCP session because of a common 5-tuple value of protocol, IP addresses, and port numbers, with the confident assurance that this value is stable for the life of the TCP session.

QUIC offers no such assurances. The 5-tuple load-balancing approach can work, but if the client is behind a NAT that performs what could be called "aggressive" rebinding, then any such load-balancing approach will be thrown. The reason why is that UDP does not provide session signalling to a NAT, so there is no a priori assurance that the NAT bindings (and the presented source address and port) will remain constant for the entire QUIC session.

Now in theory IPv6 could invoke the *Flow-ID* to provide a proxy persistent field that remains constant for a flow, but the Flow-ID is of limited size and has no assurances of uniqueness, as well as evidence of highly variable treatment by IPv6 network infrastructure and end hosts.

This topic touches upon a major assumption in today's high-capacity server infrastructure on the public Internet. Data streams use TCP and the DNS uses UDP. Using UDP to carry sustained high-volume streams may not match the internal optimisations used in server content-delivery networks.

DDoS Defence

The next issue here is exposure to *Distributed Denial-of-Service* (DDoS) attacks. An attacker can send a large volume of packets to the server and cause the server to perform work to attempt to decrypt the packet. For this capability to be successful in TLS over TCP the attacker must make a reasonable guess of the TCP sequence number and window size for the packet to be accepted and passed to the TLS decoder. QUIC has no lightweight packet filter before the decoder is invoked.

On the other hand, the session encryption uses symmetric crypto algorithms, which are less of a load on the receiver to decode than asymmetric encryption. Is this difference enough to allow large-scale QUIC platforms that are DDoS resistant to be constructed? I'm unsure if there are clear answers here, but it seems that it's part of the cost of having a more complete encryption framework, which in itself appears to be sorely needed on the public Internet.

Private QUIC

For private contexts, can QUIC negotiate a “null” TLS encryption algorithm? There is a bigger world out there beyond the public Internet, and in many private data centre environments the overheads of encrypting and decrypting packets may appear to be unnecessary. While QUIC can present some clear advantages in terms of suitability to complex application behaviours in the data centre that can leverage QUIC's multi-stream capability, the cost of encryption may be too high.

Of course, there is nothing stopping an implementation using a null encryption algorithm, but such an implementation could talk only to other implementations of itself. Strictly speaking, if you remove encryption, then it's no longer QUIC and it won't interoperate with anything else that is QUIC.

QUIC and OpenSSL

It is useful to ask that if QUIC has such clear advantages over TCP, then why hasn't the adoption of QUIC been rapid? Metrics of QUIC use tend to point to a use rate of some 30% of web sessions (such as Cloudflare's Radar report^[9]).

However, if you alter the measurement to measure the extent to which browsers on end systems are capable of supporting a QUIC session, then the measurement jumps to 60%^[10].

There are a couple of reasons why QUIC use is far lower than QUIC capability. The first is that the Chrome browser still relies on the content-level switch to QUIC, so the client has to visit the site for the first time using HTTP/2 (TCP/TLS) and thereby receive an indication if the server can support QUIC, and then on the second visit the client may use QUIC. It's not quite as simple as this, as HTTP/2 uses persistent connection, so if the second visit is sufficiently close in time to the first, then the HTTP/2 session will remain open and still be used. The Safari browser is capable of using QUIC on first use because it is triggered by the *Service Binding* (SVCB) record in the DNS, but the market share of Safari is relatively small in comparison to QUIC.

The second reason lies in the web server environment. Many servers rely on the *OpenSSL* TLS library^[11], and so far, (November 2022) *OpenSSL* does not include support for QUIC. QUIC is supported in *BoringSSL*^[12], but as the notes for *BoringSSL* state, *BoringSSL* is a fork of *OpenSSL* that is designed to meet Google's needs, and while it works for Google, it may not work for everyone else. Google does not recommend that third parties depend on *BoringSSL*. There is also *QuicTLS*, a fork of *OpenSSL* that Akamai and Google support^[13]. This fragmentation of *OpenSSL* is not exactly helpful, and the result is that many server environments are waiting for *OpenSSL* to incorporate a QUIC library. This effort was delayed by the work on *OpenSSL* release 3.0.0, and then the *OpenSSL* folks announced their intention to provide a fully functional QUIC implementation, and this development of a new QUIC protocol stack may further delay QUIC support in *OpenSSL* by months, if not years. This impediment may well be the major one behind the very large-scale deployment of QUIC in the guise of HTTP/3 across the Internet.

Conclusions

We can draw a few conclusions from this effort with QUIC:

Any useful public communications medium needs to safeguard the privacy and integrity of the communications that it carries. The time when open protocols represented an acceptable compromise between efficiency, speed, and privacy are over, and these days all network transactions in the public Internet need to be protected by adequate encryption. The QUIC model of wrapping a set of transactions between a client and a server in a single encryption state represents a sensible design decision.

Encryption is no longer an expensive luxury, but a required component for all transactions over the public Internet. The added imposition is that adding encryption into a network transaction should impose no additional performance penalty in terms of speed and responsiveness.

Network transactions come in many forms, and TCP and UDP tend to represent two ends of a relatively broad spectrum. UDP is just too susceptible to abuse, so we've heaped everything onto TCP. The problem is TCP was designed as an efficient single streaming protocol, and retro-fitting multiple sessions, short transactions, shared congestion state, and shared encryption state have proved to be extremely challenging.

Applications are now dominant in the Internet ecosystem, while platforms and networks are being commoditised. We are seeing users losing patience with platforms that provide common transport services for the application that they host, and a new model where the application comes with its own transport service. This model is not just the HTTP client/server model; it has been extended into application-specific *Domain Name System* (DNS) name resolution with DNS over HTTPS. It's highly likely that this trend will continue for the moment.

Taking an even broader perspective, the context of the Internet's success lies in shifting the responsibility for providing service from the network to the end system. This shift allowed us to make more efficient use of the common network substrate and push the cost of this packetization of network transactions over to end systems. It shifted the innovation role from the large and lumbering telco operators into the more nimble world of platform software. The success of Microsoft with its Windows product was not an accident by any means. QUIC takes this success one step further, and pushes the innovation role from platforms to applications, just at the time when platforms are declining in relative importance within the ecosystem. From such a perspective, the emergence of an application-centric transport model that provides faster services and a larger repertoire of transport models, the encompassing of comprehensive encryption was an inevitable development.

References and Further Reading

- [0] Geoff Huston, "A Quick Look at QUIC," *The Internet Protocol Journal*, Volume 22, No. 1, March 2019.
- [1] Wesley Eddy, Ed., "Transmission Control Protocol (TCP)," RFC 9293, August 2022.
- [2] Jon Postel, Ed., "Transmission Control Protocol," RFC 793, September 1981.
- [3] Eric Rescorla, "The Transport Layer Security (TLS) Protocol Version 1.3," RFC 8446, August 2018.
- [4] Jana Iyengar, Ed., and Martin Thomson, Ed., "QUIC: A UDP-Based Multiplexed and Secure Transport," RFC 9000, May 2021.
- [5] Tommy Pauly, Eric Kinnear, and David Schinazi, "An Unreliable Datagram Extension to QUIC," RFC 9221, March 2022.
- [6] Jana Iyengar, Ed., and Ian Swett, Ed., "QUIC Loss Detection and Congestion Control," RFC 9002, May 2021.

- [7] Tom Henderson, Sally Floyd, Andrei Gurtov, and Yoshifumi Nishida, “The NewReno Modification to TCP’s Fast Recovery Algorithm,” RFC 6582, April 2012.
- [8] gPRC – A cross-platform open source RPC framework:
<https://grpc.io/>
- [9] Cloudflare Radar: <https://radar.cloudflare.com/>
- [10] APNIC Labs, QUIC Usage Report:
<https://stsats.labs.apnic.net/quic>
- [11] OpenSSL a library for secure communication:
<https://openssl.com>
- [12] BoringSSL, an open source fork of the OpenSSL library operated by Google for internal use:
<https://boringssl.googlesource.com/boringssl>
- [13] QuicTLS, an open source fork of the OpenSSL library developed by Akamai and Microsoft as an interim Quic API:
<https://github.com/quictls/openssl>

GEOFF HUSTON, B.Sc., M.Sc. A.M., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: guh@apnic.net

Check your Subscription Details!

If you have a print subscription to this journal, you will find an expiration date printed on the back cover. For several years, we have “auto-renewed” your subscription, but now we ask you to log in to our subscription system and perform this simple task yourself. Make sure that both your postal and e-mail addresses are up-to-date since these are the only methods by which we can contact you. If you see the words “Invalid E-mail” on your copy this means that we have been unable to contact you through the e-mail address on file. If this is the case, please contact us at ipj@protocoljournal.org with your new information. The subscription portal is located here:
<https://www.ipjsubscription.org/>

Minimized DNS Resolution: Into the Penumbra

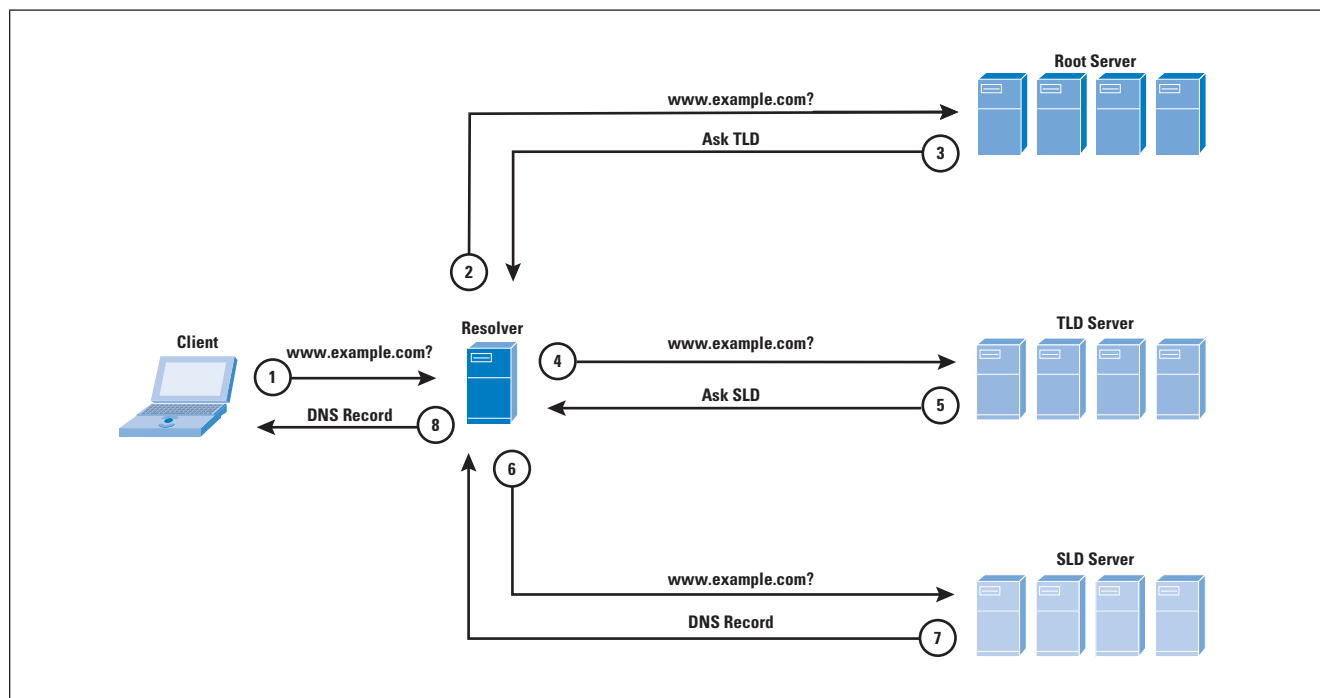
by Burton S. Kaliski Jr., Verisign

The *Domain Name System* (DNS) has long followed a traditional approach of answering queries, where resolvers send a query with the same fully qualified domain name to each name server in a chain of referrals, and, generally, apply the final answer they receive only to the domain name that was queried for. Motivated by interest in reducing both the quantity and sensitivity of information exchanged between DNS ecosystem components, DNS operators are now starting to deploy various minimization techniques that either put less information into queries or take more information out of answers, thereby reducing the need for additional queries. This article reviews four minimization techniques documented by the *Internet Engineering Task Force* (IETF), reports on their implementation status, and discusses the effects of their adoption on DNS measurement research.

DNS resolution begins with the usual occurrence that happens millions of times a second around the world: a client sends a DNS recursive resolver a query like “What is **www.example.com**’s *Internet Protocol* (IP) address?” The resolver answers, “**www.example.com**’s IP address is **93.184.216.34**.”

Many clients may use the same resolver, so the resolver may already have a response to the query in its cache. If the resolver has an empty cache, it will interact with the authoritative name-server system using a protocol flow such as shown in Figure 1.

Figure 1:Textbook DNS Resolution



1. The client asks the resolver, “What is **www.example.com**’s IP address?”
2. The resolver queries one of the DNS’s 13 root servers^[1] for an answer to question 1.
3. The root server responds with a referral-type response directing the resolver to the name server for the *Top-Level Domain* (TLD) in the query name, that is, the **.com** name server.
4. The resolver sends the query to the TLD server.
5. The TLD server refers the resolver to the name server for the *Second-Level Domain* (SLD), that is, the **example.com** name server.
6. The resolver sends the query to the SLD server.
7. The SLD server returns one or more DNS records that specify **www.example.com**’s IP address.
8. The resolver relays the DNS records to the client.

The referrals in steps 3 and 5 are a result of the delegation structure of DNS. The root zone has delegated the authority for responding to queries for domain names within existing TLDs to TLD servers. Many TLD zones have similarly delegated the authority for responding to queries within SLDs to SLD servers. In step 7, the SLD server has the authority to respond for the domain name **www.example.com**.

The DNS standard (based on RFC 1035^[2] and other documents)—as well as current practice—include many more details. For purposes of this article, the “textbook DNS” described here is an effective starting point, but two additional details may be helpful in framing the techniques that follow:

- If a name server knows that a domain name doesn’t exist, then it returns the negative response code (rcode 3), typically referred to as NXDOMAIN. (Otherwise, either the domain name exists and the name server is authoritative for it and returns a positive answer along with rcode 0; or the name server is not authoritative and returns a *referral*.)
- If a resolver and a name server implement the *Domain Name System Security Extensions* (DNSSEC)^[3], the resolver asks the name server to include DNSSEC information in its response, and if the domain name doesn’t exist, then the name server also returns an NSEC^[4] or NSEC3^[5] record that specifies two endpoints between which no other domain names exist, for some ordering of domain names. With NSEC, the ordering is based on the domain names themselves; with NSEC3, it’s based on their hash values. Either way, the resolver receives information demonstrating not only that the queried name doesn’t exist, but also that *other* domain names between the endpoints don’t exist.

(The records are formed this way so that they can be precomputed and signed when the name server is provisioned, based on domain names that do exist in a zone. The name server then already has the information it needs to respond to a query for a nonexistent domain name, without having to sign responses in real time, although some name servers do support dynamic signing.)

It's clear from a brief review of Figure 1 that textbook DNS resolution includes more information in DNS exchanges than necessary. This fact is particularly evident on the resolver-to-root exchange, where the resolver queries for a fully qualified domain name, yet the root server responds with a referral involving just the TLD. But the observation holds at other levels as well.

Forwarding fully qualified domain names may have historically simplified implementation, in that the resolver either gets the answer to a query from its cache, or it forwards the same query to a succession of name servers. This practice also minimizes the depth of the iterative resolution process, because the query includes enough information for each name server either to refer the resolver to another name server, or to answer the query itself (if the query wasn't fully qualified, then a name server might respond with a referral to itself in some cases, an unnecessary extra step). However, the textbook approach doesn't leverage all information available to the resolver, either from DNS or from other sources. Indeed, a fully qualified domain name, while convenient from an implementation perspective, may include more information than the name server needs to know.^[6]

Minimized DNS Resolution

Minimized DNS resolution encompasses an emerging set of techniques that bring the resolver-to-authoritative traffic closer to the need-to-know principle, while still facilitating DNS resolution. Four such techniques have received the most attention, each reducing the quantity and/or sensitivity of information exchanged between resolvers and authoritative name servers in a different way. Documented by the IETF's *DNS Operations* (DNSOP) working group, the techniques include:

- Query Name (or qname) Minimization, described in RFC 9156^[7];
- NXDOMAIN Cut Processing, described in RFC 8020^[8];
- Aggressive DNSSEC Caching, described in RFC 8198^[9]; and
- Local Root (sometimes called “hyperlocal”) and other locally served zones, described (in the case of the root zone) in RFC 8806^[10].

Important from an operational perspective, all four can generally be applied by a resolver on its own, without any coordinated changes by authoritative name servers, other than the participating name server conforming with previous DNS specifications. (The locally served zones technique requires that the zone data be made available.)

RFC 8932, produced by the IETF's *DNS Private Exchange* (DPRIVE) working group, encourages implementation of all four techniques to reduce both the quantity and sensitivity of “data sent onwards from the [recursive resolver] server”^[11]. (DPRIVE and other IETF working groups have also developed specifications for DNS encryption, but they are outside the scope of this article.)

The techniques can generally be adopted for interactions between resolvers and authoritative name servers for any zone. (They don't apply to the client-resolver exchange.) They are particularly beneficial for interactions with the root and TLD servers, for at least two reasons:

1. The primary purpose of the root and TLD servers is global navigational availability: referring requesters to other name servers that are actually authoritative for a response. A fully qualified domain name (or even a full set of queries) is therefore not generally needed at these servers, only enough information to make the referral, making minimization techniques appropriate options. But high-availability service is paramount, favoring techniques with low operational risk.
2. Because of the recursive, cached architecture of DNS, the sensitivity of the traffic on these exchanges is already relatively low compared to other parts of the DNS ecosystem, such as the client-to-resolver exchange. In particular, because the resolver is between the client and the authoritative name servers, its queries to the authoritative name server conceal the client's identity and instead represent aggregate interests of clients. (Moreover, although information about the client's IP address may be conveyed in a query via the “client subnet” option^[12], it is specifically recommended that this extension not be included in queries to the root and TLD servers.) Minimization techniques can therefore arguably lower the sensitivity of the information on the resolver-to-root and -TLD exchanges sufficiently that techniques with higher operational risk such as DNS encryption become questionable from a cost-benefit perspective, compared to disclosure risks on other exchanges such as client-to-resolver^[13].

Minimization techniques also can improve resolver performance, given that they enable a resolver to answer more queries on its own, and thereby respond more quickly. They can likewise improve performance for name servers, which will receive less unnecessary traffic—including attack traffic that might have leveraged a resolver as an intermediary. And as minimized traffic becomes the “new normal” on these exchanges, it may become easier for name servers to detect and deflect other types of attack traffic, which will become more “abnormal.”

Even if a resolver implements DNS encryption, it still makes sense for the resolver to implement minimization techniques to reduce the amount of information disclosed to name-server operators.

Minimization opens a new chapter in DNS resolution. With the new techniques, the traditional DNS resolution process is updated with a new approach optimized for the global DNS as it exists today, balancing confidentiality and availability objectives. The first minimization technique is perhaps the most fundamental, as it changes the most apparent nonminimized feature of textbook DNS: sending the fully qualified domain name to each name server in the chain of referrals. Note: In 2015, Verisign announced a royalty-free license to its qname minimization patents in connection with certain IETF standardization efforts. For more information, refer to IETF IPR disclosure 5197.

Query Name (Qname) Minimization

It is just a “tradition” that resolvers send the fully qualified domain name at each level of the DNS hierarchy, not a requirement of the DNS specifications. In the words of RFC 9156 (first reported by Stéphane Bortzmeyer in RFC 7816^[14]), the tradition is motivated by an early goal of minimizing the number of queries that might need to be made:

In a conversation with the author in January 2015, Paul Mockapetris explained that this tradition comes from a desire to optimise the number of requests, when the same name server is authoritative for many zones in a given name (something that was more common in the old days, where the same name servers served `.com` and the root) or when the same name server is both recursive and authoritative (something that is strongly discouraged now).

This practice, as discussed previously, can also optimize the number of requests when a name server is authoritative for only one zone.

The consequence of the tradition is that the resolver included more information than necessary in each query. Although the risk of disclosure of sensitive information on the resolver-to-root and -TLD exchanges is relatively low, as discussed previously, it would be better, per the principle of minimum disclosure, to send only as many labels as the name server needs to make a referral. Any labels beyond that point are extraneous information.

One way to reduce the amount of information disclosed is to remove one or more of the extraneous labels. In this “omitted-label” approach to reducing the amount of information included in a query to an authoritative name server, the query name `www.example.com` in the request to the root server could be replaced simply with the TLD, that is, with `.com`.

Another way is to replace one or more of the extraneous labels with random or other alternative labels. As examples of a “false-label” approach, we could replace `www.example.com` with `<r3>.<r2>.com` or with `<r2>.com`, where `<r2>` and `<r3>` are randomly generated labels. Another real-world qname minimization technique suggested replaces `www.example.com` with `_example.com`.^[7]

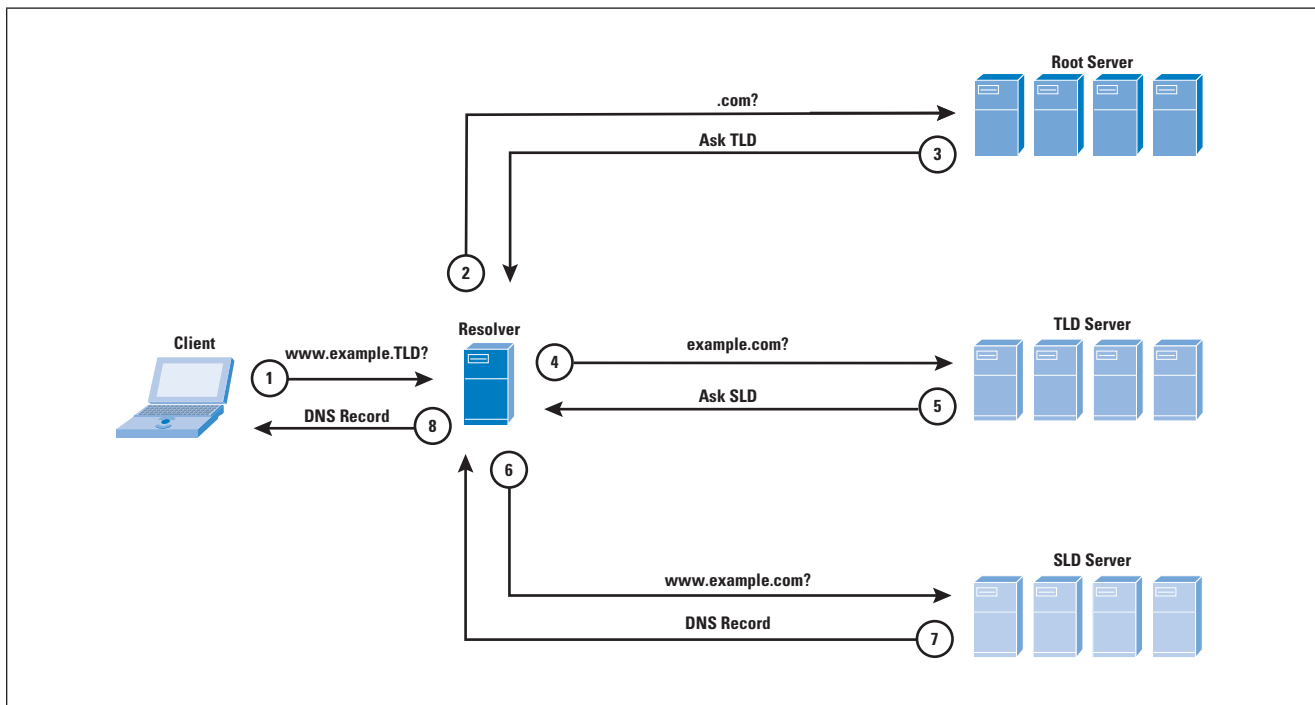
The RFC adopts the omitted-label approach for *query name* (or *qname*) minimization (or as it is spelled in the RFC, “minimisation”).

A resolver implementing *qname* minimization, as described in RFC 9156, takes advantage of information about how the DNS hierarchy is organized today at its higher, navigational levels, such as the root server delegating authority for existing TLDs to TLD servers, and typical TLD servers delegating authority for existing SLDs to SLD servers. As shown in Figure 2, when the resolver queries the root server as part of resolving a domain name, it sends only the TLD label to the root server. When it queries the TLD server, it sends only the SLD and TLD labels. It can also take advantage of the knowledge that some TLD servers delegate authority for some of their hierarchy at the third level rather than the second level, as discussed previously, thereby saving a step in those cases. The *Public Suffix List* (PSL)^[15] is a potential source for information about where these delegations or zone cuts may occur, as Geoff Huston has observed^[16].

In addition to replacing or removing labels, the resolver can also change the *Query Type* (QTYPE) from the one the client requested, to further reduce the amount of information disclosed. RFC 9156 recommends setting the QTYPE to “A” or “AAAA” regardless of the actual record type of interest, except for the final query with the full query name.

A resolver can apply *qname* minimization to its interactions with any authoritative name server at any level of the DNS hierarchy, and the name server won’t have to do anything differently. The name server will just receive queries with less information in them, except for the final name server in the chain.

Figure 2: DNS Resolution with *qname* Minimization



Qname minimization therefore provides a valuable information protection tool for both resolver operators and their users. Indeed, as Basileal Imana, Aleksandra Korolova, and John Heidemann state in their study of institutional privacy risks, “The currently available best way for institutions to reduce information leakage is to run their own resolver, and deploy query name minimization”.^[17]

Resolver operators have encountered one complication in deploying qname minimization: the *Empty Non-Terminal (ENT) Problem*, as described in RFC 7816^[14]. The problem can cause a resolver to continue to send queries during the minimized iterative resolution process, even after it should have become clear that the fully qualified domain name doesn’t exist. Retrying with the fully qualified domain name in the presence of ENTs wouldn’t disclose more information than the resolver would have been sent with traditional DNS resolution, but it would generate unnecessary additional queries. While it has become common practice simply to stop qname minimization after three labels, the underlying ENT problem remains. Resolving this complication is the focus of the next technique.

NXDOMAIN Cut Processing

NXDOMAIN, the negative answer in DNS, technically means that a domain name doesn’t exist—and therefore, by definition, that it has no subdomains.

However, because of the ENT ambiguity just mentioned, resolvers have traditionally limited their interpretation of NXDOMAIN to the domain name itself. This tradition has resulted in both additional workload for the resolver and extra traffic to the name-server system.

NXDOMAIN cut processing, described in RFC 8020^[8], expands the interpretation. As the title of the RFC states, a resolver implementing this technique interprets NXDOMAIN as “there really is nothing underneath;” the DNS tree is “cut.” In support, the RFC, authored by Stéphane Bortzmeyer and Shumon Huque, updates the DNS specifications to state that a name server must return NODATA in response to a query for an ENT, thereby resolving the ENT ambiguity.

Similar to qname minimization, a resolver can apply NXDOMAIN cut processing to its interactions with any authoritative name server. The name server doesn’t have to do anything differently as long as it handles ENT queries correctly. It will just receive less traffic.

With the root zone not having any ENTs, and with careful consideration given to the risks of ENTs in TLD zones^[18], it’s reasonable for resolvers to implement NXDOMAIN cut processing at the root and TLD levels of the DNS hierarchy, consistent with the deployment of qname minimization at those levels. Processing for additional zones can be enabled as resolver operators gain more confidence in the corresponding name servers’ handling of ENTs.

Or resolvers could simply adopt the technique unilaterally, regardless of the name server's behavior, a decision endorsed by RFC 8020:

“Such name servers are definitely wrong and have always been. Their behaviour is incompatible with DNSSEC. Given the advantages of ‘NXDOMAIN cut,’ there is little reason to support this behavior.”

NXDOMAIN cut processing helps qname minimization by enabling a resolver to stop the minimized iterative resolution process as soon as it receives an NXDOMAIN answer, meaning that the resolver will disclose less information in its traffic when a domain name doesn't exist, just as it discloses less when a domain name does exist. The combination of the two techniques can also be effective in defending against certain attacks, such as random subdomain attacks, where an adversary queries one or more resolvers for random subdomains of a common ancestor. With traditional processing, the resolvers will forward the subdomain queries to the ancestor's name server, generating a volumetric attack where the name server can't see the attack's original source. If the ancestor exists and is protected by DNSSEC, then aggressive DNSSEC caching can help a resolver reduce the number of additional subdomain queries that it forwards, as described by Petr Špaček.^[19] If the ancestor doesn't exist, then NXDOMAIN cut processing can keep the resolver from forwarding further subdomain queries after it knows that the ancestor doesn't exist.

With NXDOMAIN cut processing, a resolver broadens its interpretation of a negative answer to draw conclusions about subdomains of a domain name that it previously queried for. The next technique does something similar for negative answers in the DNSSEC case, drawing conclusions about other domain names in the zone as well.

Aggressive DNSSEC Caching

As discussed previously, negative answers in DNSSEC—in the form of NSEC and NSEC3 records—indicate that no domain names exist between two endpoints in some ordering of domain names. (One further detail: with the opt-out flag set in NSEC3, some domain names between the endpoints may actually exist, but not have DNSSEC-signed delegations. If a resolver is interested only in domain names that can be validated with DNSSEC, then the NSEC3 record is still useful information.)

Resolvers traditionally haven't taken advantage of the information these records provided about the nonexistence of other names, however.

Even though NSEC and NSEC3 records provide enough information for a resolver to conclude on its own that other domain names between the endpoints don't exist, resolvers have traditionally limited their interpretation to the domain name that was queried for.

Although authoritative name servers return NSEC or NSEC3 records in response to queries for both nonexistent domain names and ENTs, it's possible to tell the two classes apart, as detailed in Appendix B of RFC 8198^[9] for NSEC and in Sections 8.4–8.8 of RFC 5155^[5] for NSEC3.

The narrow interpretation is actually the correct one according to the original DNS specifications, not the result of an ambiguity as it was for the previous technique. RFC 4035^[3] describes the limitation as a “prudent” approach:

“In theory, a resolver could use wildcard or NSEC RRs to generate positive and negative responses (respectively) until the TTL or signatures on the records in question expire. However, it seems prudent for resolvers to avoid blocking new authoritative data or synthesizing new data on their own. Resolvers that follow this recommendation will have a more consistent view of the namespace.”

The limitation may once again result in the resolver doing more processing and sending more queries than it needs to, given the information it already has on hand.

Aggressive DNSSEC caching, described in RFC 8198^[9], takes a broader interpretation. The RFC, authored by Kazunori Fujiwara, Akira Kato, and Warren Kumari, updates the DNS specifications to state that a resolver may handle client queries for domain names that fall between the endpoints of previously received NSEC and NSEC3 records on its own. (It also allows the resolver to apply wildcard records to names between the endpoints when matching wildcard records exist.)

The technique offers an excellent illustration of the relative nature of the minimum disclosure principle, and it also improves the resolver's protection against random subdomain attacks as detailed in the previous section. Without DNSSEC, a resolver would need a name server's help for each new domain name it processes that's not a subdomain of a nonexistent domain. With DNSSEC, the resolver no longer needs as much help, so the threshold for minimum disclosure is reduced.

Similar to the two previous techniques, a resolver can apply aggressive DNSSEC caching to its interactions with any name server at any level. The name server again doesn't have a direct operational role and will just receive less traffic. The name server must handle NSEC or NSEC3 correctly, which is less of a concern than for NXDOMAIN and ENTs, given that the DNSSEC accounts for ENs.

The foregoing has three caveats:

First, as mentioned already, if an NSEC3 record has an opt-out flag, the resolver can't conclude that other domain names between the endpoints don't exist, only that they don't have secure delegations.

It therefore can't apply aggressive DNSSEC caching to such a record. Given that NSEC3 is the predominant choice for TLDs, and that the opt-out flag is commonly used^[20], aggressive DNSSEC caching will generally not help at the TLD level.

Second, the reduction in the number of queries sent assumes that the NSEC or NSEC3 endpoints actually span multiple domain names. Both techniques have variants, documented in RFC 4470^[20] and RFC 7129^[21], where the returned endpoints effectively span only the one domain name of interest, taking away the advantage of aggressive DNSSEC caching. Moreover, some implementations of these variants incorrectly report that some resource record types don't exist, possibly resulting in a resource record becoming unresolvable.^[22, 23] The "aggressive" interpretation of negative DNSSEC responses makes implementation errors more consequential as well.^[24]

Third, as Geoff Huston has observed,^[25] "the results [of aggressive DNSSEC caching] may not be that promising" for resolvers that load-balance their queries into servers with independent caches, for example, based on a hash of the query name.

These caveats aside, if the resolver were somehow to cache every NSEC or NSEC3 record in a pre-signed zone, and if there were no NSEC3 opt-outs, and if the ranges within the records collectively spanned the entire zone, then the resolver would be able to handle queries for every nonexistent domain name in the zone on its own, for as long as the records were valid.

If the resolver likewise were to cache every existing DNS record in the zone, then it could handle queries for existing domain names too.

A resolver might be able to bring all of these records into its cache if the set of queries it sends is directed, at least in part, by a carefully designed process. Using aggressive DNSSEC caching, the resolver will cache the NSEC record as evidence that domain names between the endpoints don't exist. In addition, it can cache the record as evidence that the two domain names at the endpoints do exist. Then the resolver can simply query for the DS and NS records for the two domain names at the endpoint, and it will have obtained the full DNS records for the two endpoints from this zone file. The resolver can repeat the process with other random long domain names until it has obtained a set of NSEC records that collectively span the zone. After sending queries for the zone's own DNSKEY, NS, and SOA records, the resolver will have obtained all the DNS records in the zone file. If the resolver implements NXDOMAIN cut processing and aggressive DNSSEC caching, it will then be able to answer client queries for every domain name without making further queries to the zone's authoritative name server. This process does not work well for NSEC3 and NSEC5.^[26, 27] By populating the resolver's cache in this way, the client would remove its own and other clients' interests in domain names from future resolver-to-authoritative queries.

But if the resolver just wants to avoid sending queries to a remote name server for a zone entirely, the next technique offers a more direct way to achieve the goal if the zone is appropriately configured.

Locally Served Zones

The DNS resolution processes shown in Figures 1 and 2 maintain a clear distinction between DNS ecosystem components: the client is separate from the resolver, which in turn is separate from the authoritative name servers. The separation implies a potentially global communications path between components, leading to the information disclosure concerns that have been the focus of this article.

But what if the communications path between two components was instead a local one? Such locality would not be unprecedented. Indeed, the resolver is often located within the same network as the client, which as discussed previously was one of the reasons for the relatively late standardization of an encrypted DNS protocol for the client-to-resolver exchange. An authoritative name-server instance can similarly be located within the same network as the resolver, as long as it can somehow be provisioned with a current copy of the zone file.

RFC 8806^[10], authored by Warren Kumari and Paul Hoffman, describes how to run a local instance of authoritative zone data with two constraints. First, the specification is limited to the root zone. Second, the local instance must indeed be run locally: that is, it must be accessible only to the resolver, and therefore not visible to other servers on the network. (Deploying the local instance at a loopback address, as proposed in the title to RFC 7706^[28], the predecessor to RFC 8806, is one way to ensure locality.)

ICANN's CTO organization describes the local root technique as "hyperlocal," and its OCTO-016 technical note^[29] proposes the technique as a way to "[improve] the decentralization of the root name service to mitigate risks that the [Root Server System] may face over time."

While OCTO-016 focuses on improving decentralization, and RFC 7706, per its title, on decreasing access time, it's also clear that the locally served zones technique also reduces the amount of information disclosed on the resolver-to-authoritative exchange. Indeed, RFC 8806 states that in addition to decreasing access time (particularly for negative responses), another goal of the technique is "to prevent queries and responses from being visible on the network."

A resolver can in principle get a copy of a zone file just like an authoritative name server might, via a zone-transfer protocol such as *Authoritative Transfer* (AXFR), described in RFC 5936^[30], and *Incremental Transfer* (IXFR), described in RFC 1995^[31]. These protocols give options for downloading a full zone file and for obtaining incremental updates respectively and may be enabled by a name server, depending on zone policy.

An encrypted version of these protocols, called *XFR-over-TLS* (XoT), is currently in development^[32]. Another alternative is for the zone data to be made available for download at a web address via the *Hypertext Transfer Protocol Secure* (HTTPS) protocol. For instance, ISI's *LocalRoot* project^[33] provides access to copies of the root zone, as well as the **.arpa**, **root-servers.net** and **dnssec-tools.org** zones.

In addition, the new ZONEMD record, described in RFC 8976^[34], provides a way to authenticate the integrity of a downloaded zone file (in contrast to DNSSEC, which authenticates individual sets of records).

Locally served zones and zone digests are more practical for small, slowly changing zones, such as the root zone, than for large, fast-changing ones. RFC 8976 states:

“ZONEMD is impractical for large, dynamic zones due to the time and resources required for digest calculation.”

The locally served zones technique, like others in this article, is another one that a resolver can apply to any zone at any level, in this case as long as zone data is made available for download. Its operational characteristics are similar to the other techniques: the name server doesn't need to do anything differently; the changes are all on the resolver's side (in terms of the resolution protocol); and the name server will receive less traffic (in this case, no traffic). The zone operator will need to provide a zone-transfer service, but this change is in provisioning, rather than resolution.

The technique does come with one significant caveat. The traditional DNS architecture with its resolver-to-authoritative exchanges has been optimized for the case where the operator for a zone is aware (or in the case of the root server, the multiple operators are collectively aware) of all of the name servers that are serving the zone. The operator(s) are therefore in a position where they can potentially check the consistency of the zone file information served by all these servers.

Until new mechanisms for synchronization are in place, locally served zone instances would fall outside a typical zone operator's awareness and ability to check consistency. OCTO-016 recognizes the need for additional work in stating:

“If hyperlocal were to see a significant uptake, a new system for root zone distribution would need to be devised to satisfy the reliability and scalability requirements associated with the widespread hyperlocal deployment in recursive resolvers.”

A system with these characteristics will be important if and when resolvers do adopt the locally served zones technique more broadly. But in the meantime, for resolvers that implement locally served zones, the technique will achieve the ultimate in minimum disclosure of information about client interests in domain names in the zone. The traditional resolver-to-authoritative exchange for these zones will have no conventional DNS queries at all.

Implementation Status

The minimization techniques described in the previous four sections are gradually being implemented and deployed in the DNS ecosystem. The following is a sampling of support by selected resolver operators and open source resolvers as of this writing.

A note on methodology: The distributed DNS ecosystem has tens of millions of resolvers^[35]. The ones referenced here are based on the list of “major Open DNS resolvers” in Huston’s qname minimization deployment study^[36], plus those in Mozilla’s *Trusted Recursive Resolver* (TRR) program^[37]. The determination of whether a resolver supports a technique is based primarily on public announcements. However, Huston’s study is also cited as likely evidence of qname minimization support. The open source resolver packages considered match the list included in Wouter de Vries et al.’s paper on qname minimization.^[38]

Qname Minimization

Qname minimization is included in the BIND^[39], *Knot Resolver*^[40], *PowerDNS*^[41], and *Unbound*^[42] open source resolver software packages. Cisco *Umbrella*^[43], Cloudflare’s 1.1.1.1^[44], Comcast’s *Xfinity Internet Service*^[45] (by virtue of its inclusion in Mozilla’s TRR program, which requires the capability), *Google Public DNS* (as related by Moura *et al.*^[46]), and *NextDNS*^[47] have all announced that they have implemented qname minimization. Google Public DNS has also reported that it uses a “nonce prefixes” technique where extraneous labels are replaced with a random label, an example of the “false-label” approach mentioned previously.^[48]

In addition, *dnswatch*, *dyn Recursive DNS*, *Quad9*, Neustar *Ultra-DNS Public*, and *Hurricane Electric* (HE) resolvers have been observed in Huston’s study as likely to be supporting qname minimization. The deployment of qname minimization has also been the subject of Internet measurement studies. De Vries observed that as early as April 2017, “0.9% (82 of 9,611) of RIPE Atlas probes had at least one [qname-minimizing] resolver,” and by October 2018, the percentage had grown to 11.7%.^[49] As of August 2021, NLnet Labs’ measurement dashboard showed that 47.8% of such probes interact with a qname-minimizing resolver.^[50]

Huston reported that as of mid-2020, “some 18% of users pass their queries through resolvers that actively work to minimize the extent of leakage of superfluous information in DNS queries,” adding that the percentage had increased from 3% since a year prior. Huston later clarified that the percentages likely underestimate actual adoption because the study’s active DNS measurement technique uses four-label client queries. Many resolver implementations of qname minimization revert to ordinary DNS resolution after three labels, potentially making the particular measurement technique undetectable by the study’s test servers.^[51]

In the same timeframe, according to research published by Matt Thomas,^[52] nearly half of all queries received by the **.com** and **.net** TLD servers consisted of only two labels.

The comparable percentage two years prior was 30%. The increase in two-label queries was accompanied by a similar decrease in three-label queries and thus can be taken as an indicator that qname minimization is being deployed at many resolvers. The upward trend has continued, reportedly reaching 55% as of February 2021.^[53] It should be noted, however, that many factors contribute to the composition of traffic to authoritative name servers, and the fraction of queries that have a certain number of labels may not be directly reflective of the fraction of resolvers that support qname minimization, nor with the fraction of users who interact with such resolvers.

NXDOMAIN Cut Processing

Knot Resolver,^[54] PowerDNS,^[55] and Unbound^[56] all support NXDOMAIN cut processing. BIND lists the technique as supported but made obsolete by Aggressive DNSSEC Caching.^[57] No announcements by recursive DNS operators have been found as of this writing. However, it is likely that many do support the technique, given that, as discussed previously, NXDOMAIN cut processing is not a new feature but rather the lack of accommodation for an old bug.

Aggressive DNSSEC Caching

Aggressive DNSSEC caching is included in BIND^[58], Knot Resolver^[59], PowerDNS^[60], and Unbound.^[61, 62] Cloudflare has reported that it has implemented aggressive DNSSEC caching,^[44] as well as Google.^[63]

Hyperlocal Zones

BIND^[64], Knot Resolver^[65] (following a “pre-filling” technique that RFC 8806 reports is consistent with the RFC’s requirements, but which diverges from the technique specified in RFC 7706), and Unbound^[56] all support hyperlocal zones. No announcements by recursive DNS operators were found.

Impact on DNS Measurement Research

The resolver-authoritative exchange has historically given authoritative name servers at all levels of the DNS hierarchy insights into the domain names being queried by a resolver’s clients. While the recursive, cached architecture of the DNS ecosystem conceals the identity of the specific client that originated a query, the receipt of a fully qualified domain name by an authoritative name server nevertheless reveals that *some* client is interested in the name. With the traditional DNS resolution process, that information potentially reaches all levels, starting with root and TLD.

One of the studies facilitated by this information was the DNS community’s research into name collisions related to the introduction of new *generic TLDs* (gTLDs) to the global DNS.

Root-server traffic already had shown significant evidence that resolvers (and therefore clients) were making many queries for domain names in TLDs that were not part of the global DNS^[66]. The root servers had historically, and correctly, responded that such domain names didn't exist, leading clients to query for different domain names (or to give up). But if a new TLD were added to the global DNS, the root servers (together with other servers) might begin to respond positively to client queries for domain names in the TLD. That change might then cause legacy clients to connect, inadvertently, to new, external servers—a name collision.

Because root servers had information about non-existent TLDs of interest to clients, as well as fully qualified domain names, researchers were able to determine not only which new gTLDs were already being queried for, but also which domain names within those new gTLDs were being queried. One of the sources for this information was the *Day in the Life* (DITL) exercise run annually by the *DNS Operations Analysis and Research Center* (DNS-OARC)^[67]. Researchers also performed additional analysis based on their own data sources and reported findings at a workshop on name collisions^[68].

The insights from root-server query data led to the identification of various network and client configurations that might be at risk if a new gTLD were delegated. For example, researchers identified vulnerabilities related to the *Web Proxy Auto-Discovery Protocol* (WPAD)^[69, 70, 71]. Researchers also found an operating-system vulnerability that did not involve new gTLDs based on their review of root-server query data^[72]. Verisign later conducted an outreach program that mitigated a broad range of name collision risks, again drawing from the query data^[73].

It is quite possible that if the minimization techniques described in this article had been broadly adopted a decade ago, researchers would not have been as able to study name collision risks as effectively, at least based on analyzing root-server data. One of the co-discoverers of independent vulnerability, *simMachines*—co-discoverer of the bug—is quoted in a blog post on qname minimization^[6] as stating that the “analysis would have been partially impacted” if fully qualified domain names had not been visible in root-server traffic.

The loss of visibility is exactly what should be expected, inasmuch as the goal of each of the minimization techniques is to reduce root and TLD servers' visibility into clients' interests in domain names. Adoption of the techniques will impact DNS measurement research at root and TLD servers in different ways.

- Qname minimization and NXDOMAIN cut processing, which amplify one another, reduce root and TLD servers' visibility into the lower-level domains that a resolver (and by implication, its clients) may be interested in. As more resolvers adopt qname minimization with an omitted-label approach, the overall query traffic to the root servers will trend toward single labels, while the traffic to the TLD servers will trend toward two or three labels depending on the delegation structure.

If these techniques had been in place at a given resolver when the name collisions research was performed, the root-server data associated with this resolver would only have indicated the TLDs the resolver was interested in, not the fully qualified domain names. Potential collisions between legacy systems and new gTLDs might have been highlighted, but some of the detail that helped determine the reason for the query and the impact of a positive response may have been obscured.

- Aggressive DNSSEC caching similarly reduces root and TLD servers' visibility into a resolver's interests in non-existent domain names that happen to be between the NSEC or NSEC3 endpoints obtained in response to another recently queried non-existent domain name. If aggressive DNSSEC caching had been in place at a resolver during the name collisions research, the root-server data associated with the resolver may have provided only partial information about the non-existent TLDs the resolver and its clients were interested in. This limitation on visibility may also have made it harder to assess the degree of risk associated with a given new gTLD.
- Finally, hyperlocal zones reduce the visibility of a name server into a participating resolver's interests entirely. ICANN's CTO organization, in its technical analysis of the hyperlocal root-zone technique^[74], aptly summarizes the impact on telemetry as follows: "...one likely consequence of significant hyperlocal root service deployment will be a general decrease in knowledge about how the global DNS operates."

We could make similar observations about other observations and actions motivated by root-server data. For instance, Matt Thomas' and Duane Wessels' study of DNS traffic to the root generated by Chromium-based browsers^[75] depends on statistics about queries to the root servers for non-existent TLDs. While qname minimization would not affect the statistics (the queries are already a single label), aggressive DNSSEC caching might. And the "mysterious root query data"^[76] reported by Duane Wessels and Christian Huitema, which includes many query names consisting of random 12- and 13-character SLDs followed by existing TLDs, would not have been seen if the resolver(s) that sent the queries had implemented qname minimization with an omitted-label approach. (To be fair, initial community feedback^[77] attributes the data to a different approach to reducing the amount of information in queries to the root server: the "nonce prefixes" technique previously mentioned in connection with Google Public DNS^[48].)

As minimization techniques are applied to the resolver-to-root and -TLD exchanges, researchers will need to expand their use of data sets from other parts of the DNS ecosystem—appropriately anonymized and summarized for sharing—if they want to maintain a larger view of the types of queries that clients are making. There are already numerous approaches for sharing data outside the resolver-to-root and -TLD exchanges.

DNS-OARC already collects research data from other “busy and interesting DNS servers,” not just root servers. Passive DNS tools^[78] offer an alternative approach for analyzing DNS query traffic patterns at an ecosystem level. And query data specific to security vulnerabilities can be shared with general threat-indicator tools.

The resolver-to-root and -TLD exchanges themselves will likely still have interesting data for researchers as well. Indeed, studies of these exchanges will provide important insights into the deployment of minimization techniques, which will be a gradual process over many years. Such studies may give even more information about the configuration of individual resolvers than was previously available when resolver behavior was more uniform.

Researchers may also be able to infer statistical information about the resolver selections of certain client environments, by measuring how known changes in these environments are filtered through the resolvers of different configurations. One potentially fruitful area for such research: the new HTTPS resource record^[79]. The record is gradually being introduced with early support by Apple’s iOS 14 and macOS 11 operating-system betas^[80]. Clients that support the HTTPS record will typically make three queries to their resolver, for the A, AAAA, and HTTPS record types. Traditional resolvers will then forward queries of all three types to the root and TLD servers. But resolvers that implement qname minimization may send only minimized A type queries to get a referral to the server that is actually authoritative for all three. The presence of HTTPS queries on the resolver-to-root and -TLD exchanges for a given resolver will therefore be an indicator not only that the resolver likely isn’t yet applying qname minimization, but also that a portion of the clients that query for the HTTPS record type are using the resolver.

Just as minimization techniques represent a new chapter in DNS protocol evolution, they also will bring a new era in DNS measurement research. DNS resolution will still be taking place, although in different ways, and data analysis will still be possible, but with alternate arrangements. Such alternatives will likely depend more on active measurement techniques where clients send queries that are designed to be detectable even if minimized resolution is taking place. Both the practice and the study of DNS will go on.

Conclusion: Into the Penumbra

For the past few decades, as DNS resolution has followed the textbook DNS approach shown in Figure 1, DNS operators have had significant visibility into aggregate client interests in domain names. While the visibility, as noted earlier, has not included information about specific client identities, it has included fully qualified domain names, forwarded to each authoritative name server in the chain of referrals.

As minimization techniques are deployed, less information will be sent on the resolver-to-authoritative exchange, especially at the root and TLD levels, both because individual queries will include less information (for example, because of qname minimization), and because fewer queries will be sent (because of the other techniques). That's a gain for the need-to-know principle, which is the primary motivation for the change. But it's also a loss for DNS measurement research—at least for the passive measurement research based on assumptions that textbook DNS is deployed.

Because DNS resolvers are gradually deploying minimization techniques, rather than adopting all at once, they are like an eclipse: a slow and steady occlusion of the information content of the resolver-to-authoritative DNS exchange. The minimization eclipse likely will never be a total one, as many legacy DNS resolvers will continue doing what they've been doing all along. But its effects will be noticeable, and, inasmuch as the change in visibility will be novel—minimization techniques haven't been broadly deployed before—the effects will also be a motivation for new research.

Astronomical eclipses, too, have been a source of inspiration to researchers, perhaps most notably the famous Eddington experiment of 1919 (ironically, for the time of this present writing, in the midst of another global pandemic). Eclipses had long been studied, but the change in visibility of stars, or more precisely, of the observed location of starlight passing the Sun, had not been measured. Arthur Stanley Eddington and Frank Watson Dyson organized expeditions to Principe and Sobral to record the location of the Hyades, a group of stars, during a solar eclipse^[81]. The starlight's degree of deflection by the Sun's gravity confirmed Einstein's theory of General Relativity.

Whereas Eddington's team understandably focused on a single group of stars, the DNS community will have millions of resolvers to watch. Eventually, perhaps, minimization will reach a practical maximum. But in the meantime, each resolver will be impacted in its own ways by minimization techniques. Each will also provide unique insights about the global DNS, given the aggregate characteristics of its clients and how they use DNS. Each step along the way is therefore well worth studying. For DNS and Internet protocol researchers, the minimization eclipse is just starting, and the shadows are still partial. DNS resolution is entering the *penumbra*^[9].

Acknowledgements

- The idea for this article emerged from multiple conversations with my Verisign colleagues about the history and future of qname minimization. Special thanks to Danny McPherson for his foundational work and strategic direction in this area, and to Duane Wessels and Matt Thomas for their expert technical guidance on both details and data of the techniques.
- Geoff Huston, Chief Scientist of APNIC, Vladimír Čunát and Ladislav Lhotka of CZ.NIC, Puneet Sood of Google, Victoria Risk and her colleagues at ISC, and Benno Overeinder of NLnet Labs all gave generously of their time to review drafts and respond to questions.
- The article would not have reached final form without Zaid Albanna's leadership in arranging multiple rounds of internal and external reviews, and in coordinating the revisions based on reviewers' helpful feedback. A thank you also to Kim Kelly for her careful technical editing.
- Finally, a thank you to Ole Jacobsen for persevering with IPJ and to IPJ's anonymous reviewers for their important work in reviewing this and many other contributions.

References

- [0] Penumbra: A partially shaded area around the edges of a shadow, especially an eclipse. (Source: [Wiktionary.org](https://en.wiktionary.org/wiki/penumbra))
- [1] IANA, "Root Servers,"
<https://www.iana.org/domains/root/servers>
- [2] Paul V. Mockapetris, "Domain names – implementation and specification," RFC 1035, November 1987.
- [3] Roy Arends, Rob Austein, Matt Larson, Dan Massey, and Scott Rose, "Protocol Modifications for the DNS Security Extensions," RFC 4035, March 2005.
- [4] Roy Arends, Rob Austein, Matt Larson, Dan Massey, and Scott Rose, "Resource Records for the DNS Security Extensions," RFC 4034, March 2005.
- [5] Ben Laurie, Geoffrey Sisson, Roy Arends, and David Blacka, "DNS Security (DNSSEC) Hashed Authenticated Denial of Existence," RFC 5155, March 2008.
- [6] Burton Kaliski, "Minimum Disclosure: What Information Does a Name Server Need to Do Its Job?" in *Verisign Blog*, March 2015, <https://blog.verisign.com/security/minimum-disclosure-what-information-does-a-name-server-need-to-do-its-job/>
- [7] Stephane Bortzmeyer, Ralph Dolmans, and Paul Hoffman, "DNS Query Name Minimisation to Improve Privacy," RFC 9156, November 2021.

- [8] Stephane Bortzmeyer and Shumon Huque, “NXDOMAIN: There Really Is Nothing Underneath,” RFC 8020, November 2016.
- [9] Kazunori Fujiwara, Akira Kato, and Warren Kumari, “Aggressive Use of DNSSEC-Validated Cache,” RFC 8198, July 2017.
- [10] Warren Kumari and Paul Hoffman, “Running a Root Server Local to a Resolver,” RFC 8806, June 2020.
- [11] Sara Dickinson, Benno Overeinder, Roland van Rijswijk-Deij, and Allison Mankin, “Recommendations for DNS Privacy Service Operators,” RFC 8932, October 2020.
- [12] Carlo Contavalli, Wilmer van der Gaast, David Lawrence, and Warren Kumari, “Client Subnet in DNS Queries,” RFC 7871, May 2016.
- [13] Geoff Huston, “A Look at DNS Trends and What the Future May Hold,” in *CircleID Blog*, October 2020, <https://circleid.com/posts/20201028-a-look-at-dns-trends-and-what-the-future-may-hold/>
- [14] Stephane Bortzmeyer, “DNS Query Name Minimization to Improve Privacy,” RFC 7816, March 2016.
- [15] Mozilla Foundation, “Public Suffix List,” <https://publicsuffix.org/>
- [16] Geoff Huston, “DNS Query Privacy revisited,” in *APNIC Blog*, September 2020, <https://blog.apnic.net/2020/09/11/dns-query-privacy-revisited/>
- [17] Basileal Imana, Aleksandra Korolova, and John Heidemann, “Institutional privacy risks in sharing DNS data,” in *ANRW ’21: Proceedings of the Applied Networking Research Workshop*, pp. 69–75, ACM, July 2021.
- [18] Vincent Levigneron, “ENT was here!!!”, presented at OARC 25, Dallas, October 2016, <https://indico.dns-oarc.net/event/25/contributions/403/>
- [19] Petr Špaček, “Measuring Efficiency of Aggressive Use of DNSSEC-Validated Cache (RFC 8198),” presented at OARC 28, San Juan, March 2018, <https://indico.dns-oarc.net/event/28/contributions/509/>
- [20] Sam Weiler and Johan Ihren, “Minimally Covering NSEC Records and DNSSEC On-line Signing,” RFC 4470, April 2006.
- [21] R. (Miek) Gieben and W. (Matthijs) Mekking, “Authenticated Denial of Existence in the DNS,” RFC 7129, February 2014.
- [22] Peter Van Dijk, “DVE-2018-0003: inaccurate NSEC3 answer results in domain unreachability if the resolver applies aggressive negative caching,” September 2018, <https://github.com/dns-violations/dns-violations/blob/master/2018/DVE-2018-0003.md>

- [23] Peter Van Dijk, “DVE-2021-0001: inaccurate NSEC3 answer results in domain unreachability if the resolver applies aggressive negative caching,” June 2021,
<https://github.com/dns-violations/dns-violations/blob/master/2021/DVE-2021-0001.md>
- [24] Petr Špaček, “Error in DNSSEC implementation on F5 BIG-IP load balancers” October 2019,
<https://en.blog.nic.cz/2019/07/10/error-in-dnssec-implementation-on-f5-big-ip-load-balancers/>
- [25] Geoff Huston, “NSEC Caching Revisited,” presented at OARC 31, Austin, October 2019,
<https://indico.dns-oarc.net/event/32/contributions/717/>
- [26] Daniel J. Bernstein, “Breaking DNSSEC,” presented at *3rd Usenix Workshop on Offensive Technologies* (WOOT’09), August 2009,
<https://www.usenix.org/legacy/events/woot09/tech/>
Slides: <https://cr.yp.to/talks/2009.08.10/slides.pdf>
- [27] Sharon Goldberg, Moni Naor, Dimitrios Papadopoulos, Leonid Reyzin, Sachin Vasant, and Asaf Ziv, “NSEC5: Provably Preventing DNSSEC Zone Enumeration,” in *2015 Network and Distributed System Security Symposium*, February 2015.
- [28] Warren Kumari and Paul Hoffman, “Decreasing Access Time to Root Servers by Running One on Loopback,” RFC 7706, November 2015.
- [29] ICANN Office of the Chief Technology Officer, “ICANN’s Root Name Service Strategy and Implementation,” OCTO-016, October 2020, <https://www.icann.org/en/system/files/files/octo-016-26oct20-en.pdf>
- [30] Edward Lewis and Aalfred Hoenes, Ed., “DNS Zone Transfer Protocol (AXFR),” RFC 5936, June 2010.
- [31] Masataka Ohta, “Incremental Zone Transfer in DNS,” RFC 1995, August 1996.
- [32] Willem Toorop, Sara Dickinson, Shivan Sahib, Pallavi Aras, and Allison Mankin, “DNS Zone Transfer-over-TLS,” RFC 9103, August 2021.
- [33] USC/ISI, “LocalRoot – Serve Yourself the Root,”
<https://localroot.isi.edu/>
- [34] Duane Wessels, Piet Barber, Matt Weinberg, Warren Kumari, and Wes Hardaker, “Message Digest for DNS Zones,” RFC 8976, February 2021.
- [35] Marc Kührer, Thomas Hupperich, Jonas Bushart, and Christian Rossow, “Going wild: Large-scale classification of open DNS resolvers,” in *2015 Internet Measurement Conference Proceedings*, pp. 355–368, ACM, October 2015.

- [36] Geoff Huston, “DNS Query Privacy revisited,” in *APNIC Blog*, September 2020, <https://blog.apnic.net/2020/09/11/dns-query-privacy-revisited/>
- [37] “Security/DOH-resolver-policy,” in *MozillaWiki*, <https://wiki.mozilla.org/Security/DOH-resolver-policy>
- [38] Wouter Bastiaan de Vries, Quirin Scheitle, Moritz Müller, Willem Toorop, Ralph Dolmans, and Roland van Rijswijk-Deij, “A first look at QNAME minimization in the Domain Name System,” in *Passive and Active Measurement*, PAM 2019, pp. 147–160, Springer, March 2019.
- [39] Vicky Risk, “BIND to Add QNAME Minimization,” in *ISC Blog*, March 2018, <https://www.isc.org/blogs/bind-to-add-qname-minimization/>
- [40] Knot Resolver, “Knot Resolver 1.0.0 released,” May 2016, <https://www.knot-resolver.cz/2016-05-30-knot-resolver-1.0.0.html>
- [41] “PowerDNS Recursor 4.3.0 Released,” in *PowerDNS Technical Blog*, March 2020, <https://blog.powerdns.com/2020/03/03/powerdns-recursor-4-3-0-released/>
- [42] Ralph Dolmans, “Unbound QNAME minimization,” presented at OARC 24, Buenos Aires, March 2016, <https://indico.dns-oarc.net/event/22/contributions/332/>
- [43] Alexander Harrison, “Cisco Umbrella DNS and QNAME Minimization,” <https://support.umbrella.com/hc/en-us/articles/360032551931-Cisco-Umbrella-DNS-and-QNAME-Minimization>
- [44] Ólafur Guðmundsson, “Introducing DNS Resolver, 1.1.1.1 (not a joke),” in *The Cloudflare Blog*, April 2018, <https://blog.cloudflare.com/dns-resolver-1-1-1-1/>
- [45] “Comcast’s Xfinity Internet Service Joins Firefox’s Trusted Recursive Resolver Program,” in *Firefox News*, June 2020, <https://blog.mozilla.org/en/products/firefox/firefox-news/comcasts-xfinity-internet-service-joins-firefoxs-trusted-recursive-resolver-program/>
- [46] Giovane C. M. Moura, Sebastian Castro, Wes Hardaker, Maarten Wullink, and Cristian Hesselman, “Clouding up the Internet: how centralized is DNS traffic becoming?” in *IMC ’20: Proceedings of the 2020 Internet Measurement Conference*, pp. 42–49, ACM, October 2020, <https://dl.acm.org/doi/10.1145/3419394.3423625>
- [47] NextDNS, “Privacy Policy,” <https://nextdns.io/privacy>
- [48] Google Public DNS, “Prepending nonce labels to query names,” https://developers.google.com/speed/public-dns/docs/security?hl=en#nonce_prefixes

- [49] Wouter Bastiaan de Vries, “Improving Anycast with Measurements,” PhD dissertation, University of Twente, December 2019, <https://research.utwente.nl/en/publications/improving-anycast-with-measurements>
- [50] NLnet Labs, “Qname Minimization,” <https://dnsthought.nlnetlabs.nl/#qnamemin>
- [51] Geoff Huston, private communications, October 2021.
- [52] Matt Thomas, “Maximizing Qname Minimization: A New Chapter in DNS Protocol Evolution,” in *Verisign Blog*, September 2020, <https://blog.verisign.com/security/maximizing-qname-minimization-a-new-chapter-in-dns-protocol-evolution/>
- [53] Burton Kaliski, “Standardizing Confidentiality Protections for Domain Name System (DNS) Exchanges: Multiple Approaches, New Functionality,” in *IEEE Communications Standards Magazine*, September 2021.
- [54] Petr Špaček, “NXNSAttack: Upgrade Resolvers to Stop New Kind of Random Subdomain Attack,” in *RIPE Labs Blog*, May 2020, https://labs.ripe.net/author/petr_spacek/nxnsattack-upgrade-resolvers-to-stop-new-kind-of-random-subdomain-attack/
- [55] “Third alpha release of PowerDNS Recursor 4.3.0,” in *PowerDNS Technical Blog*, October 2019, <https://blog.powerdns.com/2019/10/29/third-alpha-release-of-powerdns-recursor-4-3-0/>
- [56] NLnet Labs, “Unbound RFC Compliance,” <https://nlnetlabs.nl/projects/unbound/rfc-compliance/>
- [57] ISC, “BIND 9.19.x,” <https://gitlab.isc.org/isc-projects/bind9/-/issues/53>
- [58] Ray Bellis, “Aggressive NSEC caching in BIND 9.12,” in *APNIC Blog*, February 2018, <https://blog.apnic.net/2018/02/06/aggressive-nsec-caching-bind-9-12/>
- [59] CZ.NIC Labs, “Knot Resolver 2.0.0 (2018-01-31),” January 2018, <https://knot-resolver.readthedocs.io/en/stable/>
- [60] “First Beta Release of PowerDNS Recursor 4.5.0,” in *PowerDNS Technical Blog*, March 2021, <https://blog.powerdns.com/2021/03/26/first-beta-release-of-powerdns-recursor-4-5-0/>
- [61] NLnet Labs, “Aggressive NSEC,” <https://unbound.docs.nlnetlabs.nl/en/latest/topics/privacy/aggressive-nsec.html>

- [62] Ralph Dolmans, “Aggressive use of the DNSSEC-Validated cache in Unbound,” in *NLnet Labs Blog*, April 2018, <https://medium.com/nlnetlabs/aggressive-use-of-the-dnssec-validated-cache-in-unbound-1ab3e315d13f>
- [63] Google Public DNS, <https://developers.google.com/speed/public-dns/docs/security#dnssec>
- [64] Edward Winstead, “Running a local copy of the root zone,” presented at APRICOT 2017/APNIC 43, February 2017, https://2017.apricot.net/assets/files/APIC674/localdnszone_1488009521.pdf
- [65] CZ.NIC Labs, “Root on loopback (RFC 7706),” <https://knot-resolver.readthedocs.io/en/v5.0.1/modules-rfc7706.html>
- [66] ICANN Security and Stability Advisory Committee, “Invalid Top Level Domain Queries at the Root Level of the Domain Name System,” SAC045, November 2010, <https://www.icann.org/en/system/files/files/sac-045-en.pdf>
- [67] DNS-OARC, “Day In The Life of the Internet,” <https://www.dns-oarc.net/oarc/data/dit1>
- [68] Verisign, “Workshop on Root Causes and Mitigation of Name Collisions (WPNC)–March 2014.”
- [69] US-CERT, “WPAD Name Collision Vulnerability,” Alert TA-144A, revised October 2016, <https://us-cert.cisa.gov/ncas/alerts/TA16-144A>
- [70] Qi Alfred Chen, Eric Osterweil, Matthew Thomas, and Z. Morley Mao, “MitM attack by name collision: Cause analysis and vulnerability assessment in the new gTLD era,” in *2016 IEEE Symposium on Security and Privacy (SP)*, IEEE, 2016, <https://ieeexplore.ieee.org/abstract/document/7546529>
- [71] Qi Alfred Chen, Matthew Thomas, Eric Osterweil, Yulong Cao, and Jie You, “Client-side Name Collision Vulnerability in the New gTLD Era: A Systematic Study,” in *CCS ’17: Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pp. 941–956, ACM, October 2017, <https://doi.org/10.1145/3133956.3134084>
- [72] Carnegie-Mellon University Software Engineering Institute, “Microsoft Windows domain-configured client Group Policy fails to authenticate servers,” Vulnerability Note VU#787252, February 2015, <https://www.kb.cert.org/vuls/id/787252>
- [73] Matt Thomas, “Verisign Outreach Program Remediates Billions of Name Collision Queries,” in *Verisign Blog*, January 2021, <https://blog.verisign.com/domain-names/verisign-outreach-program-remediates-billions-of-name-collision-queries/>

- [74] Roy Arends and Nicolas Antonello, ICANN Office of the CTO, “Hyperlocal Root Zone Technical Analysis,” OCTO-027, August 2021, <https://www.icann.org/en/system/files/files/octo-027-25aug21-en.pdf>
- [75] Matt Thomas, “Chromium’s impact on root DNS traffic,” in *APNIC Blog*, August 2020, <https://blog.apnic.net/2020/08/21/chromiums-impact-on-root-dns-traffic/>
- [76] Duane Wessels and Christian Huitema, “More Mysterious Root Query Traffic from a Large Recursive Operator,” presented at OARConline 35a, September 2021, <https://indico.dns-oarc.net/event/39/contributions/864/>
- [77] Geoff Huston, “Another DNS OARC meeting,” in *APNIC Blog*, September 2021, <https://blog.apnic.net/2021/09/14/another-dns-oarc-meeting/>
- [78] Florian Weimer. “Passive DNS replication,” in *FIRST Conference on Computer Security Incident Handling*, Volume 98, 2005, <https://www.first.org/resources/papers/conference2005/florian-weimer-paper-1.pdf>
- [79] Benjamin M. Schwartz, Mike Bishop, and Erik Nygren, “Service binding and parameter specification via the DNS (DNS SVCB and HTTPS RRs),” Internet Draft, work in progress, October 2022, draft-ietf-dnsop-svcb-https.
- [80] Tommy Pauly, “Encrypted DNS support in iOS and macOS,” IETF ADD Working Group mailing list, June 2020, https://mailarchive.ietf.org/arch/msg/add/MbOOWPVHRHM_wvbKhfHuzUTwimI
- [81] Watson Dyson, Arthur Stanley Eddington, and Charles Rundle Davidson, “A Determination of the Deflection of Light by the Sun’s Gravitational Field, from Observations Made at the Total Eclipse of May 29, 1919,” *Philosophical Transactions of the Royal Society of London, Series A, Containing Papers of a Mathematical or Physical Character*, 220(571–581), pp. 291–333, 1920. <https://royalsocietypublishing.org/doi/pdf/10.1098/rsta.1920.0009>

Also available from:

<https://w.astro.berkeley.edu/~kalas/labs/documents/dyson1920.pdf>

BURTON S. KALISKI JR. (bkaliski@verisign.com) is senior vice president and chief technology officer of Verisign. He leads Verisign’s long-term research program and is responsible for the company’s industry standards engagements, university collaborations, and technical community programs. He previously served as the founding director of the EMC Innovation Network, as vice president of research at RSA Security, and as the founding scientist of RSA Laboratories, where his contributions included the development of the *Public-Key Cryptography Standards* (PKCS). He received a doctorate, master’s degree, and bachelor’s degree in computer science from the Massachusetts Institute of Technology.

KINDNS

by Adiel Akplogan, ICANN

In September 2022, the *Internet Corporation for Assigned Names and Numbers* (ICANN) launched <https://kindns.org> to support its *Knowledge-Sharing and Instantiating Norms for Domain Name System and Naming Security* (KINDNS) initiative. KINDNS was developed to improve DNS operations by promoting voluntary adherence to a clear set of security best practices tailored to authoritative and recursive DNS operators.

This initiative aligns with ICANN’s strategic goal to “Strengthen DNS coordination in partnership with the DNS stakeholders to improve the shared responsibility for upholding the security and stability of the DNS.” In other words, ICANN plans to actively promote DNS ecosystem security and relevant best practices. The KINDNS initiative is one of many programs and projects ICANN supports to help make the DNS safer.

The Domain Name System plays a crucial role in connecting users to services on the Internet. Like the Internet itself, the underlying protocols or rules that govern DNS operation are open. This is one of the greatest strengths of the DNS and the Internet: These open protocols are responsible for connecting billions of devices instantaneously all over the world.

This strength can also be a weakness. The DNS was not designed for security. Actors may snoop on DNS traffic, forge DNS traffic, and engage in denial-of-service attacks on DNS operations, among other activities. Similarly, the security systems and best practices of the DNS, which support the Internet’s operation, are characteristically open and voluntary. A clear example of this is the uneven global deployment of the *Domain Name System Security Extensions* (DNSSEC), a security enhancement specification developed by the *Internet Engineering Task Force* (IETF) more than 20 years ago.

DNS security challenges, however, are not unique. Key to Internet security, or security issues in general, is the necessity to coordinate behaviors across systems. Security challenges call for collective action and the voluntary adherence to a set of ever-evolving behaviors and technologies.

ICANN is uniquely positioned, in close collaboration with its many partners and peers, to help in collectively mitigating specific forms of DNS security threats, that fall within its mission (see below for more details). The organization’s greatest strength, in many ways, is its partnerships, active community participation, and global engagement activities. This network of technical organizations and experts from various backgrounds can provide a formidable tool in helping make the global DNS safer.

What is KINDNS?

KINDNS is an ICANN initiative tailored to authoritative and recursive DNS operators to promote voluntary adherence to a clear set of security best practices. A challenge facing DNS security is implementing and maintaining security at the same level for all authoritative and recursive operators in the DNS ecosystem. Smaller operators struggle to keep abreast of the latest security measure improvements, while large operators may implement only measures which help them achieve their professional business goals. As a result, a patchwork of varying security practices among DNS operators has led to weaknesses that malicious actors may find and use.

Currently, KINDNS has three areas of focus:

- Promoting the adoption of DNS security practices through the operator community. This includes maintaining a dynamic information portal that promotes KINDNS practices, helps operators to self-assess their practices and offers guidelines on how to implement them.
- Soliciting and gathering feedback on the KINDNS guidelines to refine and identify areas of improvement and emerging best practices that may be candidates for future additions to the framework.
- Developing advanced tools for operators to conduct self-assessments and an observatory platform around key DNS security indicators that can help measure and assess the impact of KINDNS.

KINDNS Portal

ICANN has worked with its community to develop a baseline level of security operations, a relatively small set of mutually agreed practices that operators of any size can easily implement. ICANN published the initial version of KINDNS in September 2022 with its own dedicated website: <https://kindns.org>

Voluntary Self-Assessment Form

To learn more about the initiative, KINDNS offers operators a simple self-assessment tool to check their current security practices and offers suggestions where they could improve. The current version of the self-assessment tool does not run any code on an operator's server or make any real-time measurements. KINDNS only asks questions and generates a compliance report upon completion of the survey. The form is found here: <https://kindns.org/assessments-tools>

KINDNS Guidelines and Practices

The guidelines and practices promoted by KINDNS were designed to help operators identify and implement critical security best practices. ICANN, however, does not view the KINDNS initiative as the only source of information on DNS security. On the contrary, it encourages operators to also check and work in accordance with the security dictated by the construct of their own infrastructure.

Currently, the KINDNS Framework covers the following categories:

Guidelines for Authoritative Server Operations

TLD and Critical Zones and Other Second-Level Domains (SLD) Zones: There are two types of best practices for operators of authoritative servers: DNS security and DNS availability and resilience. To learn more about these guidelines see:

<https://kindns.org/critical-zones>

and

<https://kindns.org/other-sld-zones>

Guidelines for Recursive Server Operators

These guidelines are tailored for three specific categories of recursive server operators. Depending on a company's policy or business practice, it may be operating one or more of these types of resolvers. These three categories include:

- *Private Resolvers*: these are not publicly accessible and cannot be reached over the open Internet. They are typically found in corporate networks or other restricted-access networks. Private resolvers in some cases are part of a trusted computing domain (for example, *Active Directory*).
- *Shared Private Resolver Operators*: these are typically *Internet Service Providers* (ISPs) or similar hosting service providers. They offer DNS resolution services to their customers (including for mobile, cable, DSL, fiber residential and commercial users, and hosted servers and applications).
- *Public resolvers*: this category includes both open and closed public resolvers. Closed public resolvers are typically commercial DNS filtering or scrubbing services. These service providers are typically not ISPs and the clients sending queries to them are located on remote networks. Note, some operators of closed public resolvers may also offer a free tier service, which also makes them open public resolvers.

To learn more about these guidelines visit:

<https://kindns.org/recursive-server-operators>

Guidelines for Platform Hardening

KINDNS recommends that all operators pay careful attention to practices for hardening the platforms their DNS services use. There are three types of hardening practices:

- *Network Security*: these best practices are aimed at preventing unauthorized network access to DNS servers and ensuring internal traffic does not leak onto other networks.
- *Host and Service Security*: these best practices are aimed at improving the security of hosts running DNS services to reduce the likelihood of a host compromise, denial-of-service attack or other attack.
- *Customer-Facing Portal and Service Security*: these practices are aimed at supporting the security needs of customers.

To learn more about these guidelines, visit:

<https://kindns.org/platform-hardening>

Conclusion

The website launch marks the completion of this initial phase of the KINDNS initiative, where ICANN's Technical Engagement team worked closely with the operator community and DNS experts to identify and document DNS security threats and their mitigation measures, which are the basis for the current guidelines. We hope and expect to evolve them as the DNS and the Internet continue to evolve. We invite anyone interested in participating in this initiative to join the KINDNS mailing list at:

<https://mm.icann.org/mailman/listinfo/kindns-discuss>

If you have any questions, please contact the KINDNS team at:

kindns-info@icann.org

ADIEL AKPLOGAN is Vice President, Technical Engagement at ICANN. With more than 25 years' experience in the ICT industry (20 specifically in the Internet Technology Industry), Adiel previously served as CEO for AFRINIC (The African Network Information Centre), IT Director for Symbol Technology in France (2001–2003) and Director of New Technology at CAFÉ Informatique in Togo (1994–2000). He earned a graduate degree in Electrical Engineering and holds a Master's degree in E-Business and New Technology Management from Paris Graduate School of Management. Recognized as one of the Internet technology pioneers in Africa, he has contributed to technical capacity building and deployment of some of the first private Internet Service Providers in Africa from 1996 to 1999. He can be reached at: **adiel.akplogan@icann.org**

The Internet Protocol Journal is published under the "CC BY-NC-ND" Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an "as-is" basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Internet
Society



Diamond Sponsors

Your logo here!

Ruby Sponsors



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Kjetil Aas	Václav Brožík	Karlheinz Dölger	Barry Greene	Michael Jones
Fabrizio Accatino	Christophe Brun	Michael Dragone	Jeffrey Greene	Amar Joshi
Michael Achola	Gareth Bryan	Joshua Dreier	Richard Gregor	Javier Juan
Martin Adkins	Ron Buchalski	Lutz Drink	Martijn Groenleer	David Jump
Melchior Aelmans	Paul Buchanan	Aaron Dudek	Geert Jan de Groot	Anders Marius Jørgensen
Christopher Affleck	Stefan Buckmann	Dmitriy Dudko	Ólafur Guðmundsson	Merike Kao
Scott Aitken	Caner Budakoglu	Andrew Dul	Christopher Guemez	Andrew Kaiser
Jacobus Akkerhuis	Darrell Budic	Joan Marc Riera	Gulf Coast Shots	Christos Karayiannis
Antonio Cuñat Alario	BugWorks	Duocastella	Sheryll de Guzman	Daniel Karrenberg
William Allaire	Scott Burleigh	Pedro Duque	Rex Hale	David Kekar
Nicola Altan	Chad Burnham	Holger Durer	Jason Hall	Stuart Kendrick
Shane Amante	Colin Butcher	Mark Eanes	James Hamilton	Robert Kent
Marcelo do Amaral	Jon Harald Bøvre	Andrew Edwards	Darow Han	Jithin Kesavan
Matteo D'Ambrosio	Olivier Cahagne	Peter Robert Egli	Handy Networks LLC	Jubal Kessler
Selva Anandavel	Antoine Camerlo	George Ehlers	Stephen Hanna	Shan Ali Khan
Jens Andersson	Tracy Camp	Peter Eisses	Martin Hannigan	Nabeel Khatri
Danish Ansari	Ignacio Soto Campos	Torbjörn Eklöv	John Hardin	Dae Young Kim
Finn Arildsen	Brian Candler	Y Ertur	David Harper	William W. H. Kimandu
Tim Armstrong	Fabio Caneparo	ERNW GmbH	Edward Hauser	John King
Richard Artes	Roberto Canonico	ESdatCo	David Hauweele	Russell Kirk
Michael Aschwanden	David Cardwell	Steve Esquivel	Marilyn Hay	Gary Klesk
David Atkins	Richard Carrara	Jay Etchings	Headcrafts SRLS	Anthony Klopp
Jac Backus	John Cavanaugh	Mikhail Evstiounin	Hidde van der Heide	Henry Kluge
Jaime Badua	Lj Cemerar	Bill Fenner	Johan Helsingius	Michael Kluk
Bent Bagger	Dave Chapman	Paul Ferguson	Robert Hinden	Andrew Koch
Eric Baker	Stefanos Charchalak	Ricardo Ferreira	Asbjørn Højmark	Ia Kochiashvili
Fred Baker	Molly Cheam	Kent Fichtner	Damien Holloway	Carsten Koempe
Santosh Balagopalan	Greg Chisholm	Armin Fisslthaler	Alain Van Hoof	Richard Koene
William Baltas	David Chosrova	Michael Fiumano	Edward Hotard	Alexander Kogan
David Bandinelli	Marcin Cieslak	The Flirble Organisation	Bill Huber	Matthijs Koot
Benjamin Barkin-Wilkins	Lauris Cikovskis	Gary Ford	Hagen Hultzs	Antonin Kral
Feras Batainah	Brad Clark	Jean-Pierre Forcioli	Kauto Huopio	Robert Krejčí
Michael Bazarewsky	Narelle Clark	Susan Forney	Kevin Iddles	Mathias Körber
David Belson	Horst Clausen	Christopher Forsyth	Mika Ilvesmaki	John Kristoff
Richard Bennett	James Cliver	Andrew Fox	Karsten Iwen	Terje Krogdahl
Matthew Best	Guido Coenders	Craig Fox	Joseph Jackson	Bobby Krupczak
Hidde Beumer	Joseph Connolly	Fausto Franceschini	David Jaffe	Murray Kucherawy
Pier Paolo Biagi	Steve Corbató	Valerie Fronczak	Ashford Jaggernaut	Warren Kumari
Arturo Bianchi	Brian Courtney	Tomislav Futivic	Thomas Jalkanen	George Kuo
John Bigrow	Beth and Steve Crocker	Laurence Gagliani	Martijn Jansen	Dirk Kurfuerst
Orvar Ari Bjarnason	Dave Crocker	Edward Gallagher	Jozef Janitor	Darrell Lack
Tyson Blanchard	Kevin Croes	Andrew Gallo	John Jarvis	Andrew Lamb
Axel Boeger	John Curran	Chris Gamboni	Dennis Jennings	Richard Lamb
Keith Bogart	André Danthine	Xosé Bravo Garcia	Edward Jennings	Yan Landriault
Mirko Bonadei	Morgan Davis	Oswaldo Gazzaniga	Aart Jochem	Edwin Lang
Roberto Bonalumi	Jeff Day	Kevin Gee	Nils Johansson	Sig Lange
Lolke Boonstra	Julien Dhallenne	Greg Giessow	Brian Johnson	Markus Langenmair
Julie Bottorff Photography	Freek Dijkstra	John Gilbert	Curtis Johnson	Fred Langham
Gerry Boudreaux	Geert Van Dijk	Serge Van Ginderachter	Richard Johnson	Tracy LaQuey Parker
Leen de Braal	David Dillow	Greg Goddard	Jim Johnston	Alex Latzko
Kevin Breit	Richard Dodsworth	Tiago Goncalves	Jonatan Jonasson	Jose Antonio Lazaro
Thomas Bridge	Ernesto Doelling	Ron Goodheart	Daniel Jones	Lazaro
Ilia Bromberg	Michael Dolan	Octavio Alfageme	Gary Jones	Antonio Leding
Lukasz Bromirski	Eugene Doroniuk	Gorostiaga	Jerry Jones	Rick van Leeuwen

Simon Leinen	Mohammad Moghaddas	Andrew Potter	Timothy Schwab	Fabrizio Tivano
Robert Lewis	Roberto Montoya	Ian Potts	Roger Schwartz	Peter Tomsu Fine Art
Christian Libérale	Charles Monson	Eduard Llull Pou	SeenThere	Photography
Martin Lillepuu	Andrea Montefusco	Tim Pozar	Scott Seifel	Joseph Toste
Roger Lindholm	Fernando Montenegro	David Raistrick	Paul Selkirk	Rey Tucker
Link Light Networks	Joel Moore	Priyan R Rajeevan	Yury Shefer	Sandro Tumini
Chris and Janet Lonvick	John More	Balaji Rajendran	Yaron Sheffer	Angelo Turetta
Sergio Loreti	Maurizio Moroni	Paul Rathbone	Doron Shikmoni	Michael Turzanski
Eric Louie	Brian Mort	William Rawlings	Tj Shumway	Phil Tweedie
Adam Loveless	Soenke Mumm	Mujtiba Raza Rizvi	Jeffrey Sicuranza	Steve Ulrich
Josh Lowe	Tariq Mustafa	Bill Reid	Thorsten Sideboard	Unitek Engineering AG
Guillermo a Loyola	Stuart Nadin	Petr Rejhon	Greipur Sigurdsson	John Urbanek
Hannes Lubich	Michel Nakhla	Robert Remenyi	Fillipe Cajaiba da Silva	Martin Urwaleck
Dan Lynch	Mazdak Rajabi Nasab	Rodrigo Ribeiro	Andrew Simmons	Betsy Vanderpool
David MacDuffie	Krishna Natarajan	Glenn Ricart	Pradeep Singh	Surendran Vangadasalam
Sanya Madan	Naveen Nathan	Justin Richards	Henry Sinnreich	Ramnath Vasudha
Miroslav Madić	Darryl Newman	Rafael Riera	Geoff Sisson	Philip Venables
Alexis Madriz	Thomas Nikolajsen	Mark Risinger	John Sisson	Buddy Venne
Carl Malamud	Paul Nikolich	Fernando Robayo	Helge Skrivervik	Alejandro Vennera
Jonathan Maldonado	Travis Northrup	Michael Roberts	Terry Slattery	Luca Ventura
Michael Malik	Marijana Novakovic	Gregory Robinson	Darren Sleeth	Scott Vermillion
Tarmo Marners	David Oates	Ron Rockrohr	Richard Smit	Tom Vest
Yogesh Mangar	Ovidiu Obersterescu	Carlos Rodrigues	Bob Smith	Peter Villemoes
John Mann	Tim O'Brien	Magnus Romedahl	Courtney Smith	Vista Global Coaching
Bill Manning	Mike O'Connor	Lex Van Roon	Eric Smith	& Consulting
Harold March	Mike O'Dell	Marshall Rose	Mark Smith	Dario Vitali
Vincent Marchand	John O'Neill	Alessandra Rosi	Tim Sneddon	Rüdiger Volk
Normando Marcolongo	Jim Oplotnik	David Ross	Craig Snell	Jeffrey Wagner
Gabriel Marroquin	Carl Ötne	William Ross	Job Snijders	Don Wahl
David Martin	Packet Consulting	Boudhayan	Ronald Solano	Michael L Wahrman
Jim Martin	Limited	Roychowdhury	Asit Som	Laurence Walker
Ruben Tripana Martin	Carlos Astor Araujo	Carlos Rubio	Ignacio Soto Campos	Randy Watts
Timothy Martin	Palmeira	Rainer Rudigier	Evandro Sousa	Andrew Webster
Carles Mateu	Alexis Panagopoulos	Timo Ruiters	Peter Spekrijse	Tim Weil
Juan Jose Marin Martinez	Gaurav Panwar	RustedMusic	Thayumanavan Sridhar	Jd Wegner
Ioan Maxim	Chris Parker	Babak Saberi	Paul Stancik	Westmoreland
David Mazel	Manuel Uruena Pascual	George Sadowsky	Ralf Stempfner	Engineering Inc.
Miles McCredie	Ricardo Patara	Scott Sandefur	Matthew Stenberg	Rick Wesson
Brian McCullough	Dipesh Patel	Sachin Sapkal	Martin Štěpánek	Peter Whimp
Joe McEachern	Alex Parkinson	Arturas Satkovskis	Adrian Stevens	Russ White
Alexander McKenzie	Craig Partridge	PS Saunders	Clinton Stevens	Jurrien Wijlhuizen
Jay McMaster	Dan Paynter	Richard Savoy	John Streck	Derick Winkworth
Mark Mc Nicholas	Leif Eric Pedersen	John Sayer	Martin Streule	Pindar Wong
Olaf Mehlberg	Rui Sao Pedro	Phil Scarr	David Strom	Makarand Yerawadekar
Carsten Melberg	Juan Pena	Gianpaolo Scassellati	Colin Strutt	Phillip Yialeloglou
Kevin Menezes	Chris Perkins	Elizabeth Scheid	Viktor Sudakov	Janko Zavernik
Bart Jan Menkveld	Michael Petry	Jeroen Van Ingen	Edward-W. Suor	Bernd Zeimet
Sean Mentzer	Alexander Peuchert	Schenau	Vincent Surillo	Muhammad Ziad
William Mills	David Phelan	Carsten Scherb	Terence Charles	Ziayuddin
David Millsom	Harald Pilz	Ernest Schirmer	Sweetser	Tom Zingale
Desiree Miloshevic	Derrell Piper	Benson Schliesser	T2Group	Jose Zumalave
Joost van der Minnen	Rob Pirnie	Philip Schneek	Roman Tarasov	Romeo Zwart
Thomas Mino	Marc Vives Piza	James Schneider	David Theese	廖明沂.
Rob Minshall	Jorge Ivan Pincay Ponce	Peter Schoo	Douglas Thompson	
Wijnand	Victoria Poncini	Dan Schrenk	Kerry Thompson	
Modderman-Lenstra	Blahoslav Popela	Richard Schultz	Lorin J Thompson	



Follow us on Twitter and Facebook

@protocoljournal



<https://www.facebook.com/newipj>

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

John Crain, Senior Vice President and Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Principal – Internet Technology, Policy, and Advocacy
The Internet Society

Dr. Jun Murai, Founder, WIDE Project
Distinguished Professor, Keio University
Co-Director, Keio University Cyber Civilization Research Center, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

June 2023

Volume 26, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

FROM THE EDITOR

In This Issue

From the Editor	1
ALTO	2
Wi-Fi Privacy	12
Twenty-Five Years Later.....	23
Supporters and Sponsors	49
Thank You!	50

Twenty-five years ago, we published the first issue of *The Internet Protocol Journal* (IPJ). Since then, 87 issues for a total of 3,316 pages have been produced. Today, IPJ has about 20,000 subscribers all around the world. In the early days of IPJ, most of our readers preferred the paper edition, but over time preferences have shifted steadily to a situation where only some 1,200 print subscribers remain. The rest are downloading the PDF version. This shift in reading habits is likely related to the changes in technology that have taken place in the last 25 years. Lower costs and higher-resolution displays and printers, as well as improvements in Internet access technologies, have made the online “experience” a lot better than it was in 1998.

In this issue, we will first look at two areas of work taking place in the *Internet Engineering Task Force* (IETF). The *Application-Layer Traffic Optimization* (ALTO) protocol aims to make network state such as topology, link availability, routing policies, and path cost metrics information available to applications in a standardized manner. The next article concerns the thorny topic of *tracking* of users and their devices on the Internet. This area is complex, with many potential solutions, including the use of randomized *Media Access Control* (MAC) addresses as described by members of the *MAC Address Device Identification for Network and Application Services* (MADINAS) Working Group in the IETF.

Our final article is a look back at the last 25 years of Internet technology development. As we did with our 10th anniversary issue in 2008, we asked Geoff Huston to provide an overview of the many changes that have taken place in this period. At the end of his article, you will find a list of previously published articles from IPJ on numerous aspects of Internet technologies. All back issues are, of course, available from our website.

Let me take this opportunity to thank all the people who make IPJ possible. We are grateful to all our sponsors and donors, without whose generous support this publication would not exist. Our authors deserve a round of applause for carefully explaining both established and emerging technologies. They are assisted by an equally insightful set of reviewers and advisors who provide feedback and suggestions on every aspect of our publications process. The process itself relies heavily on two individuals: Bonnie Hupton, our copy editor, and Diane Andrada, our designer. Thanks go also to our printers and mailing and shipping providers. Last, but not least, our readers provide encouragement, suggestions, and feedback. This journal would not be what it is without them.

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

The IETF ALTO Protocol

Optimizing Application Performance by Increasing Network Awareness

by Qin Wu, Mohamed Boucadair, and Jordi Ros-Giralt

In today's Internet, network-related information (for example, topology, link availability, routing policies, and path cost metrics) are usually hidden from the application layer. As a result, endpoints make network-unaware decisions that may lead to suboptimal service placement and selection decisions, sometimes resulting in poor user experience and unnecessary inter-*Internet Service Provider* (ISP) traffic. Previous approaches to this problem space have considered snooping on the lower layers to determine the state and capabilities of the network, but such techniques require applications to be aware of lower-layer components (for example, routing protocols) and, furthermore, if left unspecified, can potentially overload key network resources.

To overcome this challenge, it is necessary to gather and expose network state information (for example, the bandwidth and latency properties between two network endpoints) to applications that do not interact directly with their underlying network protocols, without increasing the risk of network service disruption. For instance, empowered with such information, service providers can safely optimize the placement of their applications in locations of the network that provide higher capacity and lower latency to the clients they intend to serve. Similarly, with such information, client applications can also optimize the selection of the server instances they decide to attach to, while relying upon a variety of cost metrics.

This article provides an overview of how the *Internet Engineering Task Force* (IETF) *Application-Layer Traffic Optimization* (ALTO) protocol enables applications with improved network awareness to overcome these challenges, and reports on some of the implementations and deployments of the ALTO protocol.

The ALTO Approach and Architecture

The IETF ALTO protocol defines a client/server network service that applications can use to gain insightful information about the current state of the network. As defined in the base protocol^[1], each ALTO server maintains a “my-Internet” view of the network it represents. In its simplest form, this view consists of a set of endpoints and costs between pairs of endpoints for each possible cost type (for example, hop count, latency, or bandwidth). An application seeking to gain this information to make optimized decisions can use an ALTO client to connect to an ALTO server using the *Hypertext Transfer Protocol* (HTTP)-based protocol defined in the ALTO base specification.

The ALTO protocol uses a *Representational State Transfer* (RESTful) design and encodes its requests and responses using *JavaScript Object Notation* (JSON) objects.

An ALTO request carries a set of source-destination endpoints and a cost type. The triggered ALTO response provides the cost value for each given source-destination endpoint. To improve scalability (for example, to reduce the load of an ALTO server) and privacy (for example, to avoid revealing sensitive topology information), ALTO introduces the concept of *groups*, which specify sets of endpoints that are close to each other from a network connectivity standpoint. In larger-scale networks, this aggregation leads to greater scalability without losing critical information. A group may be represented as an IP prefix, a *Point of Presence* (PoP), a type of access connectivity (wireless, fiber, etc.), an *Autonomous System* (AS), or a set of ASes. The entity that operates an ALTO server, called the *ALTO Service Provider*, is responsible for assigning a unique *Provider-defined Identifier* (PID) to each group.

Another generalization of the endpoint object is enabled using the concept of *Abstract Network Element* (ANE). This concept provides an abstract representation of a component in a network that handles data packets and whose properties can potentially affect the end-to-end performance of an application^[14]. ANEs can include not only endpoints, but also switches and routers that connect them.

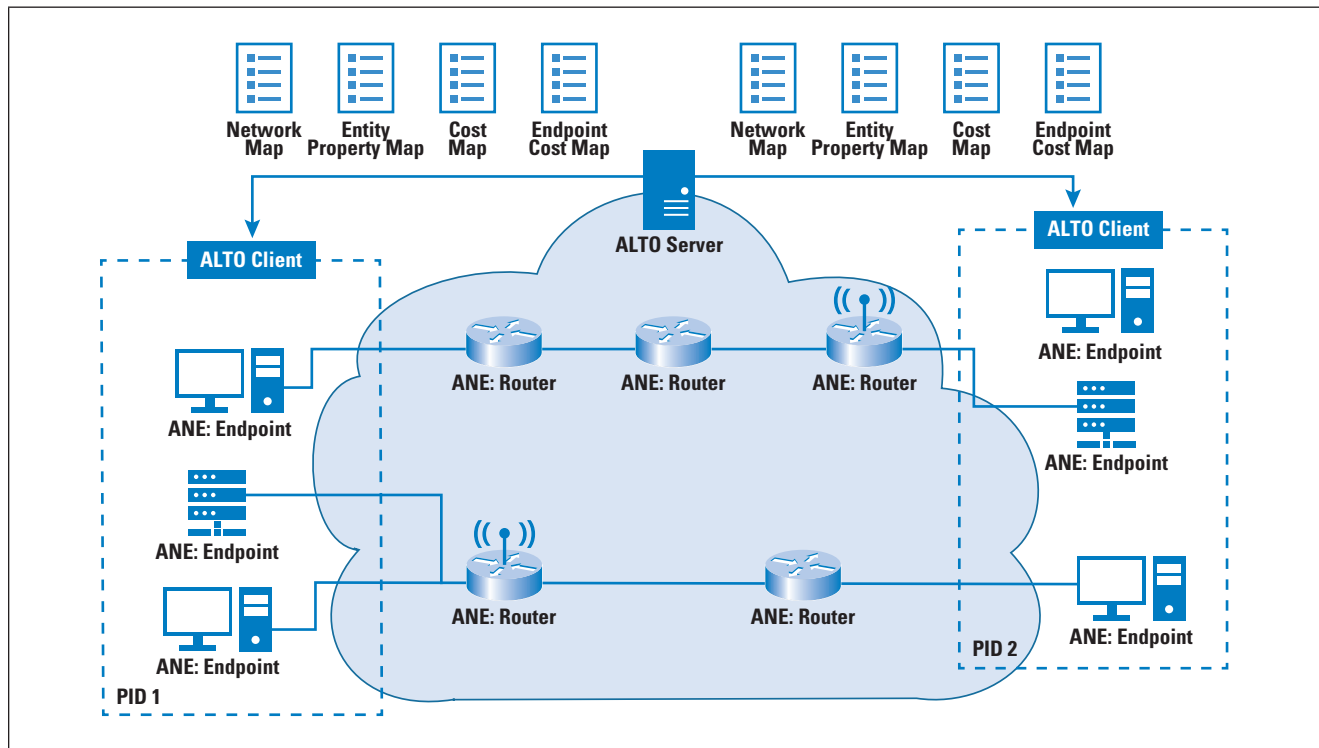
Figure 1, on the next page, depicts the main ALTO abstract objects involved in a network. In this figure, application endpoints are represented as ANEs clustered in two groups with provider-defined identifiers “PID1” and “PID2.” An endpoint can select to communicate with another endpoint based on the network properties that the ALTO server exposes. For instance, in a *Content Delivery Network* (CDN), ANEs correspond to client and server hosts, and a specific content (for example, a movie) is in general replicated in more than one server instance. A client host can decide to retrieve the content by selecting the server instance that provides the higher communication bandwidth according to the exposed ALTO information. Each of the abstract objects that are illustrated in Figure 1 is further elaborated in the following sections.

ALTO Maps

An ALTO server organizes the network information using the concept of *maps*. Maps can be constructed from physical information, logical information, or a combination thereof. ALTO supports four types of maps, as shown in Figure 1:

- The *Network Map* lists all the endpoint groups that the ALTO server tracks. This map includes PIDs that uniquely identify each group.
- The *Entity Property Map* describes the properties of each ANE in the network, including the geolocation or the connectivity type (for example, fiber or wireless) of an ANE.
- The *Cost Map* provides the cost information (for example, hop count, latency, or bandwidth) between each pair of PIDs enclosed in the network map, where a PID identifies a group of endpoints.
- The *Endpoint Cost Map* provides finer-grained cost information between specific endpoints.

Figure 1: Base ALTO Abstract Objects



ALTO Extensions

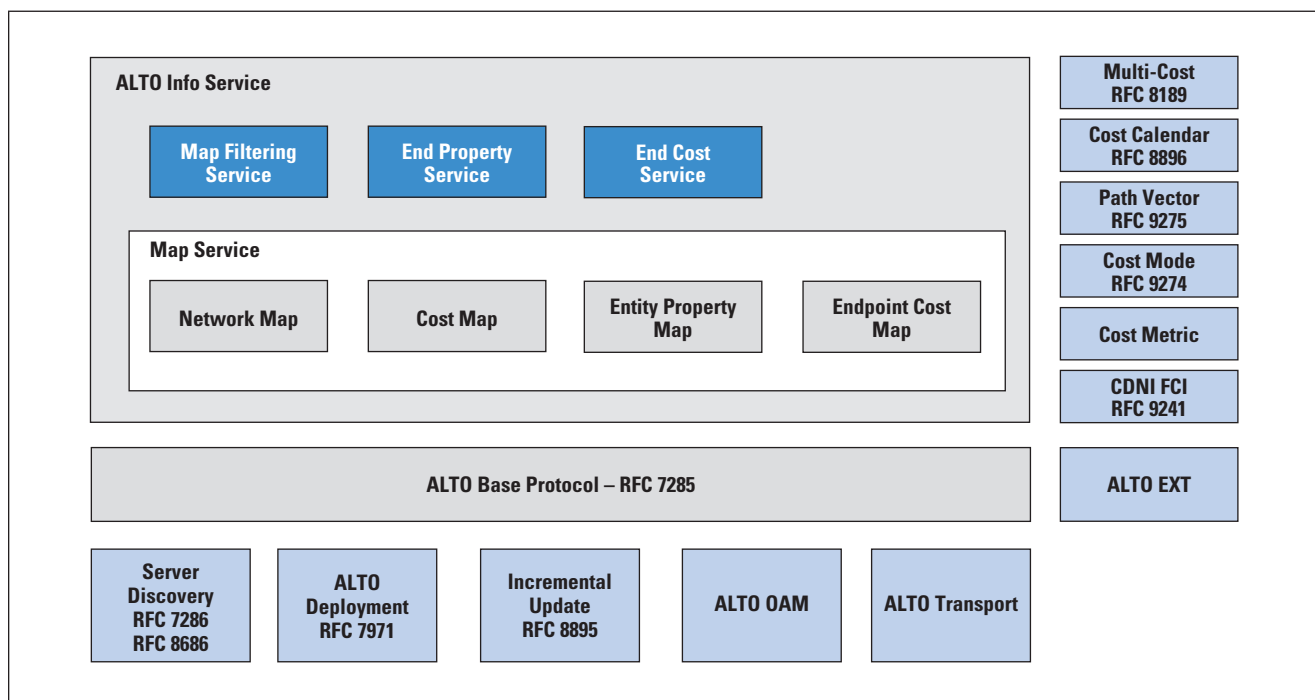
In addition to its base abstractions and maps, the ALTO protocol supports the following extensions to enable a richer network-aware application experience:

- The *Information Resource Directory* (IRD) lists the services an ALTO server provides and the locations from where you can access these services.
- The *Cost Calendar* provides a set of cost values as a function of time, allowing applications to know not only where to connect to, but also when.
- Incremental Updates using *Server-Sent Events* (SSEs) allow an ALTO server to expose cost values as delta updates, reducing the amount of server-client data exchanged.
- The CDNI Advertisement exposes a *CDNI Footprint and Capability Advertisement Interface* (FCI)^[26].
- The *Path Vector extension* exposes the set of ANEs along the path between two endpoints and the performance properties of these ANEs.
- The *Extended Performance Cost Metrics* enrich ALTO with advanced metrics such as network one-way delay, one-way delay variation, one-way packet-loss rate, hop count, and bandwidth.
- The *Entity Properties* generalize the concept of ALTO endpoint properties by presenting them as entity property maps.

Figure 2 shows the ALTO protocol core services as they are documented by the IETF as well as some of the related ALTO documents. The previously mentioned ALTO extensions are marked with a light-blue shading. The core services are organized as part of the *ALTO Information Service* consisting of the *Map Filtering*, *End Property*, and *End Cost* services, along with the *Map Service*, which is itself broken into separate map services as previously described. All of the services are dependent on the base protocol, which is documented in [1]. The ALTO protocol is enhanced through *Server Discovery*^[5, 25], and extensions for *Incremental Updates*^[9], *Operations and Management (OAM)*^[23], and support for carrying the ALTO protocol over more modern transport protocols^[22]. The practical understanding of how you can use the ALTO protocol together with a set of deployment recommendations is documented in [13].

Additional ALTO features, for example, cost manipulation^[7, 8], are shown on the right side of Figure 2.

Figure 2: Overview of ALTO Core and Extensions



History of ALTO

The ALTO Working Group was established in 2008 with an initial charter to develop a request/response protocol that would allow hosts to extract enough state information from a network to make optimized server selection decisions. The working group's first charter focused on the optimization of *Peer-to-Peer* (P2P) applications, with the first four RFCs introducing the problem statement^[11] and requirements^[4], the base protocol^[1], and support for server discovery^[5].

The working group was then rechartered in 2014 to support a broader set of applications that included CDNs and data centers.

That stage led to the development of five RFCs: Deployment recommendations^[13], protocol extensions for reducing the volume of on-the-wire data exchange^[7, 9], server discovery for multi-domain environments^[25], and a cost calendar capability to allow applications to identify the optimal times to connect to a service^[8].

The current ALTO Working Group charter was approved in 2021 with the goal to focus on three operational areas: (1) support for modern transport protocols such as HTTP/2 and HTTP/3^[22]; (2) development of OAM mechanisms^[23], and (3) collection of deployment experiences^[24]. These three areas constitute the current highest priorities of the ALTO Working Group.

Four additional RFCs that had originated from the second charter have also been published since then: (1) support of property maps for generalized entities^[10], (2) a new *Footprint and Capabilities Advertisement Interface* (FCI) protocol for CDNI^[12], (3) a new *Internet Assigned Numbers Authority* (IANA) registry for tracking cost modes supported by ALTO^[3], and (4) extensions to the cost map and ALTO property map services to allow the application to identify optimized paths^[14].

ALTO Deployments and Implementations

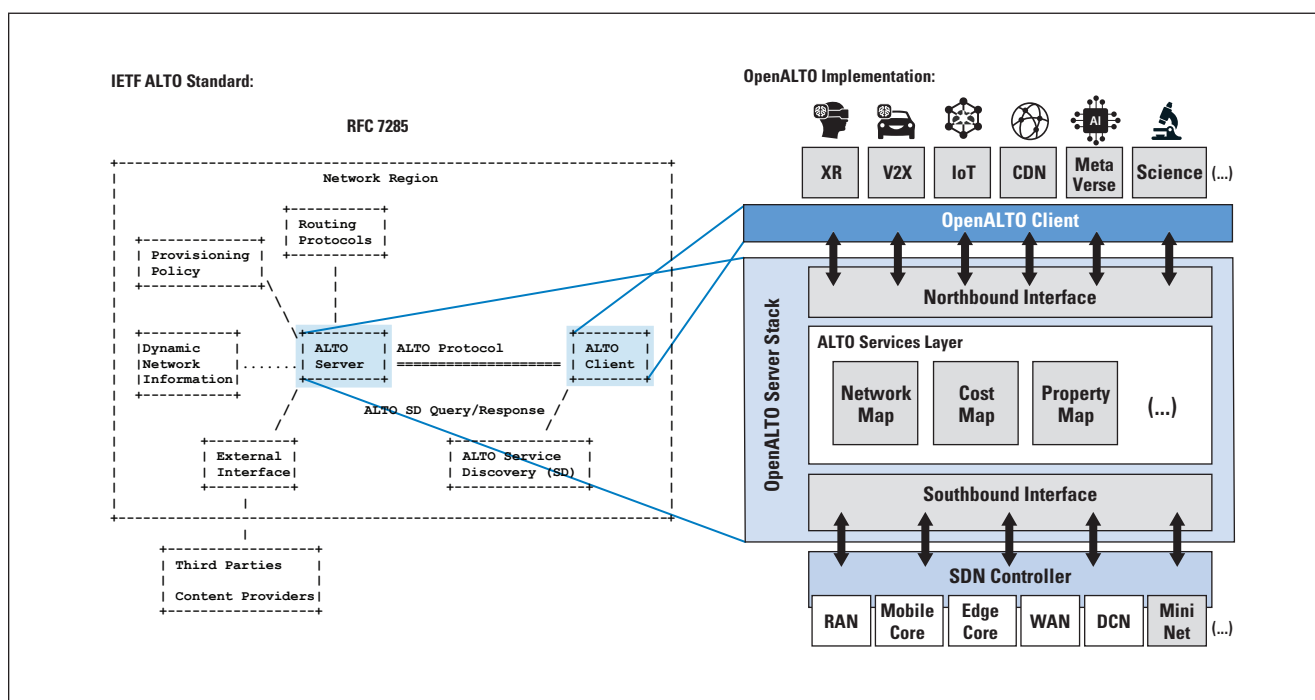
The ALTO base protocol was first implemented by Korea Telecom^[16], NEC^[17], Benocs^[18], Alcatel-Lucent Bell Labs^[19], and Nokia. An open-source implementation of the ALTO stack was also made available via the *OpenDaylight* (ODL) Project^[6]. Starting in 2020, China Mobile, Tencent^[20], and Telefonica^[21] have been actively involved in ALTO and initiated trials in their mobile networks and CDNs. Qualcomm Technologies, Inc. also joined the ALTO effort in 2021 with a focus on evaluating the fit of ALTO for exposing network state information in the context of edge computing. At the time of this writing (2023), two new deployments of ALTO are being initiated to support the networks from CERN (LHCONE) in Europe and the Network Research Platform in the United States^[27].

A further open-source initiative, the *OpenALTO* Project^[2], was initiated in 2021 to provide a standalone implementation of the ALTO stack independently of the ODL Project. The OpenALTO Project is an initiative spawning from the IETF ALTO Working Group, which focuses on developing an open-source implementation of the ALTO specifications, including the latest Internet Drafts that have been moved to Working Group Last Call to support modern transport protocols^[22] and OAM^[23]. As shown in Figure 3, the architecture maps the IETF ALTO server and ALTO client onto the OpenALTO software stack as follows:

- The OpenALTO Server stack includes three core building blocks: The Application-facing Interface, the Network-facing Interface, and the ALTO Services Layer.
- The Application-facing Interface provides an *Application Programming Interface* (API) that applications can query to retrieve the state of the network.

- The Network-facing Interface implements a variety of network plugins to support the retrieval of network state information; each plugin supports a different type of network. To facilitate the development of OpenALTO, this interface also includes plugins for simulation and emulation environments such as *Mininet*.
- The ALTO Services Layer provides the core ALTO functions by implementing [1]. This layer currently includes the Network Map, the Cost Map, and the Property Map services.
- The OpenALTO client is a thin layer that implements the HTTP-based client-side protocol described in [1] and [9]. The ALTO client is installed as a library in the same device in which the application is being run, and the application uses it to retrieve the network state from the ALTO server.

Figure 3: Mapping of RFC 7285 Entities onto the OpenALTO Software Architecture



Future Perspectives

The ALTO protocol initially started with the goal of supporting the optimization of P2P applications in 2008, then evolved to incorporate extensions for the support of CDNs in 2014, and today it is well-positioned to support the requirements of new advanced edge computing applications such as augmented reality, vehicle networks, and the metaverse, among others. Because this new class of applications requires stringent *Quality of Experience* (QoE) performance, the ALTO protocol becomes a key component to enable collaborative application/network schemes.

Specifically, ALTO contributes to the optimization of service placement and selection decisions based on the communication properties of the network. In this regard, and as its current charter is being finalized, proposals are being made to extend the protocol towards supporting edge computing applications in three possible areas: (1) extending ALTO metrics to include information about the compute resources (for example, *Central Processing Unit* [CPU], *Graphics Processing Unit* [GPU], memory, and storage) found in the distributed edge computing network, (2) incorporating protocol semantics for the sharing of state between ALTO servers in multi-domain networking environments, helping applications gain a global end-to-end view of the network, and (3) potentially incorporating information about the level of trust offered by each ANE along a communication path to improve the security of new advanced applications such as the metaverse.

Beyond the IETF, several other *Standards Development Organizations* (SDOs) such as the *3rd Generation Partnership Project* (3GPP) are also investigating solutions for exposing network capabilities to enable the optimization of new advanced applications. These solutions can naturally take advantage of ALTO, and there is the potential for IETF technology to become an important enabler of Internet capabilities demanded by developments arising in other SDOs. To enable these cross-SDO synergies, the ALTO protocol needs to be further socialized inside and outside the IETF with a focus on illustrating how it can provide the intended exposure features.

Future directions of the ALTO protocol are currently being discussed in the WG mailing list (<https://mailarchive.ietf.org/arch/browse/alto/>). The WG welcomes your participation to help identify the key priorities towards supporting the newly arising edge computing applications.

References and Further Reading

- [1] Richard Alimi, Ed., Reinaldo Penno, Ed., Richard Yang, Ed., Sebastian Kiesel, Stefano Previdi, Wendy Roome, Stanislav Shalunov, and Richard Woundy, “Application-Layer Traffic Optimization (ALTO) Protocol,” RFC 7285, September 2014.
- [2] OpenALTO Project, available at <https://github.com/openalto>
- [3] Mohamed Boucadair and Qin Wu, “A Cost Mode Registry for the Application-Layer Traffic Optimization (ALTO) Protocol,” RFC 9274, July 2022.
- [4] Sebastian Kiesel, Stefano Previdi, Martin Stiernerling, Richard Woundy, and Richard Yang, “Application-Layer Traffic Optimization (ALTO) Requirements,” RFC 6708, September 2012.
- [5] Sebastian Kiesel, Martin Stiernerling, Nico Schwan, Michael Scharf, and Haibin Song, “Application-Layer Traffic Optimization (ALTO) Server Discovery,” RFC 7286, November 2014
- [6] ODL ALTO, available at:
<https://wiki.opendaylight.org/display/ODL/ALTO>

- [7] Sabine Randriamasy, Wendy Roome, and Nico Schwan, “Multi-Cost Application-Layer Traffic Optimization (ALTO),” RFC 8189, October 2017.
- [8] Sabine Randriamasy, Richard Yang, Qin Wu, Lingli Deng, and Nico Schwan, “Application-Layer Traffic Optimization (ALTO) Cost Calendar,” RFC 8896, November 2020.
- [9] Wendy Roome and Richard Yang, “Application-Layer Traffic Optimization (ALTO) Incremental Updates Using Server-Sent Events (SSE),” RFC 8895, November 2020.
- [10] Wendy Roome, Sabine Randriamasy, Richard Yang, Jingxuan Zhang, and Kai Gao, “An Extension for Application-Layer Traffic Optimization (ALTO): Entity Property Maps,” RFC 9240, July 2022.
- [11] Jan Seedorf and Eric Burger, “Application-Layer Traffic Optimization (ALTO) Problem Statement,” RFC 5693, October 2009.
- [12] Jan Seedorf, Richard Yang, Kevin Ma, Jon Peterson, and Jingxuan Zhang, “Content Delivery Network Interconnection (CDNI) Footprint and Capabilities Advertisement Using Application-Layer Traffic Optimization (ALTO),” RFC 9241, July 2022.
- [13] Martin Stiernerling, Sebastian Kiesel, Michael Scharf, Hans Seidel, and Stefano Previdi, “Application-Layer Traffic Optimization (ALTO) Deployment Considerations,” RFC 7971, October 2016.
- [14] Qin Wu, Richard Yang, Young Lee, Dhruv Dhody, Sabine Randriamasy, and Luis Contreras, “ALTO Performance Cost Metrics,” Internet Draft, Work in Progress, **draft-ietf-alto-performance-metrics-28**, March 2022.
- [15] Richard Barnes, “Use Cases and Requirements for JSON Object Signing and Encryption (JOSE),” RFC 7165, April 2014.
- [16] Choongul Park, Yeongil Seo, Kun-youll Park, and Youngseok Lee, “The concept and realization of context-based content delivery of NGSON,” in *IEEE Communications Magazine*, Volume 50, No. 1, pp. 74–81, January 2012.
- [17] Marcus Schöller, Martin Stiernerling, Andreas Ripke, and Roland Bless, “Resilient deployment of virtual network functions,” 2013 *5th International Congress on Ultra Modern Telecommunications and Control Systems and Workshops (ICUMT)*, Almaty, Kazakhstan, 2013.
- [18] Enric Pujol, Ingmar Poese, Johannes Zerwas, Georgios Smaragdakis, and Anja Feldmann, “Steering Hyper-Giants’ Traffic at Scale,” In *Proceedings of the 15th International Conference on Emerging Networking Experiments and Technologies (CoNEXT ’19)*. Association for Computing Machinery, New York, 2019.

- [19] Michael Scharf, Thomas Voith, Wendy Roome, Bob Gaglianella, Moritz Steiner, Volker Hilt, and Vijay K. Gurbani, “Monitoring and abstraction for networked clouds,” *16th International Conference on Intelligence in Next Generation Networks*, Berlin, Germany, 2012.
- [20] Yuhang Jia, Yunfei Zhang, Richard Yang, Gang Li, Yixue Lei, Yunbo Han, and Sabine Randriamasy, “MoWIE for Network Aware Application,” Internet Draft, Work in Progress, **draft-huang-alto-mowie-for-network-aware-app-05**, November 2022.
- [21] Jingxuan Zhang et al., “Sextant: Enabling Automated Network-aware Application Optimization in Carrier Networks,” *2021 IFIP/IEEE International Symposium on Integrated Network Management (IM)*, Bordeaux, France, 2021.
- [22] Roland Schott, Richard Yang, Kai Gao, Lauren Delwiche, and Lachlan Keller, “The ALTO Transport Information Publication Service,” Internet Draft, Work in Progress, **draft-ietf-alto-new-transport-07**, April 2023.
- [23] Jingxuan Zhang, Dhruv Dhody, Kai Gao, Roland Schott, and Q. Ma, “YANG Data Models for the Application-Layer Traffic Optimization (ALTO) Protocol,” Internet Draft, Work in Progress, **draft-ietf-alto-oam-yang-06**, April 2023.
- [24] ALTO IETF Wiki on Collecting Deployment Experiences, available at: <https://wiki.ietf.org/en/group/ALTO/deployment>
- [25] Sebastian Kiesel and Martin Stiernerling, “Application-Layer Traffic Optimization (ALTO) Cross-Domain Server Discovery,” RFC 8686, February 2020.
- [26] Jan Seedorf, Jon Peterson, Stefano Previdi, Ray van Brandenburg, and Kevin Ma, “Content Delivery Network Interconnection (CDNI) Request Routing: Footprint and Capabilities Semantics,” RFC 8008, December 2016.
- [27] Jacob Dunefsky, Mahdi Soleimani, Ryan Yang, Jordi Ros-Giralt, Mario Lassnig, Inder Monga, Frank K. Würthwein, Jingxuan Zhang, Kai Gao, and Y. Richard Yang, “Transport control networking: optimizing efficiency and control of data transport for data-intensive networks,” *ACM SIGCOMM*, August 2022.

QIN WU is an expert on Network Management Architecture with Huawei's Data Communication in China. He is also responsible for enterprise networking innovation and standards work such as IoT, Security, SD-WAN, and Edge Computing. Involved in strategic standards development, engaging with some related open-source projects such as FD.io and ONAP for more than 10 years, Qin has held various positions in IETF, ITU-T, and CCSA. He has over 16 years of experience on network architecture and protocol design, starting from mobility management, performance measurement, IPTV to SDN, NFV, network management automation and YANG, telemetry, AIOPs, etc. Currently he focuses on promoting digital twin networking and Network and Application collaboration. He used to chair the IETF L3SM and L2SM working groups in the OPS area; he currently chairs the ALTO Working Group in the Transport area, serves as a member of the OPS-DIR Directorate, and has coauthored more than 52 RFCs spanning six IETF areas (OPS, SEC, RTG, TSV, RAI, and INT). He received his PhD degree of Control Theory and Engineering from Nanjing University of Science and Technologies. Qin is a member of the Internet Architecture Board.

E-mail: bill.wu@huawei.com

MOHAMED BOUCADAI is a Senior Network Architect within the "Network of the Future" team in Orange Innovation. He worked at the Orange corporate division, where he was responsible for making recommendations on the evolution of IP/MPLS core networks. Mohamed is the author/editor of several books, including *Design Innovation and Network Architecture for the Future Internet* (ISBN: 9781799876465), *Emerging Automation Techniques for the Future Internet* (ISBN: 9781522571469), *Redesigning the Future of Internet Architectures* (ISBN: 9781466683716), *Solutions for Sustaining Scalability in Internet Growth* (ISBN: 978-1466643055), *IP Telephony Interconnection Reference: Challenges, Models and Engineering* (ISBN: 978-1439851784), *Recent Advances in Providing QoS and Reliability in Future Internet Backbone* (ISBN: 978-1617618581), *Inter-Asterisk Exchange (IAX) Deployment Scenarios in SIP-Enabled Networks* (ISBN: 978-0470770726), and *Service Automation and Dynamic Provisioning Techniques in IP/MPLS Environments* (ISBN: 978-0470018291).

E-mail: mohamed.boucadair@orange.com

JORDI ROS-GIRALT is a Director of Engineering at Qualcomm Europe, Inc., where he leads the high-performance networking team as part of AI Research focusing on the area of accelerating application performance for 5G, 6G, and the Edge Cloud. Jordi has published upwards of 75 articles in scientific conferences and journals, and is the inventor and developer of several communication network algorithms and technologies, most of which have been included in commercial products. Jordi received his PhD in Computer Science, an MBA from the University of California, and a BSc in Telecommunications Engineering from the Barcelona Tech University (UPC).

E-mail: jros@qti.qualcomm.com

The Internet Protocol Journal is published under the "CC BY-NC-ND" Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an "as-is" basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

The Impact of Randomized Layer-2 Addresses on Privacy and Applications

by Carlos J. Bernardos, Juan-Carlos Zuñiga, Jerome Henry, and, Alain Mourad

Wi-Fi technology has revolutionized communication and become the preferred technology and sometimes the only networking technology used by devices such as smartphones, tablets, and *Internet-of Thing* (IoT) devices.

On the other hand, Internet privacy is becoming a huge concern, as more and more mobile devices are connecting to the Internet. This ubiquitous connectivity, together with not-very-secure protocol stacks and the lack of proper education about privacy, make it very easy to track/monitor the location of users and/or eavesdrop on their physical and online activities. The cause of this situation has many factors, such as the vast digital footprint that users leave on the Internet; for example: information sharing on social networks, the cookies that browsers and servers use to provide a better navigation experience, connectivity logs that allow tracking of a user's Layer 2 *Media Access Control* (L2 MAC) or Layer 3 (L3) address, web trackers, etc., and/or the weak (or even null in some cases) authentication and encryption mechanisms used to secure communications.

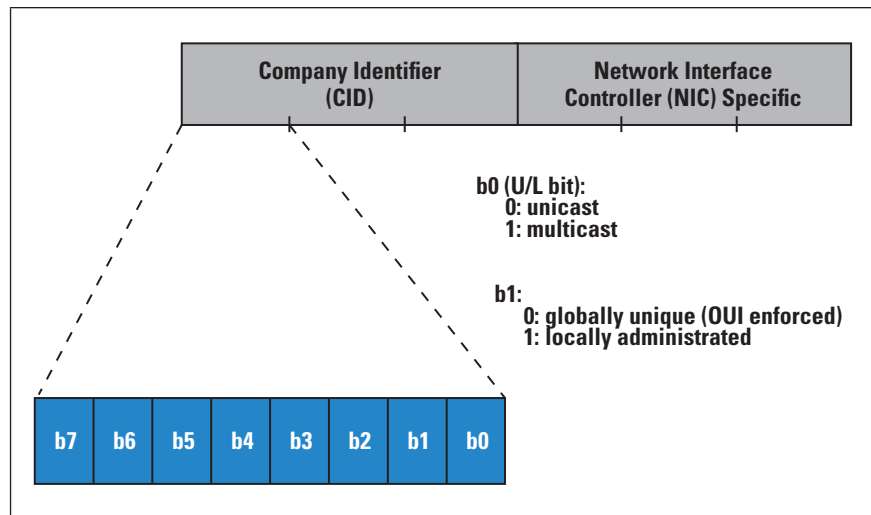
This privacy concern affects all layers of the protocol stack, from the lower layers involved in the actual access to the network (for example, you can use the L2 and L3 addresses to obtain a user's location) to higher-layer protocol identifiers and user applications^[1]. In particular, IEEE 802 MAC addresses have historically been an easy target for tracking users^[2]. Attackers who are equipped with surveillance equipment can “monitor” Wi-Fi packets and track the activity of Wi-Fi devices. After the association between a device and its user is made, identifying the device and its activity is sufficient to deduce information about what the user is doing, without the user's consent.

IEEE 802.11 (Wi-Fi) interfaces, as any other kind of IEEE 802-based network interface, like Ethernet (that is, IEEE 802.3), have a Layer 2 address, also referred to as the MAC address, which anybody who can receive the signal transmitted by the network interface can see. Figure 1 shows the format of these addresses.

A third party, such as a passive device listening to communications in the same network, can easily observe MAC addresses. In an 802.11 network, a station exposes its MAC address in two different situations:

- While unassociated and actively scanning for available networks, the MAC address is used in the *Probe Request* frames that the device sends (aka IEEE 802.11 STA).
- After it is associated to a given *Access Point* (AP), the MAC address is used in frame transmission and reception, as one of the addresses used in the address fields of an IEEE 802.11 frame.

Figure 1: IEEE 802 MAC Address Format



MAC addresses can be either universally or locally administered. A MAC address is identified as being locally administered when the second-last significant bit of the most significant octet of the address (the U/L bit) is set. The MAC address is identified as globally unique when the U/L bit is unset.

A universally administered address is uniquely assigned to a device by its manufacturer (and is called the *burned-in address*). Most physical devices are provided with a universally administered address, which is composed of two parts: (i) the *Company Identifier* (CID), which is the first three octets in transmission order, identified by the organization that issued it, and (ii) the *Network Interface Controller* (NIC)-specific address, which is the following three octets, assigned by the organization that manufactured the controller, in such a way that the resulting MAC address is globally unique. Since universally administered MAC addresses are by definition globally unique, when a device uses this MAC address to transmit data—especially over the air—it is relatively easy to track this device by simple medium observation. This possibility poses a privacy concern^[3] when the device is directly associated to a single user (for example, smartphones, etc.).

Locally administered addresses can override the burned-in address, and they can be set up by either the network administrator or the *Operating System* (OS) of the device to which the address pertains. This feature allows you to generate local addresses without the need for any global coordination mechanism to ensure that the generated address is still unique within the local network. You can use this feature to generate random addresses, which decouple the globally unique identifier from the device and thereby make it more difficult to track a user device from its MAC/L2 address^[4]. There are initiatives at the IEEE 802 and other organizations to specify ways in which these locally administered addresses should be assigned, depending on the use case.

To reduce the risks of correlation between a device activity and its owner, multiple vendors have started to implement *Randomized and Changing* MAC (RCM) addresses. With this scheme, an end device implements a different RCM over time when exchanging traffic over a wireless network. By randomizing the MAC address, the persistent association between a given traffic flow and a single device is made more difficult, assuming no other visible unique identifiers are in use.

However, such address changes may affect the user experience and the efficiency of legitimate network operations. For a long time, network designers and implementers relied on the assumption that a given machine in a network implementing IEEE 802 technologies would be represented by a unique network MAC address that would not change over time, despite the existence of tools to flush out the MAC address to bypass some network policies. When this assumption is broken, elements of network communication may also break.

For example, sessions established between the end device and network services may be lost and packets in translation may suddenly be without a clear source or destination. If multiple clients implement fast-paced RCM rotations, network services may be over-solicited by a small number of stations that appear as many clients.

At the same time, some network services rely on the client station providing an identifier, which can be the MAC address or another value. If the client implements MAC rotation but continues sending the same static identifier, then the association between a stable identifier and the station continues despite the RCM scheme. There may be environments where such continued association is desirable, but others where the user privacy has more value than any continuity of network service state.

Application and Network Scenarios That RCM Can Affect

Device identity is important in scenarios where the network needs to know the device or user identity in order to offer, operate, and maintain certain valid services. Currently, many use cases and applications make an implicit assumption that a device is represented by an IEEE 802 L2 permanent and unique MAC address. This assumption is being used in both control- and data-plane functions and protocols. RCM breaks this assumption. This paradigm shift requires updating applications to function across MAC address changes.

When a device changes its MAC address, other devices on the same LAN may fail to recognize that the same machine is attempting to communicate with them. Additionally, multiple layers implemented at upper layers have been designed with the assumption that each node on the LAN, using these services, would have a MAC address that would stay the same over time (a *persistent* MAC address).

This assumption sometimes adds to the *Personally Identifiable Information* (PII) confusion, for example in the case of *Authentication, Association, and Accounting* (AAA) services authenticating the user of a machine and associating the authenticated user to the device MAC address. Other services focus solely on the machine, for example, the *Dynamic Host Configuration Protocol* (DHCP), but still expect each machine to use a persistent MAC address, for example to re-assign the same IP address to a returning device. Changing the MAC address may disrupt these services and the user experience.

The impact of using a persistent or a randomized and changing MAC address very much depends on the environment where the device operates (that is, the use case), on the presence and nature of other devices in the environment, and on the type of network the device is communicating through. Therefore, a device can use a MAC address that can persist over time if trust with the environment is established, or that can be temporal if that address is going to be used as an identity for a service in an environment where trust has not been established. Note that this trust is not binary, and it ranges from: (i) full trust: environments where a personal device establishes a trust relationship and can share a persistent device identity with the access network devices, without the fear of that identity being shared beyond the L2 broadcast domain; (ii) selective trust: environments where the device may not be willing to share a persistent identity with some elements of the Layer 2 broadcast domain but may be willing to do it with other elements; and (iii) zero trust: environments where the device may not be willing to share any persistent identity with any local entity reachable through the AP and may express a temporal identity to each of them.

This trust relationship naturally depends on the relationship between the user of the personal device and the operator of the service. Thus, it is useful to enumerate some scenarios (which can be easily translated into use cases) where the use of RCM might have an impact:

- *Residential settings under the control of the user*: this case is typical of a home network with Wi-Fi in the LAN and an Internet connection. In this environment, traffic over the Internet does not expose the MAC address if it is not copied to another field before routing happens. The user controls the wire segment within the broadcast domain, and this segment, therefore, is usually not at risk of hosting an eavesdropper. Full trust is typically established at this level among users and with the network elements. The device trusts the AP and all Layer 2 domain entities beyond the AP. However, unless the user has full access control over the physical space where the Wi-Fi transmissions can be detected, there is no guarantee that an eavesdropper would not be observing the communications. As such, it is common to assume that, even in this environment, full trust cannot be achieved.

- *Managed residential settings*: examples of this type of environment include shared living facilities and other collective environments where an operator manages the network for the residents. The over-the-air exposure is similar to that of a home. A number of devices larger than in a standard home may be present, and the operator may be requested to provide IT support to the residents. Therefore, the operator may need to identify the activity of a device in real time, but may also need to analyze logs so as to understand a past reported issue. For both activities, a device identification associated with the session is needed. Full trust is often established in this environment, at the scale of a series of a few sessions, not because it is assumed that no eavesdropper would observe the network activity, but because it is a common condition for the managed operations.
- *Public guest networks*: public hotspots, such as in shopping malls, hotels, stores, train stations, and airports, are typical of this environment. The guest network operator may be legally mandated to identify devices or users or may have the option to leave all devices and users untracked. In this environment, trust is commonly not established with any element of the Layer 2 broadcast domain (zero trust model by default).
- *Enterprises (with Bring Your Own Device [BYOD])*: users may be provided with corporate devices or may bring their own. The devices are not directly under the control of a corporate IT team. Trust may be established as the device joins the network. Some enterprise models mandate full trust; others, considering the BYOD nature of the device, allow selective trust.
- *Managed enterprises*: in this environment, users are typically provided with corporate devices, and all connected devices are managed, for example through a *Mobile Device Management* (MDM) profile installed on the device. Full trust is created as the MDM profile is installed.

Ongoing Efforts/Approaches Regarding RCM

Practical experiences of RCM in live devices helped researchers fine-tune their understanding of attacks against randomization mechanisms^[5]. At IEEE 802.11 these research experiences eventually formed the basis for a specified mechanism introduced in the IEEE 802.11aq in 2018, which recommends mechanisms to avoid pitfalls when using randomized MAC addresses^[6].

More recent developments include turning on MAC randomization in mobile operating systems by default, which affects the ability of network operators to personalize or customize services^[7]. Therefore, follow-on work in the IEEE 802.11 mapped effects of a potentially large uptake of randomized MAC identifiers on many commonly offered operator services in 2019^[8].

In the summer of 2020, this work resulted in two new standards projects with the purpose of developing mechanisms that do not decrease user privacy and enable an optimal user experience when the MAC address of a device in an *Extended Service Set* is randomized or changes^[9] and user privacy solutions applicable to IEEE Std 802.11^[10].

The IEEE 802.1 Working Group has also published a specification that defines a local MAC address space structure, known as the *Structured Local Address Plan* (SLAP). This structure designates a range of local MAC addresses for protocols using a CID assigned by the IEEE Registration Authority. Another range of local MAC addresses is designated for assignment by administrators. The specification recommends a range of local MAC addresses for use by IEEE 802 protocols^[11].

Work within the IEEE 802.1 Security Task Group on privacy recommendations for all IEEE 802 network technologies has also looked into general recommendations on identifiers, reaching the conclusion that temporary and transient identifiers are preferable in network technology designs if there are no compelling reasons of service quality for a newly introduced identifier to be permanent. This work has been specified in the recently published IEEE P802E: “Recommended Practice for Privacy Considerations for IEEE 802 Technologies”^[12]. The IEEE P802E specification will form part of the basis for the review of user privacy solutions applicable to IEEE Std 802.11 (aka Wi-Fi) devices as part of the RCM^[7] efforts.

Currently, two task groups in IEEE 802.11 are addressing issues related to RCM:

- The IEEE 802.11bh Task Group, looking at mitigating the repercussions that RCM creates on 802.11 networks and related services, and
- The IEEE 802.11bi Task Group, which will define modifications to the IEEE Std 802.11 MAC specification to specify new mechanisms that address and improve user privacy.

At the *Wireless Broadband Alliance* (WBA), the *Testing and Interoperability Working Group* has been looking at the issues related to MAC address randomization and has identified a list of potential impacts of these changes to existing systems and solutions, mainly related to Wi-Fi identification. As part of this work, WBA has documented a set of use cases that a Wi-Fi Identification Standard should address in order to scale and achieve longer-term sustainability of deployed services. A first version of this document has been liaised with the IETF as part of the *MAC Address Device Identification for Network and Application Services* (MADINAS) activities through the “Wi-Fi Identification in a post MAC Randomization Era v1.0” paper^[13].

Several IP address assignment mechanisms such as the IPv6 *Stateless Address Auto-Configuration* (SLAAC) techniques^[14] generate the *Interface Identifier* (IID) of the address from its MAC address (via EUI64), which then becomes visible to all IPv6 communication peers. This feature potentially allows for global tracking of a device at L3 from any point on the Internet. Besides, the prefix part of the address provides meaningful insights into the physical location of the device in general, which together with the MAC address-based IID, makes it easier to perform global device tracking.

Some solutions might mitigate this privacy threat, such as the use of temporary addresses^[15] and opaque IIDs^[16,17]. Additionally, [18] proposes an extension to DHCPv6 that allows a scalable approach to link-layer address assignments where preassigned link-layer address assignments (such as by a manufacturer) are not possible or unnecessary. [19] proposes extensions to DHCPv6 protocols to enable a DHCPv6 client or relay to indicate a preferred SLAP quadrant to the server, so that the server may allocate MAC addresses in the quadrant requested by the client or relay.

Not only can you use MAC and IP addresses for tracking purposes, but some DHCP options allow you to also carry unique identifiers. These identifiers can enable device tracking even if the device administrator takes care of randomizing other potential identifications such as link-layer addresses or IPv6 addresses. [20] introduces anonymity profiles, designed for clients that wish to remain anonymous to the visited network. The profiles provide guidelines on the composition of DHCP or DHCPv6 messages, designed to minimize disclosure of identifying information.

Existing Solutions

One possible solution is to use 802.1X with *Wi-Fi Protected Access 2/3* (WPA2/WPA3). At the time of association to a Wi-Fi access point, 802.1X authentication coupled with WPA2 or WPA3 encryption schemes allows for the mutual identification of the client device or of the user of the device and an authentication authority. The authentication exchange is protected from eavesdropping. In this scenario, you can obfuscate the identity of the user or the device from external observers. However, the authentication authority is in most cases under the control of the same entity as the network access provider, thus making the identity of the user or device visible to the network owner. This scheme is therefore well-adapted to enterprise environments, where a level of trust is established between the user and the enterprise network operator.

A different approach is the Wireless Broadband Alliance *OpenRoaming* standard, which introduces an intermediate trusted relay between local venues and sources of identity. The federation structure also extends the type of authorities that can be used as identity sources (compared to the traditional enterprise-based 802.1X scheme for Wi-Fi), and facilitates the establishment of trust between a local venue and an identity provider.

Such a procedure dramatically increases the likelihood that one or more identity profiles for the user or the device will be accepted by a local venue. At the same time, authentication does not occur to the local venue, thus offering the possibility for the user or device to keep their identity obfuscated from the local network operator, unless that operator specifically expresses the requirement to disclose such identity (in which case the user has the option to accept or decline the connection and associated identity exposure). The OpenRoaming scheme therefore seems well-adapted to public Wi-Fi and hospitality environments, allowing for the obfuscation of the identity from unauthorized entities, while also permitting mutual authentication between the device or the user and a trusted identity provider.

It is also worth mentioning that most evolved client device OSes already offer RCM schemes, enabled by default (or easy to enable) on client devices. With these schemes, the device can change its MAC address, when not associated, after using a given MAC address for a semi-random duration window. These schemes also allow for the device to manifest a different MAC address in different *Service Set Identifiers* (SSIDs). Different OSes follow slightly different approaches, which are also evolving with the new releases. Such a randomization scheme enables the device to limit the duration of exposure of a single MAC address to observers. In the IEEE 802.11-2020 specification, MAC address rotation is not allowed during a given association session, and thus rotation of MAC address can occur only through disconnection and reconnection.

Ongoing Work in the IETF

The MADINAS Working Group in the IETF is addressing documentation of the services that may be affected by RCM and evaluation of possible solutions to maintain the quality of user experience and network efficiency in the presence of RCM, while user privacy is reinforced.

The group will generate documents regarding the state of affairs of RCM, and a *Best Current Practices* (BCP) document recommending a means to reduce the impact of RCM on the documented use cases while ensuring that the privacy achieved with RCM is not compromised. For scenarios where device identity stability is desirable, the BCP document will recommend existing protocols that you can use to protect the request and exchange of identifiers between the client and the service provider.

The MADINAS Working Group is focused on coordination with other IETF Working Groups (for example, DHC and IntArea). In addition, it actively liaises with other relevant organizations, such as IEEE 802 and the Wireless Broadband Alliance. The objective is to coordinate on the different recommendations, as well as planning potential follow-up activities within or outside the IETF.

It is expected that the outcome from the coordinated efforts among these standards organizations will enable the use of RCM in the different scenarios previously analyzed, providing both privacy and the operational characteristics about the quality of user experience and network efficiency that each one of the scenarios requires.

References

- [1] Carlos J. Bernardos, Juan Carlos Zúñiga, and Piers O’Hanlon, “Wi-Fi Internet Connectivity and Privacy: Hiding your tracks on the wireless Internet,” *IEEE Conference on Standards for Communications and Networking (CSCN)*, October 2015.
- [2] James Vincent, “London’s bins are tracking your smartphone,” *The Independent*, August 2013, <https://www.independent.co.uk/life-style/gadgets-and-tech/news/updated-london-s-bins-are-tracking-your-smartphone-8754924.html>
- [3] Piers O’Hanlon, Joss Wright, and Ian Brown, “Privacy at the link layer,” Contribution at W3C/IAB workshop on *Strengthening the Internet Against Pervasive Monitoring (STRINT)*, February 2014. <https://www.w3.org/2014/strint/papers/35.pdf>
- [4] Marco Gruteser and Dirk Grunwald, “Enhancing Location Privacy in Wireless LAN Through Disposable Interface Identifiers: A Quantitative Analysis,” *Mobile Networks and Applications*, Volume 10, No. 3, pp. 315-325, 2005.
- [5] Jeremy Martin, Travis Mayberry, Colin Donahue, Lucas Foppe, Lamont Brown, Chadwick Riggins, Eric C. Rye, and Dane Brown, “A Study of MAC Address Randomization in Mobile Devices and When It Fails,” arXiv:1703.02874v2 [cs.CR], 2017.
- [6] Institute of Electrical and Electronics Engineers (IEEE), “IEEE 802.11aq-2018 - IEEE Standard for Information Technology—Telecommunications and Information Exchange Between Systems Local and Metropolitan Area Networks—Specific Requirements Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications Amendment 5: Preassociation Discovery,” IEEE 802.11, 2018.
- [7] Institute of Electrical and Electronics Engineers (IEEE), “IEEE 802.11 Randomized and Changing MAC Addresses Study Group CSD on User Experience Mechanisms,” doc.:IEEE 802.11-20/1346r1, 2020.
- [8] Institute of Electrical and Electronics Engineers (IEEE), “IEEE 802.11 Randomized and Changing MAC Addresses Topic Interest Group Report,” doc.:IEEE 802.11-19/1442r9, 2019.
- [9] Institute of Electrical and Electronics Engineers (IEEE), “IEEE 802.11 Randomized and Changing MAC Addresses Study Group PAR on User Experience Mechanisms,” doc.:IEEE 802.11-20/742r5, 2020.

- [10] Institute of Electrical and Electronics Engineers (IEEE), “IEEE 802.11 Randomized and Changing MAC Addresses Study Group PAR on Privacy Mechanisms,” doc.:IEEE 802.11-19/854r7, 2020.
- [11] Institute of Electrical and Electronics Engineers (IEEE), “IEEE 802c-2017 - IEEE Standard for Local and Metropolitan Area Networks: Overview and Architecture—Amendment 2: Local Medium Access Control (MAC) Address Usage,” IEEE 802c, 2017.
- [12] Institute of Electrical and Electronics Engineers (IEEE), “IEEE 802E-2020 - IEEE Recommended Practice for Privacy Considerations for IEEE 802 Technologies,” IEEE 802E, 2020.
- [13] Wide Band Alliance, “Wi-Fi Identification Scope for Liaising - In a post MAC Randomization Era,” doc.:WBA Wi-Fi ID Intro: Post MAC Randomization Era v1.0 - IETF liaison, March 2020.
- [14] Susan Thomson, Thomas Narten, and Tatuya Jinmei, “IPv6 Stateless Address Autoconfiguration,” RFC 4862, September 2007.
- [15] Richard Draves and Dave Thaler, “Default Router Preferences and More-Specific Routes,” RFC 4191, November 2005.
- [16] Fernando Gont, “A Method for Generating Semantically Opaque Interface Identifiers with IPv6 Stateless Address Autoconfiguration (SLAAC),” RFC 7217, April 2014.
- [17] Fernando Gont, Alissa Cooper, Dave Thaler, and Will Liu, “Deprecating EUI-64 Based IPv6 Addresses,” Internet Draft, Work in Progress, **draft-gont-6man-deprecate-eui64-based-addresses-00**, October 2013.
- [18] Bernie Volz, Tomek Mrugalski, and Carlos J. Bernardos, “Link-Layer Address Assignment Mechanism for DHCPv6,” RFC 8947, December 2020.
- [19] Carlos J. Bernardos and Alain Mourad, “Structured Local Address Plan (SLAP) Quadrant Selection Option for DHCPv6,” RFC 8948, December 2020.
- [20] Christian Huitema, Tomek Mrugalski, and Suresh Krishnan, “Anonymity Profiles for DHCP Clients,” RFC 7844, May 2016.

CARLOS J. BERNARDOS received a Telecommunication Engineering degree in 2003, and a PhD in Telematics in 2006, both from the University Carlos III of Madrid, where he works as a Full Professor. He teaches different undergraduate and masters degree courses, including the Master and Specialist in NFV and SDN. His research interests include IP mobility management, network virtualization, cloud computing, vehicular communications, and experimental evaluation of mobile wireless networks. He has published over 100 scientific papers in international journals and conferences. Carlos has been an active contributor to IETF since 2005, being co-author of more than 30 contributions and several standards, he has co-chaired the IETF P2PSIP and IPWAVE Working Groups, and he currently co-chairs the MADINAS Working Group and the Internet Area Directorate (INTDIR). E-mail: cjbc@it.uc3m.es

JUAN CARLOS ZÚÑIGA is a technology and business development trilingual leader with more than 20 years of international experience. He currently leads wireless standardisation and IP efforts within the Global Technical Standards team at Cisco Systems. He has extensive experience with heterogeneous Radio Access Networks and 4G/5G Core network cloud technologies. Juan is an active thought leader and contributor in different standards and industrial fora, such as IETF, IEEE 802, ISOC, ETSI, 3GPP SA/CT/RAN, ITU, W3C, and GS1. He has strong experience developing, managing, analyzing, and developing international intellectual property (IPR) patent portfolios. He has been granted patents for more than 70 USPTO/EPO inventions and over 100 published applications. Juan holds several recognitions and awards. E-mail: juzuniga@cisco.com

JEROME HENRY is a Principal Engineer in the Enterprise Infrastructure and Solutions Group at Cisco Systems. Jerome joined Cisco in 2012. Before that time, he consulted and taught about heterogeneous networks and wireless integration with the European Airespace team, which Cisco later acquired to become its main wireless solution. Jerome holds more than 150 patents, is a member of the IEEE, where he was elevated to Senior Member in 2013, and also represents Cisco in multiple Wi-Fi Alliance working groups. With more than 10,000 hours in the classroom, Jerome was awarded the IT Training Award Best Instructor silver medal. He is based in Research Triangle Park, North Carolina. E-mail: jerhenry@cisco.com

ALAIN MOURAD is an award-winning innovator with over 20 years of experience in the Research and Development of wireless technologies spanning four generations of cellular systems (3G/4G/5G and heading towards 6G). He currently heads the Future Wireless Europe Research and Innovation Lab at InterDigital in London (UK). Prior to InterDigital, Alain was a Principal Engineer at Samsung Electronics Research and Development and a Senior Engineer at Mitsubishi Electric R&D Centre Europe, where he was active in the specification of wireless standards (3GPP, IEEE802, DVB, and ATSC). E-mail: Alain.Mourad@interdigital.com

Check your Subscription Details!

If you have a print subscription to this journal, you will find an expiration date printed on the back cover. For several years, we have “auto-renewed” your subscription, but now we ask you to log in to our subscription system and perform this simple task yourself. Make sure that both your postal and e-mail addresses are up-to-date since these are the only methods by which we can contact you. If you see the words “Invalid E-mail” on your copy this means that we have been unable to contact you through the e-mail address on file. If this is the case, please contact us at ipj@protocoljournal.org with your new information. The subscription portal is located here: <https://www.ipjsubscription.org/>

Twenty-Five Years Later

by Geoff Huston, APNIC

The Internet is not quite as young and spritely as you might have thought. Apple's iPhone, released in 2007, is now 16 years old, and YouTube is an ageing teenager at 18 after its initial release in 2005, and these two examples are relatively recent additions to the Internet. The first web browser, Mosaic, was released some 30 years ago in 1993. Going back further, the Internet emerged from its early *Advanced Research Projects Agency* (ARPA) roots in the form of the *National Science Foundation Network* (NSFNET) in 1986. At the start of 1983, the *ARPA Network* (ARPANET) had a flag day and switched over to use the *Transmission Control Protocol* (TCP). Going back further, in 1974 Vint Cerf and Bob Kahn published the first academic paper describing the protocol and the underlying architectural framework of a packet-switched network that became the Internet. This achievement was built upon earlier foundations, where numerous efforts in the late 1960s showed the viability of a *packet-switched* approach to computer networking. These packet-switched networking efforts included a program led by Donald Davies at the UK National Physics Laboratory, an effort in the US in the form of an ARPA project led by Larry Roberts, and Louis Pouzin's work in France with the CYCLADES network. This work, in turn, has some of its antecedents in work by Paul Baran at the RAND Corporation on distributed communications and packet-switched networks, published between 1960 and 1964. The Internet has managed to accumulate a relatively lengthy pedigree.

And it has been a wild ride. The Internet has undergone numerous cycles of economic boom and bust, each of which is right up there with history's finest episodes of exuberant irrational mania. It has managed to trigger a comprehensive restructuring of the entire global communications enterprise and generated a set of changes that have already altered the way in which we now work and play. That's quite a set of achievements in just 25 years!

We should start this exploration of our past some 25 years ago in 1998 at the time of publication of the first edition of the *Internet Protocol Journal*. At that time, any lingering doubts about the ultimate success of the Internet as a global communications medium had been thoroughly dispelled. The Internet was no longer just a research experiment, or an intermediate way stop on the road to adoption of the *Open Systems Interconnection* (OSI) framework. By 1998 there was nothing else left standing in the data communications landscape that could serve our emerging needs for data communications. The *Internet Protocol* (IP) was now the communications technology for the day, if not for the coming century, and the industry message at the time was a clear one that said: "adopt the Internet into every product and service or imperil your entire future in this business."

No longer did the traditional telecommunications enterprises view the Internet with some polite amusement or even overt derision. It was now time for a desperate scramble to be part of this revolution in one of the world's major activity sectors. The largest enterprises in this sector, the old-world ex-monopoly telcos, had been caught wrong-footed in one of the biggest changes of the industry for many decades, and this time the concurrent wave of deregulation and competition meant that the future of the communications industry was being handed over to a small clique of Internet players.

By the early 2000s, the Internet had finally made it into the big time. The job was apparently done, and the Internet had prevailed. But then came a new revolution, this time in mobility services, where after numerous clumsy initial efforts by others, the iPhone entered the market with a seamless blend of sleek design and astounding capability. The mobile carriage sector struggled to keep up with the new levels of rapacious demand for Internet-based mobile data. The Internet then took on the television networks, replacing the incumbent broadcast and cable systems with streaming video. But the story is not over by any means. Communications continues to drive our world, and the Internet continues to evolve and change.

The evolutionary path of any technology can often take strange and unanticipated turns and twists. At some points simplicity and minimalism can be replaced by complexity and ornamentation, while at other times a dramatic cut-through exposes the core concepts of the technology and removes layers of superfluous additions. The technical evolution of the Internet appears to be no exception, and this story contains these same forms of unanticipated turns and twists.

Rather than offer a set of unordered observations about the various changes and developments over the past 25 years, I will use the traditional protocol stack model as a template, starting with the underlying transmission media, then looking at IP, the transport layer, then applications and services, and closing with a look at the business of the Internet.

Transmission

It seems like a totally alien concept these days, but the *Internet Service Provider* (ISP) business of 1998 was still based around the technology of dial-up modems. The state-of-the-art of modem speed had been continually refined, from 9600 bps to 14.4 kbps, to 28 kbps, to finally 56 kbps, squeezing every last bit out of the phase amplitude space contained in an analogue voice circuit. Analogue modems were capricious, constantly being superseded by the next technical refinement, unreliable, difficult for customers to use, and on top of that, they were slow! Almost everything else on the Internet had to be tailored to download reasonably quickly over a modem connection. Web pages were carefully composed with compressed images to ensure a rapid download, and plain text was the dominant medium as a consequence. It could only get better.

The evolution of access networks was initially one that exposed the inner digital core of the network out to the edges. The first approach was *Integrated Services Digital Network* (ISDN), where the underlying digitised voice circuit was drawn out to the network edge. At 64 kbps, this level of improvement was inadequate, and the next major step was to use *Digital Subscriber Line* (DSL) technology. DSL used the last mile of the network for an analogue channel, but instead of running a single low-speed bearer signal, DSL layered a large collection of individual bearer signals into the single circuit, performing a form of frequency division multiplexing on the basic analogue circuit in a trellis framework. DSL relied on the combination of the telephone company's efforts to operate the copper access circuits within a base level of signal quality and noise suppression, and the modem industry's continual incremental improvements in digital-signal-processing capability. Surprisingly, DSL achieved speeds of tens of megabits per second through these legacy copper access networks. However, DSL was largely an interim holding position while the search for a viable business model that could underwrite the costs of deployment and use of an open fibre-based access networks was underway.

The transition into fixed-wire access networks based on fibre-optic cable continues. The challenge is not in finding a suitable technology for fibre, but one of finding a suitable business model than can sustain the necessary capital investment in replacing the existing copper-based infrastructure. Some national communities used a model of a public-sector program, such as the *National Broadband Network* program in Australia, while others have remained as dedicated private-sector activities, and others have taken a hybrid approach with some level of local public-sector incentives being added into a private-sector program. The issue here is that fixed wire residential access networks do not offer compelling investment opportunities in most cases, with the high initial capital costs and the generally inadequate levels of take-up across the dwellings passed by the access infrastructure acting as disinhibitory factors. It is often the case that a residential community cannot support multiple access network deployments, bringing up the related issue of local access monopolies and the challenge of permitting some level of competitive access across a single physical access network. Nevertheless, fibre access rollouts continue across many parts of the world, and the transition of the wired copper network into a fibre access network capable of sustaining hundreds of megabits per connection is still progressing, seemingly in spite of the financial barriers that exist in many scenarios.

Mobile

The mobile network has experienced a completely different evolution, and for many years now the mobile sector has been demand-driven. The first mobile data service networks, introduced in the 1980s, were little more than low-speed digital encoders working across a single voiceband circuit. These 1G networks typically delivered a 2.4-kbps data download capacity.

The next generation of mobile services, 2G, was used in the 1990s. It was still predominately a voice service, and while it could theoretically support data access at speeds of 100 kbps, this data-transfer rate was largely aspirational, and the mobile network was predominantly used by the *Short Message Service* (SMS) as an adjunct to voice. The intersection of the Internet and mobile services occurred with the introduction of 3G mobile services. The 3G architecture could push IP connectivity directly out to the handset, supporting data-transfer speeds of 1–2 Mbps. This network capability, coupled with a new generation of handsets, first with the BlackBerry in 2002 and then the iPhone in 2007, transformed the mobile service into a mass-market consumer service. The high margins available from the service captured the attention of the traditional voice industry, and the result was a large-scale opening up of radio spectrum to create an Internet access environment that quickly rivalled the wire-line access market in size, but totally outpaced it in terms of revenue margins. This massive escalation of demand created further pressures on the capacity of the mobile system, and in 2009 the mobile sector introduced 4G services, opening up more spectrum bands, and also adding *Multiple-Input Multiple Output* (MIMO) to the mobile device to achieve greater deliverable capacity to each connected device. Over time these services were to deliver peak download speeds of 50 to 100 Mbps. The industry was also selling hundreds of millions of mobile devices per year. 4G dispensed with circuit-switched services, and it exclusively used packet switching. In 2018 5G was introduced. 5G can be supported over more spectrum bands, including a high-band millimetre spectrum at 24–47Ghz. These higher carrier frequencies permit multi-gigabit data services, but they come at a cost of higher density of base-station towers to compensate for the lower propagation distances.

Wi-Fi and Bluetooth

A second radio technology that has also transformed the Internet emerged in 1998, and it could be argued that it has become so fundamental that it has weaved itself so naturally into our environment that it all but disappeared. The combination of low-power radio systems and unlicensed radio spectrum allocation, or *Wi-Fi*, and subsequently *Bluetooth*, has been transformational. The combination of efficient battery technology, computer chips that operate with low power consumption, and the unwiring of the last few meters in the home and office completely changed our collective of technology, and it is only because of our desire to use products that are portable, unobtrusively wearable, and powerful enough to be useful that the component technologies such as batteries and processors have been pushed in this direction over this period. While large bands of radio spectrum space have been allocated to cellular mobile service operators, the intensity of use and the utility of use of radio spectrum peaks in the unlicensed spectrum space used by Wi-Fi and Bluetooth. It could be argued that the economic value of these unlicensed spectrum bands exceeds the exclusively licensed cellular radio systems by orders of magnitude. It could also be argued that the untethering of the last meter of the Internet transformed the Internet, and digital technologies in general, from a specialist pursuit into the consumer product space.

In the 1990s we described the effort to simplify the use of technology through the term “plug and play.” Wi-Fi was the critical technical development that made that term irrelevant by removing any need for the plug, or the socket for that matter!

Satellite

Mobile data services, Wi-Fi, and Bluetooth really revolutionised the Internet, taking it from a “destination you visit” to an “always-on utility in your pocket.” The Internet was now a set of applications on a device that went with you everywhere. Always available, always connected, no matter where you might be or what you might be doing. But that was not exactly the full truth. Head out into remote country far enough, or head onto the world’s oceans, and your connection options quickly disappeared, leaving only somewhat expensive satellite-based services.

These satellite services have been in operation since the early 1960s, but the high launch costs, limited capacity, and competing interests of terrestrial providers have meant that these services were often operated at the margins of viability. The best example is Motorola’s *Iridium* project of the late 1990s, where even before the entire service constellation of satellites was launched, the \$5B Iridium project was declared bankrupt. *Starlink*, a recent entrant in the satellite service area, is using a constellation of some 4,000 low-earth-orbiting spacecraft and appears so far to have been able to break through this financial barrier. Using reusable launch vehicles, smaller (and lighter) satellites, transponder arrays on board, and a new generation of digital-signal-processing capabilities, Starlink is in a position to offer retail access services of 100 Mbps or more to individual customers. The low altitude of the spacecraft means that the Starlink service competes directly with terrestrial access services in terms of performance. The introduction of inter-spacecraft laser links means that the system can provide a service in any location, and the limiting factor, as with the Iridium effort decades ago, is obtaining the necessary clearances and licenses to have customers located in the respective national geographies. Starlink is certainly revolutionary in terms of capacity, coverage, and cost. The questions are whether it is sufficiently revolutionary and whether it can scale up to provide a high-capacity service to hundreds of millions of users. At this point in time these questions are not easy to answer, but the limitations inherent in *Low Earth Orbit* (LEO)-based services point to a potential advantage in terrestrial-based access networks. Nevertheless, Starlink is redefining the Internet access market space in many countries, and setting price/performance benchmarks that their terrestrial competitors now have to match.

If we move away from access networks to look at the changes in the “core” of the Internet over the past 25 years, then once more we can see a dramatic change. In 1998, the Internet was constructed using the margins of oversupply in the telephone networks. In 1998, the core infrastructure of most ISPs was still being built by leasing telephone trunk supergroups (E-1 and T-1 circuits, and then E-3 and T-3 as capacity needs escalated, and then OC-1 circuits).

While it was not going to stop here, squeezing even more capacity from the network was now proving to be a challenge; 622-Mbps IP circuits were being deployed, although many of them were constructed using 155-Mbps *Asynchronous Transfer Mode* (ATM) circuits using router-based load balancing to share the IP load over four of these circuits in parallel. Gigabit circuits were just around the corner, and the initial exercises of running IP over 2.5-Gbps *Synchronous Digital Hierarchy* (SDH) circuits were being undertaken in 1998.

In some ways 1998 was a pivotal year for IP transmission. Until this time, IP was still just one more data application that was positioned as just another customer of the telco's switched-circuit infrastructure. This telco infrastructure was designed and constructed primarily to support telephony. From the analogue voice circuits to the 64K digital circuit through to the higher-speed trunk bearers, IP had been running on top of the voice network infrastructure. Communications infrastructure connected population centres where there was call volume. The Internet had different demands. Internet traffic patterns did not mirror voice traffic, and IP performance is sensitive to every additional millisecond of delay. Constraining the Internet to the role of an overlay placed on top of a voice network was showing signs of stress, and by 1998 things were changing. The Internet had started to make ever larger demands on transmission capacity, and the driver for further growth in the network infrastructure was now not voice, but data. It made little sense to provision an ever-larger voice-based switching infrastructure just to repackage it as IP infrastructure, and by 1998 the industry was starting to consider just what an all-IP high-speed network would look like, building an IP network all the way from the photon in a fibre-optic cable all the way through to the design of the Internet application.

Fibre Optics

At the same time, the fibre-optic systems were changing with the introduction of *Wave Division Multiplexing* (WDM). Older fibre equipment with electro-optical repeaters and *Plesiochronous Digital Hierarchy* (PDH) multiplexors allowed a single fibre pair to carry around 560 Mbps of data. WDM allowed a fibre pair to carry multiple channels of data using different wavelengths, with each channel supporting a data rate of up to 10 Gbps. Channel capacity in a fibre strand was between 40 and 160 channels using *Dense WDM* (DWDM). Combined with the use of all-optical amplifiers, the most remarkable part of this entire evolution in fibre systems is that a cable system capable of an aggregate capacity of a terabit can be constructed today for much the same cost as a 560-Mbps cable system of the mid-1990s. That's a cost-efficiency improvement of a factor of one million in a decade. The drive to deploy these high-capacity DWDM fibre systems was never based on expansion of telephony. The explosive growth of the industry was all about supporting the demand for IP. So, it came as no surprise that at the same time as the demand for IP transmission was increasing there was a shift in the transmission model where instead of plugging routers into telco switching gear and using virtual point-to-point circuits for IP, we started to plug routers into wavelengths of the DWDM equipment and operate all-IP networks in the core of the Internet.

DWDM is not the only technology that has fundamentally changed these core transmission systems in the past 25 years. Two further technologies have been transformational in the fibre-optic area. The first is the use of optical amplifiers. *Erbium Doped Fibre Amplifiers* provide a highly efficient means of signal amplification without the need to convert the signal back into a digital form and then passing it back through a digital/analogue converter to modulate the next-stage laser driver. This technology has allowed fibre systems to support terabit-per-second capacity without necessarily having to integrate terabit-per-second digital systems. The second fundamental change was a switch in signal modulation from a basic on/off signal into a signal modulation technique that uses signal-phase amplitude and polarity to increase the total capacity of a wavelength within a fibre strand. *Digital Signal Processors* (DSPs) offered the key technology here, and as we improve on the track width of conductor tracks in these processors we can increase the gate count on a single chip, thereby allowing support for more complex signal manipulation algorithms. These algorithms can be used to increase the sensitivity of the DSP function. In 2010, we were using 40-nm track silicon chips in DSPs, supporting *Polarization Mode Quadrature Phase Shift Keying* (PM-QPSK), which allowed a cable to operate with 100-Gbps data rates in a single wavelength, or an aggregation of 8 Tbps in a fibre strand. In 2023, DSPs now use 5-nm tracks, which can support PCS-144QAM modulation of a base 190-Gbaud signal, which can support 2.2-Tbps data rates per wavelength, or 105-Tbps total capacity per fibre strand. A 12-strand cable would have a total capacity of 1.2 Pbs.

Such very-high-performance fibre cable systems are generally used in submarine cable systems to link data centres between continents. In data-centre contexts and other terrestrial scenarios, we are now using 200- and 400-G per wavelength fibre system as the common technology base. The major outcome is that, in general, transmission is no longer a scarce resource. It is in every sense of the term an abundant commodity. There is no sense in rationing access to communications capacity, be it short or long haul. This change is a major one not only in the economic framework of the communications industry, but also in phrasing the way in which we use communications. In a scarce system, we tend to use “just-in-time” delivery mechanisms, passing content across the communications system only when it is needed, while an abundant system allows us to use “just-in-case” delivery mechanisms, causing a dramatic impact on the architecture of the Internet. Indeed, this extraordinary increase in the underlying capacity of our communications infrastructure through the past 25 years is perhaps the most significant change in the entire landscape of the Internet, as we will see when we look at content networking.

Network Management

In network operations, we are seeing some stirrings of change, but it appears to be a rather conservative area, and adoption of new network management tools and practices takes time.

The Internet converged on using the *Simple Network Management Protocol* (SNMP) more than a quarter of a century ago, and despite its security weaknesses, its inefficiency, its incredibly irritating use of *Abstract Syntax Notation One* (ASN.1), and its application in sustaining some forms of *Distributed Denial-of-Service* (DDoS) attacks, it still enjoys widespread use. But SNMP is only a network monitoring protocol, not a network configuration protocol, as anyone who has attempted to use SNMP write operations can attest. The more recent *Network Configuration Protocol* (NETCONF) and *Yet Another Next Generation* (YANG) data modelling languages are attempting to pull this area of configuration management into something a little more usable than *Command-Line Interface* (CLI) scripts driving interfaces on switches.

At the same time, we are seeing orchestration tools such as *Ansible*, *Chef*, *Network Automation and Programmability Abstraction Layer with Multivendor* (NAPALM), and *Salt* enter the network operations space, permitting the orchestration of management tasks over thousands of individual components. These network operations-management tools are welcome steps forward to improve the state of automated network management, but it's still far short of a desirable endpoint. The desired endpoint of a fully automated network-management framework is still far from our reach. Surely it must be feasible to feed an adaptive autonomous control system with the network infrastructure and available resources, and allow the control system to monitor the network and modify the operating parameters of network components to continuously meet the service-level objectives of the network? Where is the driverless car for driving networks? Maybe the next 10 years will get us there.

The Internet Layer

If our transmission systems have been subject to dramatic changes in the past quarter century, then what has happened at the IP layer over the same period?

First, we need to consider the “elephant” in the Internet layer room. One fundamental change at the Internet level of the protocol stack was meant to have all happened some 20 years ago, and that's the transition to IP version 6. Twenty-five years ago, in 1998, we were forecasting that we would have consumed all the remaining unallocated IPv4 addresses by around 2025. That estimate gave us slightly more than 25 years, so there was no particular sense of urgency. We didn't need to ring the emergency bell or raise any alarms. The overall aim was to proceed in an orderly manner. Things took a different course because we failed to appreciate the true impact of the shift of the Internet to mobile devices. All of a sudden, we were dealing with an Internet with billions of users, using billions of new mobile devices, and our comfortable predictions of a stately and steady run-down of the IPv4 address pools were discarded about as quickly as you could say the word “iPhone.” From “all the time in the world” we reached “no time left to do anything” within a year.

In the 5-year period between 2005 and 2010, when mobile services exploded in volume, the total count of allocated IP addresses rose from 1.5B addresses to 3.1B, from a total address pool of 3.7B addresses. The network had doubled in size, and the time left to complete the transition had shrunk from more than 20 years to a little over 1!

At that point, all the plans for an orderly transition were discarded, and many network administrators scrambled to obtain IPv4 addresses, further depleting the IPv4 pools. The central pool of IPv4 addresses, operated by the *Internet Assigned Numbers Authority* (IANA), was exhausted in February 2011. The *Asia Pacific Network Information Centre* (APNIC) depleted its IPv4 pool in April of that year, the *Réseaux IP Européens Network Coordination Centre* (RIPE NCC) 18 months later, the *Latin America and Caribbean Network Information Centre* (LACNIC) in 2014, and the *American Registry for Internet Numbers* (ARIN) in 2015. We had expected that this situation would motivate network operators to hasten their plans for IPv6 deployment, yet, perversely, that did not happen. Less than 1% of the Internet user base was using IPv6 in 2011. Five years later, as each of the *Regional Internet Registries* (RIRs) ran down their remaining pools of IPv4 addresses, this Internet-wide IPv6 user count had increased to just 5%. In 2023, the process is still underway, and some 35% of the Internet user base has IPv6 capability. I'm not sure anyone is willing to predict how long this anomalous situation of running the IPv4 Internet "on an empty tank" will persist.

NATs

How has the Internet managed to continue to operate, and even grow, without a supply of new IPv4 addresses? In a word, the answer is "NATs." While the *Network Address Translator* (NAT) concept received little fanfare when it was first published, it has enjoyed massive deployment over the past 25 years, and today NATs are ubiquitous. The application architecture of the Internet has changed, and we are now operating a client/server framework. Servers have permanent IP addresses, while clients "borrow" a public IPv4 address to complete a transaction and return it back to a common pool when they are done. Time-sharing IP addresses, and also using the 16-bit source port field in TCP and the *User Datagram Protocol* (UDP), has managed to extend the IPv4 address space by some 20 bits, making the IPv4+NAT address space up to a million times larger than the original 32-bit IPv4 address space. In practice, the story is a little more complicated than that, and some very large service providers have reached logistical limits in using NATs to compensate for the exhaustion of IPv4 addresses. This situation has motivated these providers to transition to a dual-stack mode of operation, and they are relying on a dual-stack host behaviour that prefers to use IPv6 when possible, thus relieving the pressure on the IPv4 NAT functions

NATs have prompted a significant change at the IP level in changing the default assumption about the semantics of an IP address. An IP address is no longer synonymous with the persistent identity of the remote party, but it has assumed the role of an ephemeral session token.

The leisurely pace of the IPv6 transition is partly due to this altered role of addresses, as we no longer require every connected device to have a persistently assigned globally unique IP address.

IPv6 and NATs are not the only areas of activity in the Internet layer in the past 25 years. We have tried to change many parts of the Internet layer, but interestingly, few, if any, of the proposed changes have managed to gain any significant traction out there in the network. The functions performed at the Internet layer of the protocol stack are no different from those of 25 years ago. IP Mobility, Multicast, and *IP Security* (IPSec) are largely Internet layer technologies that have failed to gain significant levels of traction in the marketplace of the public Internet.

QoS *Quality of Service* (QoS) was a “hot topic” in 1998, and it involved the search for a reasonable way for some packets to take some form of expedited path across the network, while other packets took an undifferentiated path. We experimented with various forms of signalling, packet classifiers, queue-management algorithms, and interpretations of the *Type of Service* bits in the IPv4 packet header, and we explored the QoS architectures of *Integrated and Differentiated Services* in great detail. However, QoS never managed to get established in mainstream Internet service environments. In this case, the Internet took a simpler direction, and in response to not enough network capacity we just augmented the network to meet demand.

Again, this is an aspect of the altered mindset when the communication system shifts from scarcity and rationing to one of abundance. We have given up installing additional intricate mechanisms in the network, in host protocol stacks, and even in applications to negotiate how to share insufficient network capacity. So far, the simple approach of just adding more capacity to the network has prevailed, and QoS remains largely unused.

MPLS The switch from circuit switching to packet switching has never managed to achieve universal acceptance. We have experimented with putting circuits back into the IP datagram architecture in various ways, most notably with the *Multi-Protocol Label Switching* (MPLS) technology. This technology used the label-swapping approach that was previously used in X.25, *Frame Relay* and ATM virtual circuit-switching systems, and it created a collection of virtual paths from each network ingress to each network egress across the IP network. The original idea was that in the interior of the network you no longer needed to load up a complete routing table into each switching element, and instead of performing destination-address lookup you could perform a much smaller, and hopefully faster, label lookup. This performance differentiator did not eventuate and switching packets using the 32-bit destination address in a fully populated forwarding table continued to present much the same level of cost efficiency at the hardware level as virtual circuit label switching.

However, one aspect of MPLS and similar approaches has proved to be invaluable for many network operators. A general-purpose network utility has many disparate client networks, and a single packet-switched environment does not allow the network operator to control the way in which the common network resource is allocated to each client network. It also does not readily support segmentation of reachability. Virtual circuit overlays, such as MPLS, provide mechanisms to control resource allocation and constrain cross-network leakage, and for many network operators these reasons are adequate to head down an MPLS-like path for their network platform.

Routing

Moving sideways at this level of the protocol stack, we probably should look at the evolution of routing technologies. The early 1990s saw a flurry of activity in the routing space, and various routing protocols were quickly developed and deployed. By 1998 the conventional approach to routing was to use either *Intermediate System-to-Intermediate System* (IS-IS) or *Open Shortest Path First* (OSPF) as the interior routing protocol, and the *Border Gateway Protocol* (BGP) as the inter-domain routing protocol. This picture has remained constant right up to today. In some ways, it is reassuring to see a basic technology that can sustain a quite dramatic growth rate through many years of scaling, but in other ways it is less reassuring to see that the unresolved issues we had with the routing system in 1998 are largely still with us today.

The largest of these unresolved issues lies in the trust we place in the inter-domain routing system of the Internet. There is no overall orchestration of the routing system. Each network advertises reachability information to its adjacent networks and selects what it regards as the “best” reachability information from the set received from these same network peers. This mutual trust that each network places in all other networks can, and has, been abused in various ways. The effort to allow each routing entity to distinguish between what is a “correct” item of routing information and what is a “false” route has a rich history of initiatives that have faltered for one reason or another. The most recent effort in this space is built upon the foundations of the number system, and it uses the association of a public/private key pair with the current holders of addresses and autonomous system numbers, allowing these holders to issue signed authorities about the use of these number resources in the context of routing, and by coupling these authorities with the information being propagated in the routing system, the intention being that unauthorized use cases will be detected.

RPKI

This effort, the *Resource Public Key Infrastructure* (RPKI), has achieved some level of acceptance in the networking space, and in 2023 around one-third of all route objects have associated RPKI credentials. The work is still “in progress” because the more challenging aspect of this work is to associate verifiable credentials with the propagation route through a network that does not impose onerous burdens on the routing system and is not overly fragile in its operation.

The extended period where the routing system has operated in a state that essentially cannot be trusted has prompted the application layer to generate its own mechanisms of trust. These days it is largely left to *Transport Layer Security* (TLS) to determine whether a client has reached its intended server. Given that we have been unable to construct a secured routing system for many decades, the question arises whether there is still the same level of need for such a system that we had some 25 years ago, given that the application space sees this problem as largely solved through the close-to-ubiquitous use of TLS.

This tension between the Internet layer and the upper layers of the protocol stack is also evident in the way in which we have addressed the perennial issue of location and identity. One of the original simplifications in the IP architecture was to bundle the semantics of identity, location, and forwarding into an IP address. While that has proved phenomenally effective in terms of simplicity of applications and simplicity of IP networks, it has posed some serious challenges when considering mobility, routing, protocol transition, and network scaling. Each of these aspects of the Internet would benefit considerably if the Internet architecture allowed identity to be distinct from location. Numerous efforts have been directed at this problem over the past decade, particularly in IPv6, but so far, we really haven't arrived at an approach that feels truly comfortable in the context of IP. The problem we appear to have been stuck on for the past decade is that if we create a framework of applications that use identity as a rendezvous mechanism and use an IP layer that requires location, then how is the mapping between identity and location distributed in an efficient and suitably robust manner? The transport layer of the protocol stack has also looked at the same space and developed some interesting approaches, as we will see in the next section.

Transport

Back in 1998 the transport layer of the IP architecture consisted of UDP and TCP, and the network use pattern was around 95% TCP and 5% UDP. It has taken all of the intervening 25 years, but this picture has finally changed.

We have developed some new transport protocols in this period, such as the *Datagram Congestion Control Protocol* (DCCP) and the *Stream Control Transmission Protocol* (SCTP), which can be regarded as refinements of TCP to extend a flow-control mechanism to apply to datagram streams in the case of DCCP and a shared flow-control state over multiple reliable streams in the case of SCTP. However, in a world of transport-aware middleware that has been a constant factor over this period, the level of capability to actually deploy these new protocols in the public Internet is marginal at best. Firewalls do not recognize these more recent transport protocols, NATs and similar, and as a result, the prospects of wide-scale deployment of such protocols in the public Internet are not very good. We seem to be firmly stuck in a world of TCP and UDP.

TCP has proved to be remarkably resilient over the years, but as the network increases in capacity the ability of TCP to continue to deliver ever-faster data rates over distances that span the globe is becoming a significant issue. Much work has been done to revise the TCP flow-control algorithms so that they still share the network fairly with other concurrent TCP sessions yet can ramp up to multi-gigabit-per-second data-transfer rates and sustain those rates over extended periods of time. The mainstream TCP flow-control protocol has been shifting from the conventional Reno-styled protocol to CUBIC, which attempts to find a stable sending rate and then slowly add flow pressure to the network path to see if the network can support greater sending rates. The response to packet drop remains a dramatic rate drop, but not quite as dramatic as the rate halving of Reno, but nevertheless it is still a drop-sensitive ack-paced flow-control protocol.

However, the picture has changed with the introduction of the *Bottleneck Bandwidth and Round-Trip* (BBR) protocol. Driving the network into the point not only of network queue formation, but right to the point of queue overflow and packet loss, is a crude approach. The problem here is that packet loss represents a loss of feedback, and in a feedback-based flow-control protocol, this loss of feedback pushes the protocol into a space where it has to pull back its sending rate to re-establish a signal flow. BBR represents a different way of looking at flow control, and it attempts to drive the flow to the point of the onset of queue formation in the network rather than aiming at the point of queue collapse. This process reduces the latency of the flow and the cost of network switching equipment by reducing the very-high-speed fast memory buffer requirements.

This area is not the only one of new experimentation in changing the TCP congestion-control, paradigm. Another approach is being explored in the *Low Latency Low Loss Scalable* throughput initiative (L4S), which is looking at incorporating network signals into the flow-control algorithm. Here the packet switches use the *Explicit Congestion Notification* (ECN) signal in the IP header when standing queues start to form. The receiver of this signal is expected to back off its sending rate in a manner similar to packet loss. The advantage of this approach is that there is no loss of feedback signalling, and the flow reacts to the formation of congestion conditions rather than the end point of queue collapse. However, ECN requires the deployment of ECN-marking equipment, and the effort of synchronising network equipment and transport-protocol behaviours is far greater when compared to protocol-only approaches such as BBR.

Other initiatives in the transport space that are also worthy of note include Multipath TCP and QUIC.

The first of these initiatives is *Multipath TCP*. The observation here is based around the increasing ubiquity of both Wi-Fi and cellular radio services, and the configuration of most mobile devices to include the ability to access both of these networks.

In general, the choice of which network interface to use is a single decision made by the mobile platform for all active applications. When a usable Wi-Fi network is detected, the device will prefer to use that connection for all new connections because it is assumed that the Wi-Fi service will be cheaper for the user and will operate at a higher performance level. But if performance and resilience are issues, then can we allow a TCP session to use *all* the available networks at once, and optimise its use of these multiple network paths to the destination such that the total data throughput is optimised? This is the objective of Multipath TCP, where a single TCP session is broken into numerous sub-sessions, where each sub-session uses a different network path by using a different local network interface.

Multipath TCP allows separate TCP states to control the flows passing across each network path to optimise throughput. It also can permit flow migration, allowing a logical TCP flow to switch from one network path to another while preserving integrity. The interesting aspect of this behaviour is that the control of the multipath behaviours is, in the first instance, under the control of the application rather than the host platform. This response was an early one to recognize the increasing capacity and diversity in edge networks, and how we could respond to this situation at the transport session level.

QUIC

The second initiative, which for me is a fundamental change in transport capabilities and functions, is the introduction of the QUIC protocol. At its simplest level, you could say QUIC is a packaging of the combination of TCP and TLS into a UDP wrapping. However, I would suggest that such a description is well short of the mark. QUIC is in many ways a far more ambitious transport protocol, bringing transport to the point where it is better suited to the current application behaviour. QUIC is intended to improve the transport performance for encrypted traffic with faster session setup. QUIC allows for further evolution of transport mechanisms with support for *Remote Procedure Calls* (RPC). QUIC also has integral support for concurrent session multiplexing that avoids TCP head-of-line blocking. QUIC encrypts the payload data, but unlike TLS, QUIC also encrypts the control data (the equivalent of the TCP header) and explicitly avoids the emerging TCP ossification within the network by occluding the entirety of the control exchange from the network of the session. QUIC is address agile, in that it can react to network-level address renumbering in an active QUIC session, as can occur with the presence of NATs on the network path. You can implement QUIC in user space, so applications can control their own transport functions. There is no longer a dependence on the platform in terms of the quality of the implementation of the transport service. With QUIC the application exercises a comprehensive level of control of the way the application interacts with the network.

Numerous lessons can be drawn from the QUIC experience. Any useful public communications medium needs to safeguard the privacy and integrity of the communications that it carries.

The time when open protocols represented an acceptable compromise between efficiency, speed, and privacy are over, and these days all network transactions in the public Internet need to be protected by adequate encryption. The QUIC model of wrapping a set of transactions, including both data and control transactions between a client and a server, into an end-to-end encryption state represents a minimum level of functionality in today's networking environment.

Secondly, QUIC provides needed additional transport functionality. TCP and UDP represent just two points of transport functions within a broader spectrum of possible transport models. UDP is just too susceptible to abuse, so we have heaped everything onto TCP. The issue is that TCP was designed as an efficient single streaming protocol, and retrofitting multiple sessions, short transactions, remote procedure calls, reliable single-packet transactions, and shared congestion states have proved to be impossible to implement in TCP.

Applications are now dominant in the Internet ecosystem, while platforms and networks are being commoditised. We are seeing loss of patience with platforms that provide common transport services for the application that they host, and a new model where the application comes with its own transport service. Taking an even broader perspective, the context of the success of the Internet lies in shifting the responsibility for providing service from the network to the end system. This shifting allowed us to make more efficient use of the common network substrate and push the cost of this packetization of network transactions over to end systems.

It shifted the innovation role from the large and lumbering telco operators into the nimbler world of software. QUIC takes it one step further, and pushes the innovation role from platforms to applications, just at the time when platforms are declining in relative importance within the ecosystem. From such a perspective, the emergence of an application-centric transport model that provides faster services, a larger repertoire of transport models, and encompassing comprehensive encryption were inevitable developments.

We have pushed the responsibility for end-to-end authentication into the transport layer with the close-to-ubiquitous TLS. TLS layers themselves above TCP (or merges with the TCP-like function in the case of QUIC), and the client passes the name of the service it intends to connect with to the remote server. The server passes its public key to the client, and the client authenticates this key using its own trust anchors. The server and client then negotiate a session key and proceed with an encrypted session. TLS is robust in almost every respect. Its major weakness lies in the highly distributed trust model, where there are hundreds of different operators of trusted credentials (certification authorities) and thousands of various registration agents. These entities are placed in a highly trusted role, and they can never lie.

The problem is that they have proved to be corruptible occasionally. They typically operate using online services, and a successful attack against such platforms can be abused to allow the issuance of trusted public certificates. We have invested considerable time and effort in shoring up this trust framework, but at the same time we have been working to make these public key certificates a commodity rather than an expensive luxury. The introduction of free certification authorities has succeeded in making these certificates available to all, but at the same time the totally automated certificate issuance process is liable to various forms of abuse. Despite these considerations, we have placed the entirety of the burden of service authenticity and session encryption onto TLS, to the point that other related efforts, such as IPsec, BGP routing security, and *Domain Name System Security Extensions* (DNSSEC) in the DNS, are generally perceived as optional extras rather than basic essentials to be included the security toolkit.

Applications and Services

This layer has also seen quite profound changes over the past quarter century, tracking the progress of increasing technical capability as well as consumer demands. In the late 1990s, the Internet was on the cusp of portal mania, where *LookSmart* was the darling of the Internet boom and everyone was trying to promote their own favourite “one stop shop” for all your Internet needs.

By 1998 the *AltaVista* search engine had made its debut, and these content-collation portals were already becoming passé. This change, from compiling directories and lists to active search, completely changed the Internet. These days we simply assume that we can type any query we want to into a search engine and the search machinery will deliver a set of pointers to relevant documents. And every time it occurs our expectations about the quality and utility of search engines are reinforced. Content is also changing as a result, as users no longer remain on a *site* and navigate around the site. Instead, users are driving the search engines, and pulling the relevant pages without reference to any other material. But it has not stopped there. Search engines are morphing into “instant answer machines,” where instead of providing a set of pointers to sources where there is a high level of correlation between the source and the question, the search engine attempts to extract material from the source and show what it believes is the answer to the implicit question in the search term. Even this process is just a way point in a longer journey, and today we are seeing *Artificial Intelligence* (AI) chat bots appearing, where the underlying data set that has been indexed by the search machinery is now being used as a corpus of data to drive an AI chat engine. The interaction is based on a natural language model.

If you thought of the Internet as an information resource, then the use of AI in this manner is a disturbing step. In this AI model the responding system generates plausible, but very definitely not necessarily factual, natural language responses to the implicit question in the query.

It's challenging to see this path from indexing data sources and matching query terms to the terms that primary sources use to one of a natural language generator that produces textual responses that are not grounded in facts, nor necessarily derived from primary sources, as being progress! Despite such misgivings about the deliberate abasement of the quality of the Internet as an information resource, this shift does fit into a larger picture of the transformation of the Internet to a mass entertainment vehicle, which is much of the driving force in today's content world.

Social Media

A related area of profound change has been the rise of social media. The television, radio, film, and print industries had evolved to use content mediators, compilers, and editors to curate their content, and the widespread deployment of highly capable user devices allowed end users to directly perform content production without the need to engage with mediators or producers. This situation has transformed many societies, and the social media platforms, including YouTube, Flickr, Facebook, Instagram, and TikTok, have been rocketed into societal prominence, prompting major debates about the role of these platforms and levels of societal influence that such platforms can generate.

Underlying these changes is another significant development, namely the change in the content economy. In 1998 content providers and ISPs were eyeing each other in an effort to gain user revenue. Content providers were unable to make pay-per-view and other forms of direct financial relationships with users work in their favour and argued that ISPs should fund content. After all, they pointed out, the only reason users paid for Internet access was the perceived value of the content they found there. ISPs, on the other hand, insisted that content providers were enjoying a "free ride" across the ISP-funded infrastructure, and content providers should contribute to network costs. The model that has gained ascendancy as a result of this unresolved tension is that of advertisement-funded content services, and this model has been able to sustain a vastly richer, larger, and more compelling content environment.

However, using this model comes at a price, and in this case the price lies in the motivations of the platforms that perform ad delivery. The critical objective now is to engage the user for longer periods, so that they can present more ads and glean more information about the user's profile. Merely informing a user is largely a transactional interaction, whereas entertaining a user can be far more lucrative in terms of generating advertising revenue because of the longer attention span. This model has been highly successful for some content players, particularly the current giants of streaming content, and it's therefore unsurprising that the largest entities in the content world, such as Alphabet, Microsoft, Amazon, and Apple, are more valuable in terms of market capitalization than their counterparts in the carriage world. We are now seeing the next round of the friction between content and carriage, where the access network operators are arguing that the content players should contribute to the costs of access carriage.

The *Domain Name System* (DNS) also merits a mention in this section. From one perspective, little has changed in this space, and the DNS name-resolution protocol hasn't changed to any appreciable extent. In some sense that's true, but at the same time there have been some significant changes.

DNS

The first of these changes is adoption of *Domain Name System Security Extensions* (DNSSEC), a framework that allows DNS clients to validate the answers that they receive from the DNS. The DNS has always been a point of security vulnerability in the Internet in that it has always been prone to various forms of attack where false answers are substituted in place of the genuine answer. DNSSEC provides a digital signature record to each normal record, and also implements an interlocked chain of signatures to link to the key associated with the root zone. A client may request the signature record to be provided with the normal response, and then make further requests to construct the validation chain all the way to the root zone. Successful validation assures a client that the data provided in the original response is authentic and current. The root zone of the DNS was first signed in 2010, but adoption of DNSSEC has been slow. While the addition of such a validation mechanism is undoubtedly a step forward in protecting users against various forms of name-based interference, the cost is increased fragility of the DNS and increased resolution times. One underlying problem is that the addition of digital signatures to a DNS response is highly likely to push the DNS into sending large responses, and large responses over a UDP-based transport is prone to fragmentation-based unreliability, and the switch to use TCP also takes time. What this problem has implied is that the path to adoption of DNSSEC has been slow, despite the obvious protections it can provide regarding potential tampering with the DNS.

The second major theme of change in the DNS concerns the larger issue of pervasive monitoring in the DNS, highlighted by the Snowden revelations of 2013. Most Internet transactions start with a call to the DNS, and the meta-data contained in DNS queries and responses provides a rich real-time profile of user activity, both in general and potentially on a user-by-user basis. This situation has prompted a concerted effort to improve the privacy aspects of the DNS as a protocol. One approach has been to take the existing use of DNS across a TCP session and add TLS to the TCP session, so the contents of the interaction between the client and the DNS server are impervious to third-party inspection or manipulation. This approach can be taken a step further with DNS over *Hypertext Transfer Protocol Secure* (HTTPS)/2, where the DNS payload has a lightweight HTTP wrapper in addition to TLS. This approach allows DNS traffic to be melded in with all other HTTP traffic as a further step of obscuring DNS transactions. More recently we have seen DNS over QUIC, using QUIC faster session start times and fast open capabilities to improve the performance of the arrangement, and DNS over HTTPS/3, which combines QUIC with HTTP object semantics.

The primary focus of this work has been the part of the DNS where the client's stub resolver interacts with a recursive resolver, because this scenario identifies the client. The useful property of this part of the DNS is that the same client/server setup is used repeatedly, so either a long-held secure transport session or a fast-reopen session can amortise the high setup cost of a reliable secure session over many subsequent queries, making the overall cost of such a secure transport arrangement more palatable.

Such measures still have some security problems, as the recursive resolver is privy to both the client's identity and the DNS queries that they make. Recent work has begun on an "oblivious" model of DNS operation, where the recursive resolver function is split in two and two layers of encryption are used. The client talks to the first party, a DNS relay over an encrypted session, and passes it a query that has been encrypted using the public key of the second party, the recursive resolver. The relay resends the encrypted DNS query to the recursive resolver to resolve. The first party knows the identity of the client, but not the DNS query that is being made. The second party knows the DNS query, but not the identity of the client.

This work on DNS privacy has extended into the scenarios of the recursive resolver talking with authoritative name servers, although it's unclear as to the extent of the security benefits (because the end user is not identified directly in such queries), nor is session reuse as feasible in this scenario.

Cloud

In many ways applications and services have been the high frontier of innovation in the Internet in this period. An entire revolution in open interconnection of content elements has taken place, and content is now a very malleable concept. It is no longer the case of "my computer, my applications, my workspace" or "your server, your content" but an emerging model where not only the workspace for each user is held in the network, but where the applications and services themselves are part of the network, and all are accessed through a generic mechanism based around permutations of the HTTPS access model. This world is one of the so-called *Cloud Services*. The world of cloud services takes advantage of abundance in computation, storage, and communications resources, and rather than a network facilitating users to connect to service delivery points, the cloud model inverts the model and attempts to bring replicant copies of content and services closer to the user. If distance equates to cost and performance in the networking world, then the cloud model dramatically shortens the distance between consumer and content, with obvious implications in terms of cost and performance reductions. The clouded Internet can achieve extremely challenging performance and cost objectives by changing the provisioning model of the content and service from "just-in-time" on-demand service to "just-in-case" pre-provisioning of local caches so that the local cache is ready if a local client accesses the service.

Cyber Hostility

We still are under relentless attack at all levels. We are beset by data leaks, surveillance, profiling, disruption, and extortion.

Attacks are now commonplace. Many of them are brutally simple, relying on a tragically large pool of potential zombie devices that are readily subverted and co-opted to assist in attacks. The attacks are often simple, such as UDP reflection attacks where a single UDP query generates a large response. The source address of the query is forged to be the address of the intended attack victim, and not much more needs to be done. A small query stream can result in a massive attack. UDP protocols such as SNMP, the *Network Time Protocol* (NTP), the DNS, and *memcached* have been used in the past and doubtless will be used again.

Why can't we fix this problem? We've been trying for decades, and we just can't seem to get ahead of the attacks. Advice to network operators to prevent the leakage of packets with forged source addresses was published more than two decades ago, in 2000. Yet massive UDP-based attacks with forged source addresses still persist today. Aged computer systems with known vulnerabilities continue to be connected to the Internet and are readily transformed into attack bots.

The picture of attacks is also becoming more ominous. Although we previously attributed these hostile attacks to "hackers," we quickly realised that a significant component of them had criminal motivations. The progression from criminal actors to state-based actors is also entirely predictable, and we are seeing an escalation of this cyber warfare arena with the investment in various forms of vulnerability exploitation that are considered desirable national capabilities.

It appears that a major problem here is that collectively we are unwilling to make any substantial investment in effective defence or deterrence. The systems that we use on the Internet are overly trusting to the point of irrational credulity. For example, the public key certification system used to secure web-based transactions is repeatedly demonstrated to be entirely untrustworthy, yet that's all we trust. Personal data is continually breached and leaked, yet all we seem to want to do is increase the number and complexity of regulations rather than actually use better tools that would effectively protect users.

The larger picture of hostile attacks is not getting any better. Indeed, it's getting much worse. If any enterprise has a business need to maintain a service that is always available for use, then any form of in-house provisioning is just not enough to withstand attack. These days only a handful of platforms can offer resilient services, and even then, it's unclear whether they could withstand the most extreme of attacks.

A constant background level of scanning and probing goes on in the network, and any form of visible vulnerability is ruthlessly exploited. One could describe today's Internet as a toxic wasteland, punctuated with the occasional heavily defended citadel.

Those who can afford to locate their services within these citadels enjoy some level of respite from this constant profile of hostile attack, while all others are forced to try to conceal themselves from the worst of this toxic environment, while at the same time aware that they will be completely overwhelmed by any large-scale attack. It is a sobering thought that about one-half of the world's population are now part of this digital environment. A more sobering thought is that many of today's control systems, such as power generation and distribution, water distribution, and road-traffic-control systems are exposed to the Internet.

IoT What makes this scenario even more depressing is the portent of the so-called *Internet of Things* (IoT). In those circles where Internet prognostications abound and policy makers flock to hear grand visions of the future, we often hear about the boundless future represented by this Internet of Things. This phrase encompasses some decades of the computing industry's transition from computers as esoteric pieces of engineering affordable only by nations to mainframes, desktops, laptops, handheld devices, and now wrist computers.

Where next? In the vision of the IoT, we are going to expand the Internet beyond people and press on using billions of these chattering devices in every aspect of our world. What do we know about the "things" that are already connected to the Internet? Some of them are not very good. In fact, some of them are just plain stupid. And this stupidity is toxic, in that their sometime-inadequate models of operation and security affect others in potentially malicious ways.

If such devices were constantly inspected and managed, we might see evidence of aberrant behaviour and correct it. But these devices are unmanaged and all but invisible. Examples include the controller for a web camera, the so-called "smart" thing in a smart television, or the controls for anything from a washing machine to a goods locomotive. Nobody is looking after these devices. When we think of an IoT we think of a world of weather stations, webcams, "smart" cars, personal fitness monitors, and similar things.

But what we tend to forget is that all of these devices are built on layers of other people's software that is assembled into a product at the cheapest possible price point. It may be disconcerting to realise that the web camera you just installed has a security model that can be summarised with the phrase: "no security at all," and it's actually offering a view of your house to the entire Internet. It may be slightly more disconcerting to realise that your electronic wallet is on a device that is using a massive compilation of open-source software of largely unknown origin, with a security model that is not completely understood, but appears to be susceptible to be coerced into being a "yes, take-all-you-want" device. It would be nice to think that we have stopped making mistakes in code, and from now on our software in our things will be perfect. But that's hopelessly idealistic. It's just not going to happen. Software will not be perfect. It will continue to have vulnerabilities.

It would be nice to think that this Internet of Things is shaping up as a market where quality matters, and consumers will select a more expensive product even though its functional behaviour is identical to a cheaper product that has not been robustly tested for basic security flaws. But that too is hopelessly naive.

The IoT will continue to be a marketplace where the compromises between price and quality will continue to push us on to the side of cheap rather than secure. What is going to stop us from further polluting our environment with a huge and diverse collection of programmed unmanaged devices with inbuilt vulnerabilities that will be all too readily exploited? What can we do to make this world of these stupid cheap toxic things less stupid and less toxic? So far, we have not found workable answers to this question.

Our ability to effectively defend the network and its connected hosts continues to be, on the whole, ineffectual. Anyone who still has trust in the integrity of the systems that make up the digital world is just hopelessly naive. This space is toxic and hostile, and we still have no idea how we can shift it to a different state that can resist such erosive and insidious attacks. But somehow, we are evidently not deterred by all this information. Somehow each of us has found a way to make the Internet work for us.

The Business of the Internet

As much as the application environment of the Internet has been on a wild ride over the past 25 years, the business environment has also had its tickets on the same roller coaster ride, and the list of business winners and losers includes some of the historical giants of the telephone world as well as the Internet-bred new wave of entrants.

In 1998, despite the growing momentum of public awareness, the Internet was still largely a curiosity. Its environment was inhabited by geeks, game players, and academics, whose rites of initiation were quite arcane. As a part of the data networking sector, the Internet was just one further activity among many, and the level of attention from the mainstream telco sector was still relatively low. Most Internet users were customers of independent ISPs, and the business relationship between the ISP sector and the telco was tense and acrimonious. The ISPs were seen as opportunistic leeches on the telco industry; they ordered large banks of phone lines, but never made any calls; their customers did not hang up after 3 minutes, but kept their calls open for hours or even days at a time, and they kept on ordering ever-larger inventories of transmission capacity, yet had business plans that made scribbles on the back of an envelope look professional by comparison.

The telco was unwilling to make large long-term capital investments in additional communications infrastructure to pander to the extravagant demands of a wildcat set of Internet speculators and their fellow travellers.

The telco, on the other hand, was slow, expensive, inconsistent, ill-informed, and hostile to the ISP business. The telco wanted financial settlements and bit-level accounting while the ISP industry appeared to manage quite well with a far simpler system of peering and tiering that avoided putting a value on individual packets or flows.

This relationship was never going to last, and it resolved itself in ways that in retrospect were quite predictable. From the telco perspective, it quickly became apparent that the only reasons the telco was being pushed to install additional network capacity at ever-increasing rates were demands from the ISP sector. From the ISP perspective, the only way to grow at a rate that matched customer demand was to become one's own carrier and take over infrastructure investment. And, in various ways, both outcomes occurred. Telcos bought up ISPs, and ISPs became infrastructure carriers.

All this activity generated considerable investor interest, and the rapid value escalation of the ISP industry and then the entire Internet sector generated the levels of wild-eyed optimism that are associated only with an exceptional boom. By 2000 almost anything associated with the Internet, whether it was a simple portal, a new browser development, a search engine, or an ISP, attracted investor attention, and the valuations of Internet start-ups achieved dizzying heights. Of course, one of the basic lessons of economic history is that every boom has an ensuing bust, and in 2001 the Internet collapse happened. The bust was as inevitable and as brutal as the preceding boom was euphoric. But, like the railway boom and bust of the 1840s, after the wreckage was cleared away what remained was a viable, and indeed a valuable, industry.

By 2003 the era of the independent retail ISP was effectively over. But it reshaped itself dramatically with the introduction of mobile services. It was the old telco sector that had secured spectrum allocations in the bidding wars in the early 2000s, and while they had thought that mobile voice would be the reason why these investments would make sense, it was mobile Internet services that proved to be the lasting service model. In this period, the Internet was the amalgam of the largest of the original ISP, the transformed cable television operators, and the mobile providers. Each national regime was populated with some three to five major service providers, and the business started to stabilise around this model.

Into this world came the content world of the Internet, using cloud-based models of service delivery to circumvent communications bottlenecks in the long-haul transit sector. They embarked on service models that included advertiser-funded models of content generation and delivery and direct subscription models, and the result has been so effective that the value of this sector is far greater than the traditional ISP and carriage sector. The content world is now the major funder of subsea cable systems, and the carriage world has been very reluctantly pushed into an undistinguished commodity role as a result.

This situation is reflective of a broader process of technology permeation. The telephone world used the network as the major focus of technology and investment. The edge devices, telephone handsets, were simple, cheap devices, whereas network switches and transmission elements were built to exacting and expensive standards. As we attached computers to the edges of the network, these devices were able to tolerate a broader spectrum of network behaviours, and had a lower base-level expectation of behaviour. Consequently, value has moved out from the core of the network to its edges.

But this process has also been reflected within these edge devices. We started with a model of a highly capable and complicated operating system platform, and relatively simple applications that used platform services. Some 25 years ago the release of Windows 98 was a Big Thing, and rightly so. As these edge devices become more capable and have higher processing capability, more local storage applications have elected to take on more of the responsibility in terms of the user's experience. In doing so they no longer rely on the release schedules of the platform provider, and they are no longer as concerned about the level of control being exercised by this platform provider and gaining an essential level of self-control. Modern browsers (Chrome and a few far smaller fellow travellers) are far more complex than most operating systems, and they continue to subsume functions and roles that the platform previously carried out. With DNS over HTTPS, the task of DNS name resolution can be transformed to an application function, rather than a common platform function. With QUIC, the transport protocol itself has been subsumed into the application space.

Not only have we seen the commoditisation of the network over the past 25 years, we have also seen similar commoditisation pressures on the end-device platforms and on the operating systems used on these devices. Even the browser space has been commoditised. The brunt of competitive differentiation in this industry has been pushed up the protocol stack into the content and service economy, and there is the distinct feeling that even in that space competitive differentiation is perhaps a misnomer, and what we have is a synthetic form of competition between a select small group of digital service-delivery behemoths that in any other time and context would probably be called a cartel.

What Now?

It's been a revolutionary quarter-century for us all, and the Internet has directly or indirectly touched the lives of almost every person on this planet. Current estimates put the number of regular Internet users at one half of the world's population.

Over this period, some of our expectations were achieved and then surpassed with apparent ease, while others remained elusive. And some things occurred that were entirely unanticipated. At the same time, very little of the Internet we have today was confidently predicted in 1998, while many of the problems we saw in 1998 remain problems today.

This work-in-progress means the next quarter-century will probably see the same level of intensity of yet more structural changes to the global communications sector. And that is a somewhat scary prospect, given the collection of other challenges that we will all confront in the coming decades. At the same time, I think it would be good to believe that the debut of the Internet in our world has completely rewritten what it means to communicate, the way in which we can share our experience and knowledge, and, hopefully, the ways in which we can work together on these challenges.

References and Further Reading

The Internet Protocol Journal has published articles on all the major aspects of the technical evolution of the Internet over the past 25 years. To illustrate the extraordinary breadth of these articles, I have included as references here some pointers to articles that have been published in IPJ.

- [1] William Stallings, “SSL: Foundation for Web Security,” *The Internet Protocol Journal*, Volume 1, No. 1, June 1998.
- [2] Fred Avolio, “Firewalls and Internet Security,” *The Internet Protocol Journal*, Volume 2, No. 2, June 1999.
- [3] William Stallings, “Gigabit Ethernet,” *The Internet Protocol Journal*, Volume 2, No. 3, September 1999.
- [4] Mark Handley and Jon Crowcroft, “Internet Multicast Today,” Volume 2, No. 4, December 1999.
- [5] Geoff Huston, “Quality of Service – Fact or Fiction?,” *The Internet Protocol Journal*, Volume 3, No. 1, March 2000.
- [6] Geoff Huston, “The Future for TCP,” *The Internet Protocol Journal*, Volume 3, No. 3, September 2000.
- [7] Chris Lonvick, “Securing the Infrastructure,” *The Internet Protocol Journal*, Volume 3, No. 3, September 2000.
- [8] William Stallings, “Mobile IP,” *The Internet Protocol Journal*, Volume 4, No. 2, June 2001.
- [9] Geoff Huston, “The Middleware Muddle,” *The Internet Protocol Journal*, Volume 4, No. 2, June 2001.
- [10] William Stallings, “MPLS,” *The Internet Protocol Journal*, Volume 4, No. 3, September 2001.
- [11] Stephen Kent, “Securing BGP: S-BGP,” *The Internet Protocol Journal*, Volume 6, No. 3, September 2003.
- [12] Russ White, “Securing BGP: soBGP,” *The Internet Protocol Journal*, Volume 6, No. 3, September 2003.
- [13] Geoff Huston, “Anatomy: Inside Network Address Translators,” *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.
- [14] Daniel McCarney, “Automatic Certificate Management,” *The Internet Protocol Journal*, Volume 20, No. 2, June 2017.

- [15] Charalampos Patrikakis, Michalis Masikos, and Olga Zouraraki, “Distributed Denial of Service Attacks,” *The Internet Protocol Journal*, Volume 7, No. 4, December 2004.
- [16] David Crocker, “Challenges in Anti-Spam Efforts,” *The Internet Protocol Journal*, Volume 8, No. 4, December 2005.
- [17] Vint Cerf, “A Decade of Internet Evolution,” *The Internet Protocol Journal*, Volume 11, No 2, June 2008.
- [18] Geoff Huston, “A Decade in the Life of the Internet,” *The Internet Protocol Journal*, Volume 11, No. 2, June 2008
- [19] Thayumanavan Sridhar, “Cloud Computing – A Primer,” *The Internet Protocol Journal*, Volume 12, No. 3, September 2009, Vol. 12, No. 4, December 2009.
- [20] Bob Hinden, “The Internet of Insecure Things,” *The Internet Protocol Journal*, Volume 20, No. 1, March 2017.
- [21] Andrei Robachevsky, “Improving Routing Security,” *The Internet Protocol Journal*, Volume 22, No. 2, July 2019.
- [22] Geoff Huston, “DNS Privacy,” *The Internet Protocol Journal*, Volume 22, No. 2, July 2019.
- [23] David Strom, “So You Want to Sell Your IPv4 Address Block?,” *The Internet Protocol Journal*, Volume 23, No. 2, September 2020.
- [24] Geoff Huston, “DNS Trends,” *The Internet Protocol Journal*, Volume 24, No. 1, March 2021.
- [25] Geoff Huston, “Securing Inter-Domain Routing,” *The Internet Protocol Journal*, Volume 24, No. 3, October 2021, and Volume 25, No. 1, April 2022
- [26] Geoff Huston, “Comparing TCP and QUIC,” *The Internet Protocol Journal*, Volume 25, No. 3, December 2022.
- [27] Geoff Huston, “Protocol Basics: The Network Time Protocol,” *The Internet Protocol Journal*, Volume 15, No. 4, December 2012.
- [28] Burton S. Kaliski Jr., “Minimized DNS Resolution: Into the Penumbra,” *The Internet Protocol Journal*, Volume 25, No. 3, December 2022.

GEOFF HUSTON AM, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

Supporters and Sponsors

Supporters



Internet
Society



Diamond Sponsors

Your logo here!

Ruby Sponsors



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Kjetil Aas	Václav Brožík	Michael Dragone	Jeffrey Greene	Javier Juan
Fabrizio Accatino	Christophe Brun	Joshua Dreier	Richard Gregor	David Jump
Michael Achola	Gareth Bryan	Lutz Drink	Martijn Groenleer	Anders Marius
Martin Adkins	Ron Buchalski	Aaron Dudek	Geert Jan de Groot	Jørgensen
Melchior Aelmans	Paul Buchanan	Dmitriy Dudko	Ólafur Guðmundsson	Merike Kao
Christopher Affleck	Stefan Buckmann	Andrew Dul	Christopher Guemez	Andrew Kaiser
Scott Aitken	Caner Budakoglu	Joan Marc Riera	Gulf Coast Shots	Naoki Kambe
Jacobus Akkerhuis	Darrell Budic	Duocastella	Sheryll de Guzman	Christos Karayiannis
Antonio Cuñat Alario	BugWorks	Pedro Duque	Rex Hale	Daniel Karrenberg
William Allaire	Scott Burleigh	Holger Durer	Jason Hall	David Kekar
Nicola Altan	Chad Burnham	Karlheinz Dölger	James Hamilton	Stuart Kendrick
Shane Amante	Randy Bush	Mark Eanes	Darow Han	Robert Kent
Marcelo do Amaral	Colin Butcher	Andrew Edwards	Handy Networks LLC	Jithin Kesavan
Matteo D'Ambrosio	Jon Harald Bøvre	Peter Robert Egli	Stephen Hanna	Jubal Kessler
Selva Anandavel	Olivier Cahagne	George Ehlers	Martin Hannigan	Shan Ali Khan
Jens Andersson	Antoine Camerlo	Peter Eisses	John Hardin	Nabeel Khatri
Danish Ansari	Tracy Camp	Torbjörn Eklöv	David Harper	Dae Young Kim
Finn Arildsen	Brian Candler	Y Ertur	Edward Hauser	William W. H.
Tim Armstrong	Fabio Caneparo	ERNW GmbH	David Hauweele	Kimandu
Richard Artes	Roberto Canonico	ESdatCo	Marilyn Hay	John King
Michael Aschwanden	David Cardwell	Steve Esquivel	Headcrafts SRLS	Russell Kirk
David Atkins	Richard Carrara	Jay Etchings	Hidde van der Heide	Gary Klesk
Jac Backus	John Cavanaugh	Mikhail Evstiounin	Johan Helsingius	Anthony Klopp
Jaime Badua	Lj Cemerar	Bill Fenner	Robert Hinden	Henry Kluge
Bent Bagger	Dave Chapman	Paul Ferguson	Damien Holloway	Michael Kluk
Eric Baker	Stefanos Charchalakis	Ricardo Ferreira	Alain Van Hoof	Andrew Koch
Fred Baker	Molly Cheam	Kent Fichtner	Edward Hotard	Ia Kochiashvili
Santosh Balagopalan	Greg Chisholm	Armin Fisslthaler	Bill Huber	Carsten Koempe
William Baltas	David Chosrova	Michael Fiumano	Hagen Hultzs	Richard Koene
David Bandinelli	Marcin Cieslak	The Flirble Organisation	Kauto Huopio	Alexander Kogan
A C Barber	Lauris Cikovskis	Jean-Pierre Forcioli	Asbjørn Højmark	Matthijs Koot
Benjamin Barkin-Wilkins	Brad Clark	Gary Ford	Kevin Iddles	Antonin Kral
Feras Batainah	Narelle Clark	Susan Forney	Mika Ilvesmaki	Robert Krejčí
Michael Bazarewsky	Horst Clausen	Christopher Forsyth	Karsten Iwen	John Kristoff
David Belson	James Cliver	Andrew Fox	Joseph Jackson	Terje Krogdahl
Richard Bennett	Guido Coenders	Craig Fox	David Jaffe	Bobby Krupczak
Matthew Best	Joseph Connolly	Fausto Franceschini	Ashford Jaggernaut	Murray Kuchera
Hidde Beumer	Steve Corbató	Valerie Fronczak	Thomas Jalkanen	Warren Kumari
Pier Paolo Biagi	Brian Courtney	Tomislav Futivic	Jozef Janitor	George Kuo
Arturo Bianchi	Beth and Steve Crocker	Laurence Gagliani	Martijn Jansen	Dirk Kurfuerst
John Bigrow	Dave Crocker	Edward Gallagher	John Jarvis	Mathias Körber
Orvar Ari Bjarnason	Kevin Croes	Andrew Gallo	Dennis Jennings	Darrell Lack
Tyson Blanchard	John Curran	Chris Gamboni	Edward Jennings	Andrew Lamb
Axel Boeger	André Danthine	Xosé Bravo Garcia	Aart Jochem	Richard Lamb
Keith Bogart	Morgan Davis	Oswaldo Gazzaniga	Nils Johansson	Yan Landriault
Mirko Bonadei	Jeff Day	Kevin Gee	Brian Johnson	Edwin Lang
Roberto Bonalumi	Rodolfo Delgado-Bueno	Greg Giessow	Curtis Johnson	Sig Lange
Lolke Boonstra	Julien Dhallenne	John Gilbert	Richard Johnson	Markus Langenmair
Julie Bottoff Photography	Freek Dijkstra	Serge Van Ginderachter	Jim Johnston	Fred Langham
Gerry Boudreaux	Geert Van Dijk	Greg Goddard	Jonatan Jonasson	Tracy LaQuey Parker
Leen de Braal	David Dillow	Tiago Goncalves	Daniel Jones	Alex Latzko
Kevin Breit	Richard Dodsworth	Ron Goodheart	Gary Jones	Jose Antonio Lazaro
Thomas Bridge	Ernesto Doelling	Octavio Alfageme	Jerry Jones	Lazaro
Ilia Bromberg	Michael Dolan	Gorostiaga	Michael Jones	Antonio Leding
Lukasz Bromirski	Eugene Doroniuk	Barry Greene	Amar Joshi	Rick van Leeuwen

Simon Leinen	Mohammad Moghaddas	Blahoslav Popela	Timothy Schwab	Peter Tomsu Fine Art
Robert Lewis	Charles Monson	Andrew Potter	Roger Schwartz	Photography
Christian Liberale	Andrea Montefusco	Ian Potts	SeenThere	Joseph Toste
Martin Lillepuu	Fernando Montenegro	Eduard Llull Pou	Scott Seifel	Rey Tucker
Roger Lindholm	Roberto Montoya	Tim Pozar	Paul Selkirk	Sandro Tumini
Link Light Networks	Joel Moore	David Raistrick	Andre Serralheiro	Angelo Turetta
Chris and Janet Lonvick	John More	Priyan R Rajeevan	Yury Shefer	Michael Turzanski
Sergio Loreti	Maurizio Moroni	Balaji Rajendran	Yaron Sheffer	Phil Tweedie
Eric Louie	Brian Mort	Paul Rathbone	Doron Shikmoni	Steve Ulrich
Adam Loveless	Soenke Mumm	William Rawlings	Tj Shumway	Unitek Engineering AG
Josh Lowe	Tariq Mustafa	Mujtiba Raza Rizvi	Jeffrey Sicuranza	John Urbanek
Guillermo a Loyola	Stuart Nadin	Bill Reid	Thorsten Sideboard	Martin Urwaleck
Hannes Lubich	Michel Nakhla	Petr Rejhon	Greipur Sigurdsson	Betsy Vanderpool
Dan Lynch	Mazdak Rajabi Nasab	Robert Remenyi	Fillipe Cajaiba da Silva	Surendran Vangadasalam
David MacDuffie	Krishna Natarajan	Rodrigo Ribeiro	Andrew Simmons	Ramnath Vasudha
Sanya Madan	Naveen Nathan	Glenn Ricart	Pradeep Singh	Randy Veasley
Miroslav Madić	Darryl Newman	Justin Richards	Henry Sinnreich	Philip Venables
Alexis Madriz	Thomas Nikolajsen	Rafael Riera	Geoff Sisson	Buddy Venne
Carl Malamud	Paul Nikolich	Mark Risinger	John Sisson	Alejandro Vennera
Jonathan Maldonado	Travis Northrup	Fernando Robayo	Helge Skrivervik	Luca Ventura
Michael Malik	Marijana Novakovic	Michael Roberts	Terry Slattery	Scott Vermillion
Tarmo Mämers	David Oates	Gregory Robinson	Darren Sleeth	Tom Vest
Yogesh Mangar	Ovidiu Obersterescu	Ron Rockrohr	Richard Smit	Peter Villemoes
John Mann	Jim Oplotnik	Carlos Rodrigues	Bob Smith	Vista Global Coaching &
Bill Manning	Tim O'Brien	Magnus Romedahl	Courtney Smith	Consulting
Harold March	Mike O'Connor	Lex Van Roon	Eric Smith	Dario Vitali
Vincent Marchand	Mike O'Dell	Marshall Rose	Mark Smith	Rüdiger Volk
Normando Marcolongo	John O'Neill	Alessandra Rosi	Tim Sneddon	Jeffrey Wagner
Gabriel Marroquin	Carl Örne	David Ross	Craig Snell	Don Wahl
David Martin	Packet Consulting	William Ross	Job Snijders	Michael L Wahrman
Jim Martin	Limited	Boudhayan	Ronald Solano	Lakhinder Walia
Ruben Tripiana Martin	Carlos Astor Araujo	Roychowdhury	Asit Som	Laurence Walker
Timothy Martin	Palmeira	Carlos Rubio	Ignacio Soto Campos	Randy Watts
Carles Mateu	Gordon Palmer	Rainer Rudigier	Evandro Sousa	Andrew Webster
Juan Jose Marin Martinez	Alexis Panagopoulos	Timo Rüter	Peter Spekrijse	Jd Wegner
Ioan Maxim	Gaurav Panwar	RustedMusic	Thayumanavan Sridhar	Tim Weil
David Mazel	Chris Parker	Babak Saberi	Paul Stancik	Westmoreland
Miles McCredie	Alex Parkinson	George Sadowsky	Ralf Stempfner	Engineering Inc.
Brian McCullough	Craig Partridge	Scott Sandefur	Matthew Stenberg	Rick Wesson
Joe McEachern	Manuel Uruena Pascual	Sachin Sapkal	Martin Štěpánek	Peter Whimp
Alexander McKenzie	Ricardo Patara	Arturas Satkovskis	Adrian Stevens	Russ White
Jay McMaster	Dipesh Patel	PS Saunders	Clinton Stevens	Jurrien Wijlhuizen
Mark Mc Nicholas	Dan Paynter	Richard Savoy	John Streck	Derick Winkworth
Olaf Mehlberg	Leif Eric Pedersen	John Sayer	Martin Streule	Pindar Wong
Carsten Melberg	Rui Sao Pedro	Phil Scarr	David Strom	Makarand Yerawadekar
Kevin Menezes	Juan Pena	Gianpaolo Scassellati	Colin Strutt	Phillip Yialeloglou
Bart Jan Menkveld	Chris Perkins	Elizabeth Scheid	Viktor Sudakov	Janko Zavernik
Sean Mentzer	Michael Petry	Jeroen Van Ingen	Edward-W. Suor	Bernd Zeimet
Eduard Metz	Alexander Peuchert	Schenau	Vincent Surillo	Muhammad Ziad
William Mills	David Phelan	Carsten Scherb	Terence Charles Sweetser	Ziauddin
David Millsom	Harald Pilz	Ernest Schirmer	T2Group	Tom Zingale
Desiree Miloshevic	Derrell Piper	Benson Schliesser	Roman Tarasov	Jose Zumalave
Joost van der Minnen	Rob Pirnie	Philip Schneek	David Theese	Romeo Zwart
Thomas Mino	Jorge Ivan Pincay	James Schneider	Douglas Thompson	廖明沂.
Rob Minshall	Ponce	Peter Schoo	Kerry Thompson	
Wijnand	Marc Vives Piza	Dan Schrenk	Lorin J Thompson	
Modderman-Lenstra	Victoria Poncini	Richard Schultz	Fabrizio Tivano	



Follow us on Twitter and Facebook

@protocoljournal



<https://www.facebook.com/newipj>

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

John Crain, Senior Vice President and Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Principal – Internet Technology, Policy, and Advocacy
The Internet Society

Dr. Jun Murai, Founder, WIDE Project
Distinguished Professor, Keio University
Co-Director, Keio University Cyber Civilization Research Center, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

September 2023

Volume 26, Number 2

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

FROM THE EDITOR

In This Issue

From the Editor	1
Introduction to 5G	2
LEO Satellites for Internet Access.....	31
Letter to the Editor.....	44
Fragments.....	45
Thank You!	48
Call for Papers.....	50
Supporters and Sponsors	51

Technologies used for accessing the Internet have evolved a great deal since my very first encounter with the ARPANET in 1976. Using a 110-baud Teletype machine, I accessed a computer at SRI International in Menlo Park, California, from my laboratory at the *Norwegian Defence Research Establishment* (NDRE) at Kjeller, Norway. Today, my Internet service is delivered by fiber-optic cable at 1 Gbps. Numerous other technologies for Internet access have emerged, and in this issue we look at two of them, namely 5G mobile systems and *Low Earth Orbit* (LEO) satellites.

In simple terms 5G can be described as a new set of cellular radio frequencies to allow for much faster data connections for mobile devices. Beyond this simplified explanation lies numerous details that are described in a two-part article by William Stallings. Part One introduces the standards and specifications that define 5G and describes the usage scenarios that 5G supports. Part Two, to be published in a future issue, will provide an overview of the structure and function of 5G networks. A third article on *Network Slicing*, which is closely related to 5G, will also be published in a future edition of this journal.

The world's first communication satellite, *Telstar 1*, which was launched in July 1962, provided proof-of-concept for both live television transmission and telephone service. Since that time, satellites have been deployed for many services, ranging from weather observations, navigation systems, and more recently Internet access. LEO satellites are particularly well-suited for Internet access because they offer coverage to remote areas without introducing substantial propagation delays as compared to other alternatives. Our second article, by Dan York and Geoff Huston, provides an overview of LEO systems for Internet access.

This journal now has around 1,000 print subscribers and just over 18,000 online subscribers who download their copy from our website. Given this shift in subscriber demographics, we will no longer be printing those long and cumbersome URLs in the references section at the end of each article. Instead, you can simply click on the references themselves using the PDF copy.

As always, we welcome your feedback and suggestions on anything you read in this journal. Letters to the Editor may be edited for clarity and length and can be sent to ipj@protocoljournal.org

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

Introduction to 5G

Part One: Standards, Specifications, and Usage Scenarios

by William Stallings

5G is the fifth-generation technology for wireless cellular networks. A significant technological leap beyond the capabilities of the 4G networks that currently dominate available cellular network services, 5G delivers a substantial increase in peak and average speeds and capacity. A significant increase in download and upload speeds will enhance many existing use cases, including cloud-based storage, augmented reality, and artificial intelligence. It also will enable cell sites to communicate with a greater number of devices. Reduced latency enables edge computing and will transform *Internet of Things* (IoT) capabilities and application breadth.

This two-part article provides an introduction to 5G. Part One introduces the standards and specifications that define 5G and describes the usage scenarios that 5G supports. Part Two, to be published in a future issue, will provide an overview of the structure and function of 5G networks.

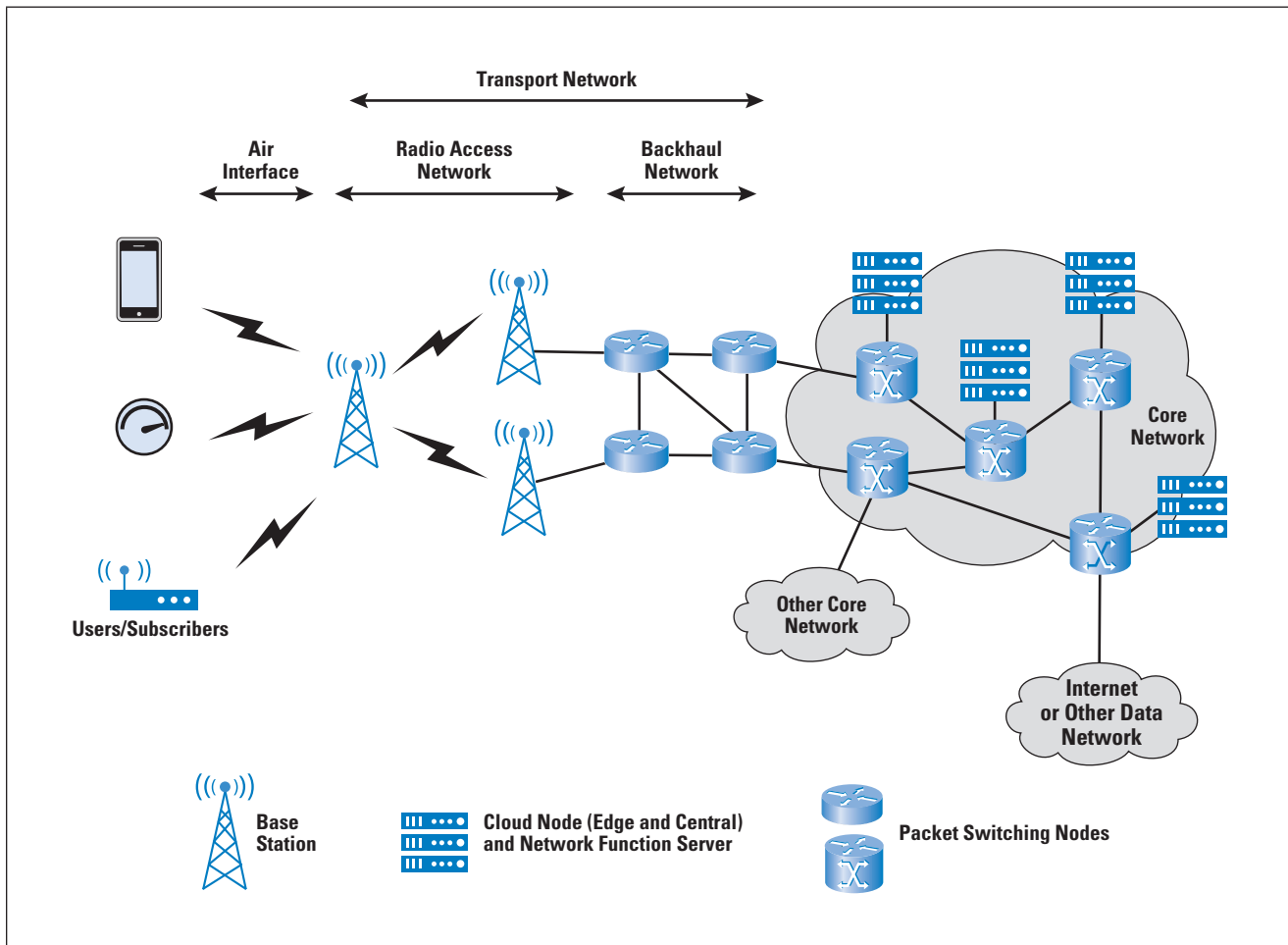
Simplified View of a 5G Network

Figure 1 shows a very simplified view of a 5G cellular network; it provides a useful framework for discussing standards and specifications for the technology. Users, or subscribers, to the network are fixed or mobile wireless devices, referred to in 5G documents as *User Equipment* (UE). Examples of fixed wireless UE include a modem that serves as an access point from a home or office Wi-Fi or Ethernet network to the 5G network, and IoT devices such as sensors or surveillance cameras. Mobile UE includes cell phones, laptops, and other mobile or portable devices equipped with 5G capability.

The *Air Interface*, also called a *Radio Interface*, is the wireless link between UE and the nearest cellular base station. The air interface specifies the method for transmitting information over the air between base stations and mobile units, including protocols, frequency range, channel bandwidth, channel coding, and the modulation scheme.

The *Radio Access Network* (RAN) consists of a collection of base stations that provide service to the UE in a geographic region. The base station provides radio transmission and reception in one or more cells to or from the user equipment. A base station can have an integrated antenna or can be connected to an antenna by feeder cables. Each base station communicates with nearby base stations, generally wirelessly, to enable handoff of UE from one base station to another as the UE moves. The RAN also includes other management and transmission elements.

Figure 1: Simplified View of 5G Network



The *Core Network* is a backbone network that provides interconnection service between RANs in different regions; it provides access for UE to the Internet or other data networks and UE on other RANs. In addition, the core network implements numerous network functions that support user- and control-plane traffic and provides for *Quality of Service* (QoS) and management and orchestration. Core networks also generally provide both edge and central cloud services for 5G users.

The *Transport Network* is the collection of communication links that interconnect nodes of the RAN, as well as communication links connecting RAN elements to the 5G core network. Links between RAN elements and UE are generally not considered part of the transport network.

Standards and Specifications for 5G

Many of the important developments in information and technology and communications, such as the Internet, IoT, Cloud Computing, and Virtualization, have been driven in part by international standards.

However, in all of these cases, much of the technology was developed and deployed in advance of universally agreed-upon standards. The case of 5G is quite different. Although a reasonably complete set of standards based on fixed specification is only just coming to fruition, the implementations and deployments that preceded these standards and specifications anticipated their final form. Throughout the 5G ecosystem, which includes device and component manufacturers, cellular network providers, network software providers, and application developers, the work done prior to the introduction of the first set of standards in 2020 closely follows what has ultimately been standardized. Going forward, there is universal agreement that 5G-related implementations will follow the standards.

Because an understanding of 5G depends on an understanding of the process by which the standards are developed and the content of those standards, the first part of this article provides an overview. It covers the two organizations that are responsible for the development of 5G: the *International Telecommunication Union* (ITU) and the *3rd Generation Partnership Project* (3GPP). In essence, the process of standards development for 5G follows this sequence:

1. The ITU has issued—and continues to issue—standards, called *Recommendations*, and other documents, called *Reports*, that define the overall concept for 5G, as well as the technical, performance, and service requirements for 5G.
2. Based on the ITU requirements, as well as requirements generated by national and regional standards organizations and market-based organizations, 3GPP has developed—and continues to develop—a detailed set of technical specifications for the implementation of 5G.
3. The ITU has translated these specifications into international standards (Recommendations) that dictate how 5G is being implemented.

This process is ongoing as further refinements and capabilities are added to the requirements and the technical specifications.

With respect to 5G, the two relevant components of ITU are the *ITU Radiocommunication* (ITU-R) Sector and the *ITU Telecommunication Standardization* (ITU-T) Sector. In general terms, ITU-R issues standards related to user requirements and the air interface. ITU-T issues standards related to the RAN, the transport network, and the core network.

ITU-R and IMT-2020

Perhaps the most prominent initiative by ITU-R is the *International Mobile Telecommunications* (IMT) project. IMT is the generic term the ITU community uses to designate broadband mobile systems. It encompasses *IMT-2000*, *IMT-Advanced*, and *IMT-2020* collectively, which correspond to 3G, 4G, and 5G, respectively.

A foundational document in the definition of IMT-2020 is *ITU-R Recommendation M.2083*^[1]. In broad strokes, this document develops a vision of the 5G mobile broadband connected society and future IMT. The two main contributions of this recommendation are a set of target values for key capabilities and a definition of usage scenarios, discussed subsequently.

M.2083 lists the eight key capabilities for IMT, together with the minimum requirements for each. The objectives that determined these target values follow:

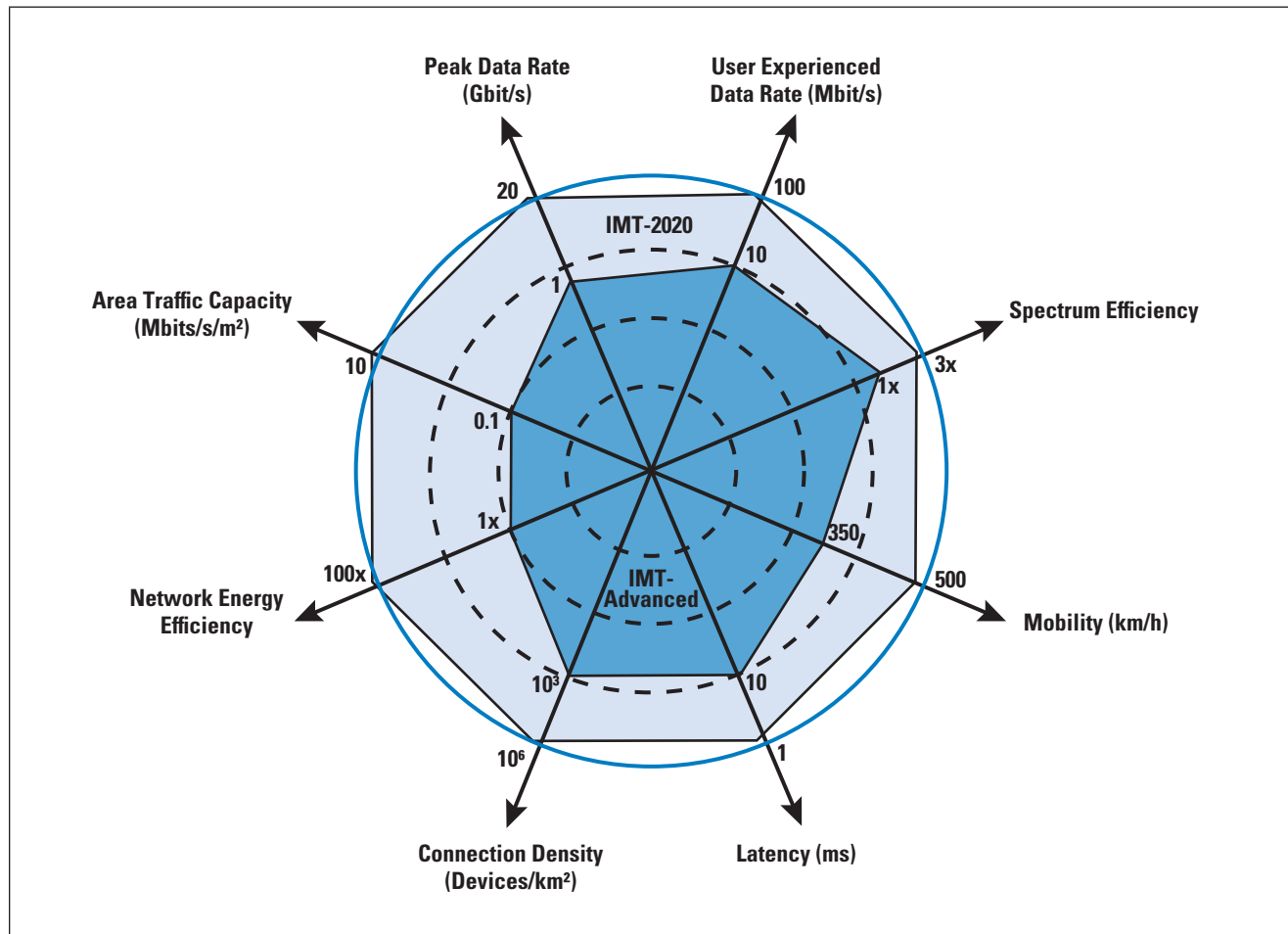
- The user experience with IMT-2020 using a mobile device should match—to the extent possible—the experience with fixed networks. The enhancement will be realized by increased peak and user experienced data rate, enhanced spectrum efficiency, reduced latency, and enhanced mobility support.
- IMT-2020 should support massive machine-to-machine interconnections, for a variety of IoT environments.
- IMT-2020 should be able to provide these capabilities without undue burden in terms of energy consumption, network equipment cost, and deployment cost to make future IMT sustainable and affordable.

Figure 2, from M.2083, compares the capability requirements of IMT-2020 (5G) to those of IMT-Advanced (4G). It shows that substantial improvements are mandated for all eight capabilities, with the most substantial required improvements in the areas of traffic capacity and network energy efficiency.

The target values were published in 2015, with the admonition that they are presented for purposes of research and development and may be revised in light of future studies and implementation experience. This list was expanded and refined in 2017 into 13 technical performance requirements in ITU-R Report M.2410^[2]. The purpose of these performance requirements is to assure that there should be a noticeable improvement of user *Quality of Experience* (QoE) for legacy 4G services and applications, and a high QoE for emerging 5G services and applications. Two terms should be distinguished:

- *Quality of Service* (QoS): The measurable end-to-end performance properties of a network service, which can be guaranteed in advance by a *Service-Level Agreement* (SLA) between a user and a service provider, so as to satisfy specific customer application requirements. Note: These properties may include throughput (bandwidth), transit delay (latency), error rates, priority, security, packet loss, and packet jitter.
- *Quality of Experience* (QoE): A subjective measure of performance in a system. QoE relies on human opinion and differs from QoS, which you can measure precisely.

Figure 2: Enhancement of Key Capabilities from IMT-Advanced to IMT-2020



In essence, the performance requirements for 5G are QoS measures designed to produce a high QoE. The M.2410 minimum technical performance requirements are as follows:

- **Peak Data Rate:** The maximum achievable data rate under ideal conditions per user/device (in Gbps). The minimum target values are downlink peak data rate of 20 Gbps and uplink peak data rate of 10 Gbps.
- **Peak Spectral Efficiency:** The maximum data rate under ideal conditions normalized by channel bandwidth (in bits/s/Hz). Another way of expressing this term is that it is the maximum data rate that can be transmitted over a given bandwidth. The relationship can be expressed as follows: $R_p = W \times SE_p$, where R_p is peak data rate, W is the available bandwidth, and SE_p is peak spectral efficiency. The minimum for peak spectral efficiencies is a downlink of 30 bps/Hz and uplink of 15 bps/Hz.
- **User-Experienced Data Rate:** The achievable data rate that is available ubiquitously across the coverage area to a mobile user/device (in Mbps or Gbps). This rate will depend the type of environment.

- *5th Percentile User Spectral Efficiency*: The 5% point of the cumulative distribution function of the normalized user throughput. The normalized user throughput is defined as the number of correctly received bits, that is, the number of bits contained in the *Service Data Units* (SDUs) delivered to Layer 3, over a certain period of time, divided by the channel bandwidth; it is measured in bits/s/Hz.
- *Average Spectral Efficiency*: The average data throughput per unit of spectrum resource and per cell (bits/s/Hz). The goal is a spectral efficiency of three times higher than IMT-Advanced.
- *Area Traffic Capacity*: The total traffic throughput served per geographic area (in Mbps/m²).
- *Latency*: Deals with transmission delays introduced by the network. Report M.2410 considers two types of latency:
 - *User-Plane Latency*: The contribution by the radio network to the time from when the source sends a packet to when the destination receives it (in ms).
 - *Control-Plane Latency*: Refers to the transition time from a most “battery efficient” state (for example, Idle state) to the start of continuous data transfer (for example, Active State). The minimum requirement is 20 ms.
- *Connection Density*: The total number of connected and/or accessible devices per unit area (per km²) that fulfills a specific QoS. The minimum requirement is 106/km².
- *Energy Efficiency*: In general terms, the relation between useful output and energy consumption. In the context of M.2410, this parameter has two aspects:
 - *Network Energy Efficiency*: Refers to the quantity of information bits transmitted to/received from users, per unit of energy consumption of the RAN (in bits/Joule). The objective is efficient data transmission when the load on the network is substantial. The energy consumption for the RAN of IMT-2020 should not be greater than for IMT-Advanced, while delivering the enhanced capabilities. The network energy efficiency should therefore be improved by a factor at least as great as the envisaged traffic capacity increase of IMT-2020 relative to IMT-Advanced.
 - *Device Energy Efficiency*: Refers to a quantity of information bits per unit of energy consumption of the communication module (in bits/Joule). The objective is low-energy consumption when no data is being sent or received.
- *Reliability*: The probability of successful transmission of a Layer 2/3 packet within a required maximum time, which is the time it takes to deliver a small data packet from the radio protocol Layer 2/3 SDU ingress point to the radio protocol Layer 2/3 SDU egress point of the radio interface at a certain channel quality.

- *Mobility*: the maximum speed at which a defined QoS and seamless transfer between radio nodes that may belong to different layers and/or radio-access technologies (multi-layer/-RAT) can be achieved (in km/h). The following classes of mobility are defined:
 - *Stationary*: 0 km/hr
 - *Pedestrian*: 0 to 10 km/hr
 - *Vehicular*: 10 to 120 km/hr
 - *High-Speed Vehicular*: 120 to 5000 km/hr
- *Mobility Interruption Time*: The smallest time delay the system supports, during which the end-user device cannot exchange packets with any base stations during transmissions. The mobility interruption time includes the time required to execute any RAN procedure, radio resource control signaling protocol, or other message exchanges between the mobile station and the RAN. The minimum requirement is 0 ms.
- *Bandwidth*: The maximum aggregated system bandwidth. The minimum requirement is 100 MHz.

ITU-T and IMT-2020 “Softwarization”

As mentioned previously, the role of ITU-T in defining requirements and developing standards for IMT-2020 is complementary to that of ITU-R. ITU-T; it specifies requirements for overall non-radio aspects of the IMT-2020 network, especially with respect to network operations and support of service requirements. ITU-T Recommendations cover the core, RAN, and transport networks. ITU-T Y.3101^[3] lists the following objectives with respect to IMT-2020:

1. Minimized dependency on access network technologies
2. Coping with traffic explosion in urban areas
3. Easy incorporation of future emerging services
4. Provision of a cost-efficient infrastructure
5. Expansion of the geographic reach of the network

The ITU-T approach to achieving these objectives depends in large part on the introduction of *network softwarization* in IMT-2020 network components. Y.3101 defines network softwarization as an overall approach for designing, implementing, deploying, managing, and maintaining network equipment and/or network components by software programming. It enables you to use modular network functions that you can deploy and scale on demand to accommodate various use cases easily and cost-efficiently.

Four aspects of network softwarization are important in 5G networks and are reflected in the ITU-T documents:

- *Software-Defined Networking* (SDN): An approach to designing, building, and operating large-scale networks based on programming the forwarding decisions in routers and switches with software from a central server. SDN differs from traditional networking, which requires configuring each device separately and relies on protocols that you cannot alter.

- *Network Functions Virtualization (NFV)*^[18]: The virtualization of compute, storage, and network functions by implementing these functions in software and running them on virtual machines. NFV decouples network functions, such as routing, firewalling, intrusion detection, and network address translation from proprietary hardware platforms and implements these functions in software. It uses standard virtualization technologies that run on high-performance hardware to virtualize network functions. It is applicable to any data-plane processing or control-plane function in both wired and wireless network infrastructures.
- *Edge Computing*: A distributed *Information Technology (IT)* architecture in which client data is processed at the periphery of the network, as close to the originating source as possible.
- *Cloud-Edge Computing*: A form of edge computing that offers application developers and service providers cloud computing capabilities, as well as an IT service environment, at the edge of a network. The aim is to deliver compute, storage, and bandwidth much closer to data inputs and/or end users.

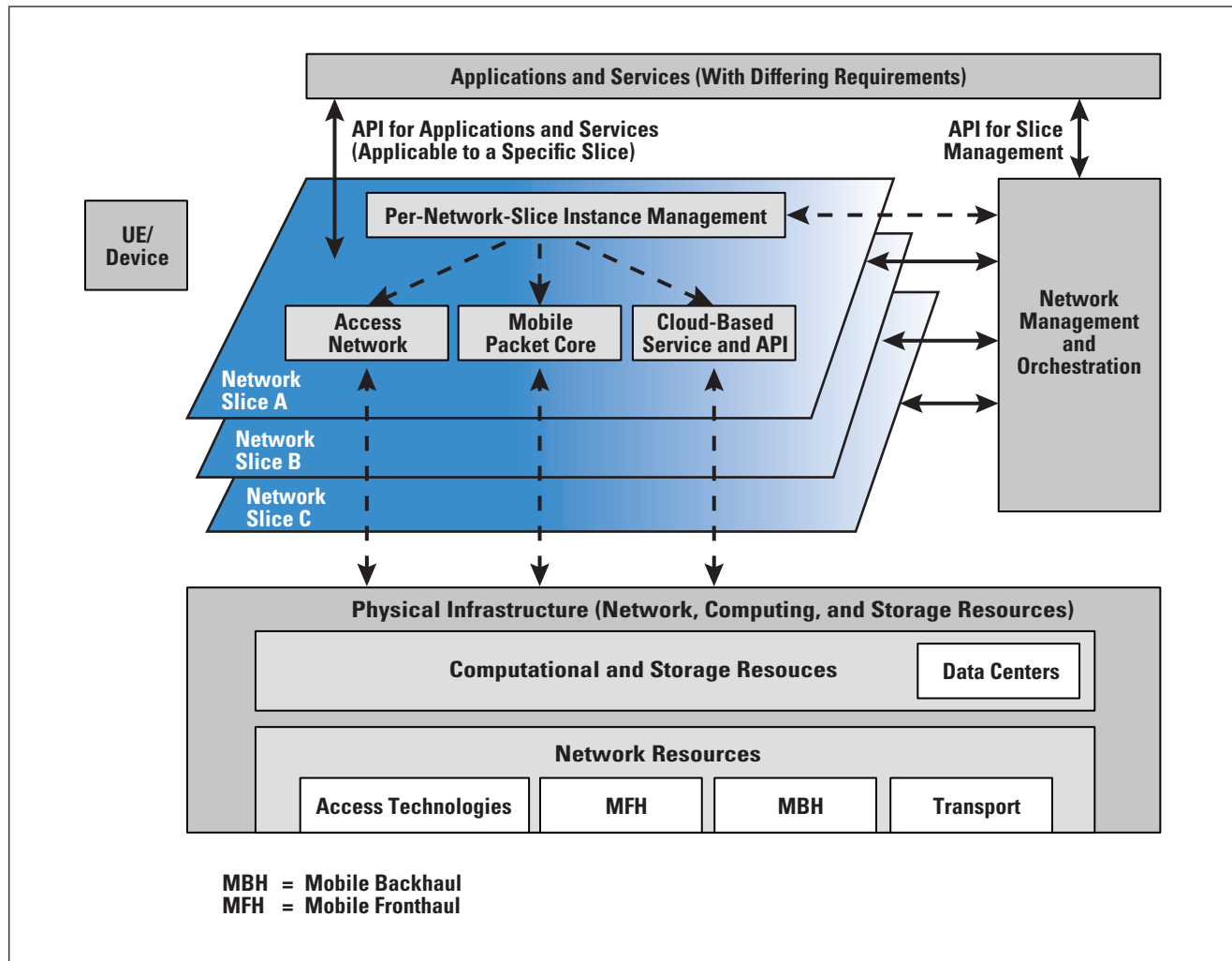
Network softwarization allows you to implement one of the essential features of 5G networks: *network slicing*. Network slicing permits you to separate a physical network into multiple virtual networks (logical segments) that can support different QoS requirements from applications and end users. Network slicing involves the selection and reservation of resources in the air interface, the RAN, the transport network, and the core network.

In essence, network slicing allows you to create multiple virtual networks atop a shared physical infrastructure. In this virtualized network scenario, physical components are secondary and logical (software-based) partitions are paramount, devoting capacity to certain purposes dynamically, according to your need. As your needs change, so can your devoted resources. Using common resources such as storage and processors, network slicing enables you to create slices devoted to logical, self-contained, and partitioned network functions. Network slicing supports the creation of virtual networks to provide a given QoS, such as guaranteed delay, throughput, reliability, and/or priority.

Figure 3, from ITU-T Recommendation Y.3150^[4], illustrates how network softwarization is incorporated in the design of IMT-2020 networks. The underlying physical infrastructure consists of a heterogeneous collection of network, computing, and storage resources. The figure shows four network resource categories. The access technologies consist of the resources at the air interface, including bandwidth, access protocol, channel coding, and modulation scheme. The mobile fronthaul refers to network paths between centralized radio controllers and remote radio units of a base-station function. The mobile backhaul refers to the network path between base-station systems and a core network. The transport resources consist of the switching hardware and software for routing data packets in the transport and core networks; an SDN controller manages these packets.

Using NFV, these underlying resources are abstracted to virtual resources used to create network slices, under the control of the management and orchestration function. Individual network slices can have specific characteristics that reflect various different requirements derived from application and services.

Figure 3: Network Softwarization for IMT-2020



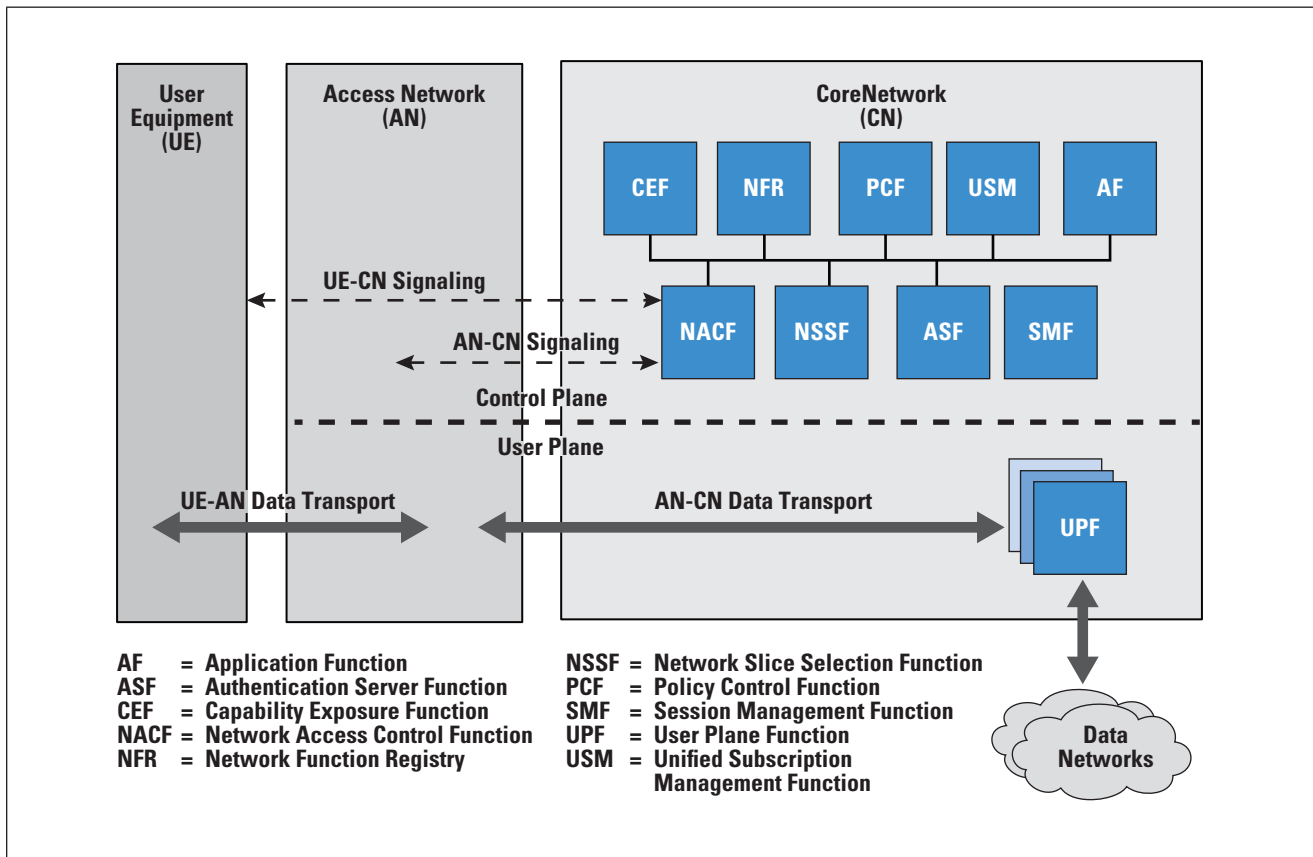
ITU-T IMT-2020 Core Network Framework

ITU-T Recommendation Y.3102^[5] provides the framework for overall non-radio aspects of the IMT-2020 network. Figure 4, from Y.3102, illustrates the interactions between the network functions for providing network service.

The framework delineates three domains. The *User Equipment* (UE) domain consists of devices that transmit and receive data over the IMT-2020 network. The *Access Network* (AN) domain is the wireless connection between the UE and the *Core Network* (CN), defined by the ITU-R radio interface recommendations.

The framework diagram also depicts the division between a control plane and a user plane, which cuts across the AN and CN.

Figure 4: Framework of the IMT-2020 Network



The *Control Plane* performs the call and connection control functions. For this purpose, a signaling connection between the UE and the CN exchanges signaling messages that manage the signaling connection and the call established for the UE. The Control-Plane functions are requested and managed via control signals that are exchanged between UE and the AN, and between the AN and the CN. Through signaling, the control plane sets up and releases connections, and may restore a connection if a failure occurs. The control plane also performs other functions that support call and connection control, such as routing information dissemination.

The core network includes the following functional elements:

- *Network Access Control Function (NACF)*: Provides access to the CN services for the AN and UE. NACF includes:
 - *Registration Management*: Enables UE to register for network access. NACF performs, but is not limited to, network slice instance selection, UE authentication, authorization of network access and network services, and network access policy control.
 - *Connection Management*: Establishes and releases a signaling connection between the UE and the core network.

- *Session Management Function Selection*: Determines the session management function that is most appropriate to establish and manage a session. In the context of IMT-2020, a session is an association between UE and a data network that provides a *Protocol Data Unit* (PDU) connectivity service.
- *Session Management Function* (SMF): Sets up and manages one or more sessions that provide connectivity between the local UE and a remote UE. This function deals with user path selection and enforcement of policies, including QoS policy and charging policy.
- *Policy Control Function* (PCF): Provides for control and management of policy rules.
- *Capability Exposure Function* (CEF): Enables the exposure of network functions and network slices as a service to third parties.
- *Network Function Registry Function* (NRF): Assists the discovery and selection of required network functions.
- *Unified Subscription Management Function* (USM): Stores and manages UE context and subscription information including, but not limited to, UE information on registration and mobility management, information on network functions that serve the UE, and information on session management. USM also provides UE authentication information to the *Authentication Server Function* (ASF).
- *Network Slice Selection Function* (NSSF): When UE requests registration with the network, NSSF sends a network slice selection request to NSSF with preferred network slice selection information. The NSSF responds with a message including the list of appropriate network slice instances for the UE.
- *Authentication Server Function* (ASF): Performs authentication between UE and the network.
- *Application Function* (AF): Interacts with application services that require dynamic policy control. AF extracts session-related information (for example, QoS requirements) from application signaling and provides it to PCF in support of its rule generation.

The user plane refers to the set of traffic forwarding components through which traffic flows. Its principal function is to provide transfer of end-user information.

The sole functional element in the user plane is the *User Plane Function* (UPF). This function includes traffic routing and forwarding, *Protocol Data Unit* (PDU) session tunnel management, and QoS enforcement. The PDU session tunnels are used between AN and UPF(s) as well as between different UPFs as user-plane data transport for PDU sessions. UPF also provides optional functions including packet inspection and collection of *User-Plane* (UP) traffic for lawful intercept. In order to accommodate the diversity of network scenarios, UPF may also provide interworking functions among different network segments, for example, interworking between the IP-based core network and the non-IP-based access network.

Y.3102 also lists the primary network services that the supported core network framework supports. They include:

- *Registration Management* (RM): Register or deregister UE with the IMT-2020 network and establish the user context in the network.
- *Connection Management* (CM): Establish and release the signaling connection between the UE and NACF.
- *Session Management* (SM): Manages PDU sessions including control of PDU session tunnel establishment, modification, and release.
- *User-Plane Management* (UPM): Forward user traffic, including user traffic rerouting between UPFs because of the serving UPF relocation and enforcement of QoS policies.
- *Mobility Management* (MM): Used to handle all aspects related to UE mobility. Mobility management aspects include, but are not limited to, UE reachability and handover management.

3GPP

The *3rd Generation Partnership Project* (3GPP) was formed in 1998 by a global consortium of regional *Standards Development Organizations* (SDOs) to develop technology specifications for 3G cellular networks. Because it involved the efforts of the world's leading national standards organizations, 3GPP became the dominant agent in the development of specifications for 3G, then 4G, and now 5G cellular networks.

3GPP began work in 2016 on defining 5G technical specifications for a new radio access technology, known as 5G NR (New Radio) and a next-generation network architecture (5G NextGen). Unlike previous generations, competing standards bodies are no longer working on potential solutions for 5G.

Figure 5 shows the key players in the 3GPP process and their relationships to one another. Within the 3GPP organization is a *Project Coordination Group* (PCG). It is responsible for overall time frame and management of technical work to ensure that the 3GPP specifications are produced in a timely manner as required by the marketplace. Subordinate to the PCG are three *Technical Specification Groups* (TSGs). Each TSG has the responsibility to prepare, approve, and maintain the specifications within its terms of reference; it may organize its work in *Working Groups* (WGs) and liaise with other groups as appropriate. The TSGs report to the PCG.

Key to the 3GPP process are the *organizational partners*. An organizational partner is a standards organization with a national, regional, or other officially recognized status (in its country or region) that has the capability and authority to publish standards nationally or regionally. Associated with organizational partners are individual members, which are member companies affiliated with one of the organizational partners.

Finally, there are *market-representation partners*, which are organizations invited to participate by the organizational partners to offer market advice to 3GPP and to bring into 3GPP a consensus view of market requirements (for example, services, features, and functions) falling within the 3GPP scope.

Figure 5: 3GPP Process

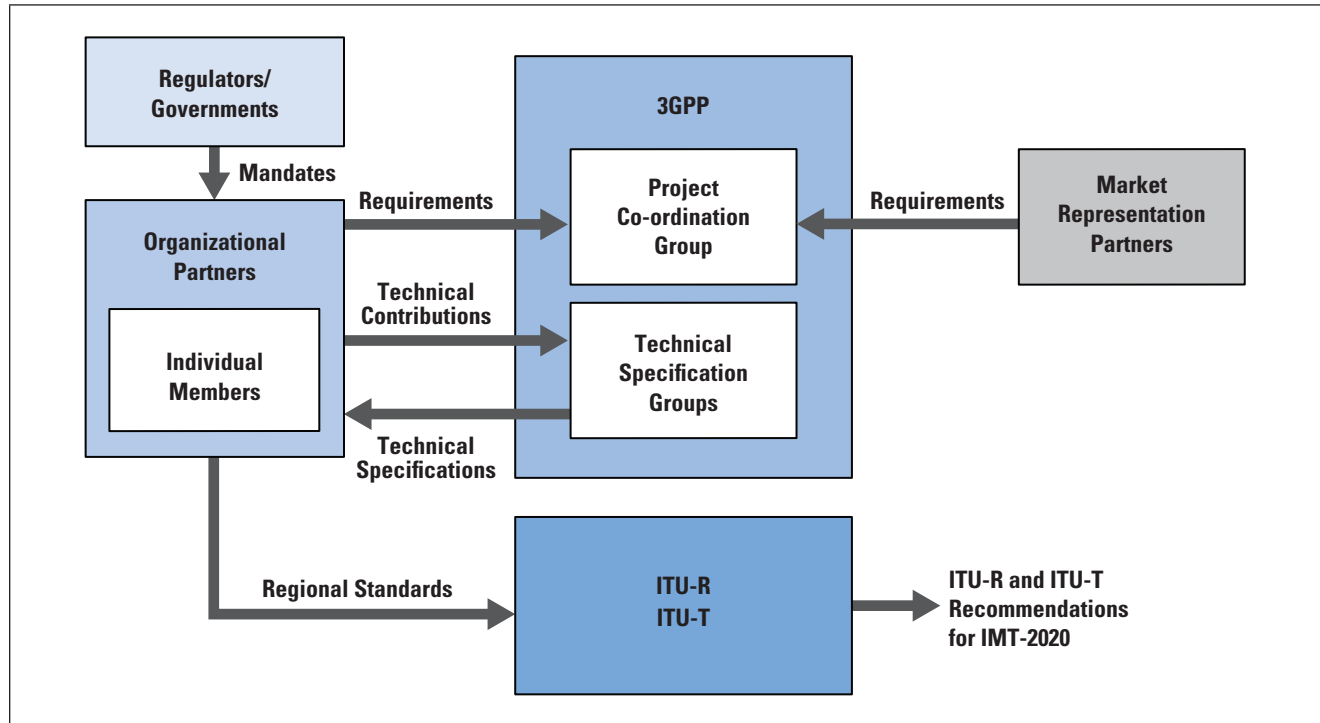


Figure 5 shows, in general terms, the flow of information between the previously mentioned entities. The PCG plans the work of 3GPP based on requirements provided by the organizational partners and the market representation partners. The organizational partners are influenced particularly by their respective national and regional governments and regulators, whereas the market representation partners generate requirements dictated by the potential market. Individual members provide technical contributions to the TSGs, which ultimately result in technical specifications. These specifications are transmitted from the TSGs to the organizational partners, who translate them into national and regional standards. Finally, these standards serve as input to ITU in the development of 5G-related Recommendations.

3GPP Releases

3GPP uses a system of parallel *Releases* that provide developers with a stable platform for the implementation of features at a given point and then allow for the addition of new functions in subsequent Releases. Releases are staggered and work is done on multiple Releases in parallel at different stages. When a Release is finalized, it means that all new features are functionally frozen and ready for implementation. Furthermore, each 3GPP Release is self-contained, meaning that you can build a cellular system based on the set of frozen specifications in that Release.

As such, Releases do not just contain the newly implemented features, but instead are introduced in a highly iterative manner that builds upon previous Releases. Table 1 provides information on the three releases relating to 5G that are completed at the time of this writing. Release 15 provided an early definition of useful 5G features to enable deployment by 2020. Subsequent releases add progressively more functions. Release 16 should closely resemble the initial set of IMT-2020 Recommendations issued by ITU in 2020.

Table 1. 3GPP Releases for 5G

Release #	Status	Functional Freeze	End Date
Release 17	Frozen	2022-03-18	2022-06-10
Release 16	Frozen	2020-07-03	2020-07-03
Release 15	Frozen	2019-03-22	2019-06-07

When a Release is frozen, the TSGs can add no additional functions to the specifications. However, detailed protocol specifications may not yet be complete. The end date shown in Table 1 is indicative only, since for each Release, a considerable number of refinements and corrections can be expected for at least two years following this date.

3GPP Requirements for 5G

The 3GPP documents include a description of 5G requirements that are significantly more detailed than those provided in the ITU documents. As such, they provide an important guide to implementers of 5G networks, components, and systems as to what the market requirements are for 5G success (Refer to Figure 6).

Figure 6: 3GPP Basic Capability Requirements

Network Slicing Diverse Mobility Magement Multiple Access Technologies Resource Efficiency Efficient User Plane Efficient Content Delivery Priority, QoS, and Policy Control Dynamic Policy Control Connectivity Models Network Capability Exposure Context Aware Network Self-Backhaul Flexible Broadcast/Multicast Service	Subscription Aspects Energy Efficiency Markets Requiring Minimal Service Levels Extreme Long-Range Coverage in Low-Density Areas Multi-Network Connectivity and Service Delivery Across Operators 3GPP Access Network Selection eV2X Aspects NG-RAN Sharing Unified Access Control QoS Monitoring Ethernet Transport Services	Non-Public Networks 5G LAN-Type Service Positioning Services Cyber-Physical Control Applications in Vertical Domains Messaging Aspects Steering of Roaming Minimization of Service Interruption UAV Aspects Video, Imaging, and Audio for Professional Applications Critical Medical Applications
--	---	--

eV2X = Enhanced Vehicle-to-Everything

UAV = Unmanned Aerial Vehicle

3GPP Technical Specification TS 22.261^[6] defines requirements for 34 basic capabilities to be provided by a 5G network; they are listed in Figure 6. For each capability, TS 22.261 provides a description and elaborates on the requirements for that capability.

TS 22.261 also lists performance requirements that are more detailed and more demanding than those defined in ITU-R Report M.2410. The requirements cover the following categories:

- *High Data Rates and Traffic Densities*: Several 5G scenarios require the support of very high data rates or traffic densities, including urban and rural areas, office and home, and special deployments (for example, massive gatherings, broadcast, residential, and high-speed vehicles).
- *Low Latency and High Reliability*: Some scenarios require the support of very low latency and very high communications service availability, which in turn implies very high reliability. The overall service latency depends on the delay on the radio interface, transmission within the 5G system, transmission to a server that may be outside the 5G system, and data processing. Some of these factors depend directly on the 5G system itself, whereas for others the impact can be reduced by suitable interconnections between the 5G system and services or servers outside of the 5G system, for example, to allow local hosting of the services. TS 22.261 provides an overview of potential scenarios and references other technical specifications for specific requirements.
- *High Accuracy Positioning*: The 5G System shall provide different 5G positioning services with configurable performances working points (for example, accuracy, positioning service availability, positioning service latency, energy consumption, update rate, and time to first fix) according to the needs of users, operators, and third parties. TS 22.261 lists quantitative requirements for numerous indoor and outdoor scenarios.
- *Key Performance Indicators (KPIs) for a 5G System with Satellite Access*: In some contexts, a 5G access network will use at least one satellite link. KPIs defined in TS 22.261 include minimum and maximum UE-to-satellite delay for various earth orbits, as well as maximum propagation delay.
- *High Availability IoT Traffic*: This requirement is concerned specifically with medical monitoring but is applicable to other scenarios that require highly reliable machine-type communication in both stationary and highly mobile settings.
- *High Data Rate and Low Latency*: This requirement defines data and latency requirements for such scenarios as audio-visual interaction, gaming, and virtual reality.
- *KPIs for UE-to-Network Relaying in 5G System*: In several scenarios, it can be beneficial to relay communication between one UE and the network via one or more other UEs. This category includes performance requirements for various scenarios.

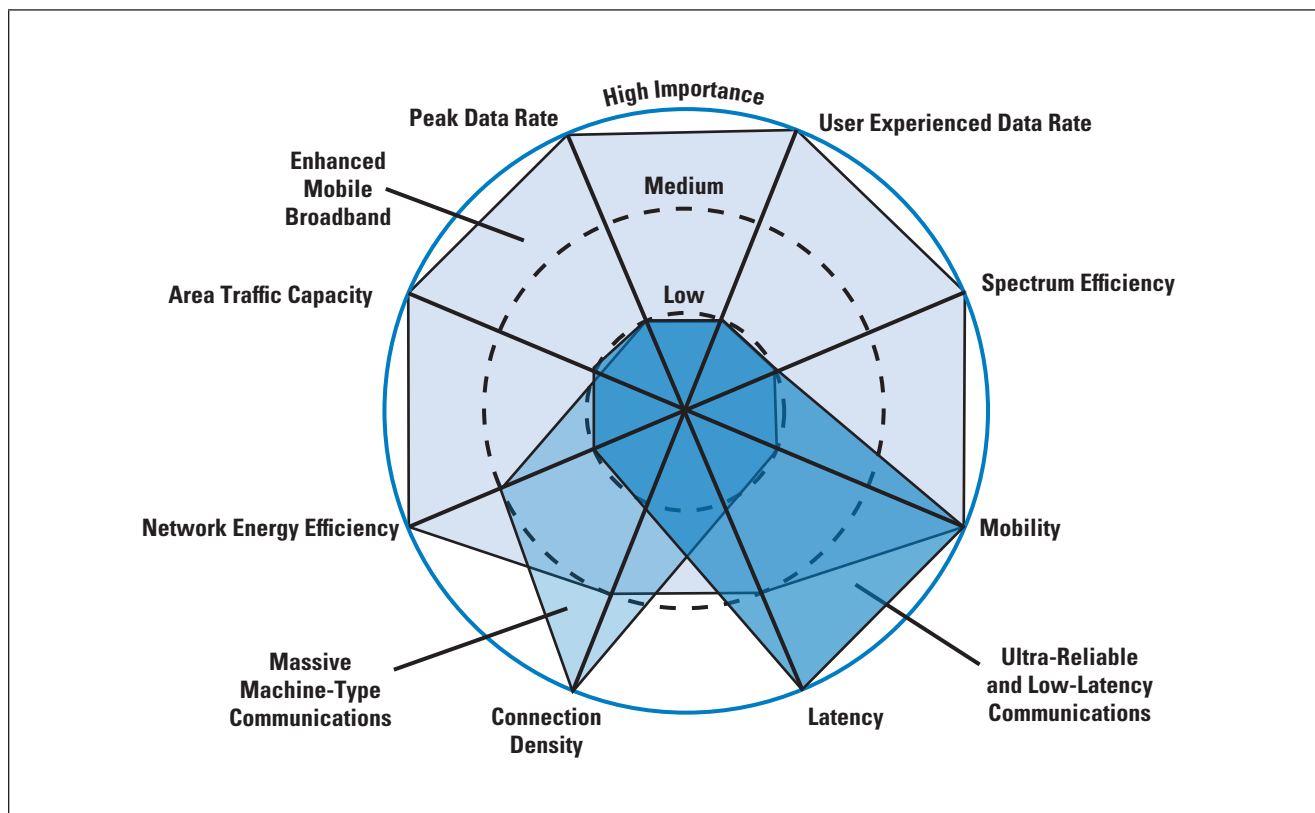
Usage Scenarios and Use Cases

Two important concepts in ITU-R Recommendation M.2083 and related documents are *Usage Scenario* and *Use Case*. No ITU document defines these terms, but the following definitions should suffice for this article.

- *Usage Scenario*: A general description of a way in which an IMT network is used. A usage scenario dictates various performance and technical requirements. A wide but nevertheless constrained variety of use cases are encompassed by a usage scenario.
- *Use Case*: A specific application or way of using an IMT network; also, a general account of a situation or course of actions that use an IMT network. It is described from the end user's perspective and illustrates fundamental characteristics. A use case dictates more specific and refined performance and technical requirements than the corresponding usage scenario.

M.2083 defines three usage scenarios: enhanced mobile broadband, massive machine-type communications, and ultra-reliable and low-latency communications. Figure 7, from M.2083, indicates the relative importance of the key capabilities for the three usage scenarios.

Figure 7: The Importance of Key Capabilities in Different Usage Scenarios



Enhanced Mobile Broadband

Enhanced Mobile Broadband (eMBB) is the 5G feature that provides a significant increase in data rate over 4G for a normal mobile Internet user. Enhanced mobile broadband services allow users to experience high-speed and high-quality multimedia services such as virtual reality, *Augmented Reality (AR)*, and 4,000-pixels horizontal resolution (4K) video, at any time and place.

These applications require reasonably low-latency and good connection density, with a high demand on the other six key capabilities. In addition to the consumption of multimedia content for entertainment purposes, eMBB supports numerous business applications. They include cloud access apps for commuters and other off-site employees, the ability of remote workers to communicate with the back office, or indeed an entire smart office where all devices are wirelessly and seamlessly connected.

Of the three usage scenarios defined in M.2083, eMBB is the only general-purpose case, and it is the one that is most familiar to current 4G users. In essence, eMBB is an enhanced version of 4G, providing improved performance and an increasingly seamless user experience.

ITU-R Report M.2410 lists three deployment options that characterize the scope of eMBB and that are used for purposes of evaluation of candidate specifications: indoor hotspot, dense urban, and rural. The remainder of this section provides a brief overview of all three.

ITU-R Report M.2412^[7] defines *Indoor Hotspot* as “...an indoor isolated environment at offices and/or in shopping malls based on stationary and pedestrian users with very high user density.” This deployment scenario focuses on small coverage per site/*Transmission and Reception Point* (TRxP) and high user throughput or user density in buildings. The key characteristics of this deployment scenario are high capacity, high user density, and consistent user experience indoors.

5G capabilities should enable a seamless interface for users moving into and out of the indoor zone, without the necessity of joining a Wi-Fi network for indoor use. Types of demand include frequent upload and download of data from a company’s servers and real-time video meetings with local as well as remote participants.

One of the main challenges for supporting 5G use cases in the indoor environment is a consequence of the use of much higher-frequency bands for 5G than are used for 4G and earlier generations. These higher bands lead to greater link losses. For example, outdoor signals on the C band will be subject to an 8- to 13-dB link loss when penetrating through one concrete wall. The signals on the higher-mm wave band will experience difficulty in penetrating through a wall as the link loss exceeds 60 dB. It is a considerable challenge for outdoor 5G macro signals to cover indoor areas, and a dedicated 5G network consisting of interconnected base stations will be required for indoor environments.

An example use case in the indoor hotspot category is the smart office. The installation of 5G networks in the office environment can enable dramatic changes in the capabilities that businesses can exploit.

Examples of features now in use or that may soon be in use in 5G-enabled workplaces include:

- Facial recognition can be used for entrance security. The employee need not carry an identification tag or use some sort of token to gain entrance.
- A 5G virtual desktop infrastructure enables workers to connect their mobile device on a docking pad to a cloud computing system.
- Workers can convene remote conferences, talking to each other's avatars in cyberspace.
- Security systems can use high-definition video to monitor in greater detail and expand the ability to scan for security threats.
- Workers have faster access to a broader selection of apps.
- 5G enables real-time collaboration between people and things, possibly including augmented reality features.
- Real-time video interaction will become standard. This access allows capabilities such as real-time troubleshooting and ad hoc meetings.
- Synchronization of local data with the cloud becomes almost instantaneous, further enhancing collaboration.
- Sensors or facial recognition can tell if people are in the building and where they are at any given moment.

In essence, the smart office use case is characterized by heavy data use, with a particular reliance on high-definition video, in an indoor environment with low mobility requirements. In this use case scenario hundreds of users require ultra-high bandwidth to serve intense bandwidth applications. To some extent, Wi-Fi supports these capabilities, but with the increasing demands for high traffic volume, high density of users, and seamless integration of local and wide-area communications, a unified 5G solution has inherent advantages over a mixed Wi-Fi/cellular environment.

ITU-R Report M.2412^[7] defines *Dense Urban* as "...an urban environment with high user density and traffic loads focusing on pedestrian and vehicular users." The dense urban microcellular deployment scenario focuses on macro TRxPs with or without micro TRxPs and high user densities and traffic loads in city centers and dense urban areas. The key characteristics of this deployment scenario are high traffic loads and outdoor-to-outdoor and outdoor-to-indoor coverage.

The dense urban environment for 5G is characterized by the use of a dense collection of small cells to supplement macro cells for two reasons^[8]:

- The concentrated collection of stationary, pedestrian, and vehicular users, with 5G use cases, generates a tremendous traffic load.
- 5G-mm Wave networks are predominantly noise-limited. The result is that only small cell sizes can be supported.

An example use case in this category is provided by the EU project METIS (*Mobile and wireless communications Enablers for the Twenty-twenty Information Society*)^[9]. It refers to the connectivity and data rates required for users of high-volume services at any place and at any time in a dense urban environment, including both user interaction with cloud services and data- and device-centric services.

5G enables enhanced cloud services beyond the traditional services of web browsing, file download, and social media. Enhanced services include high-definition video streaming and video sharing. Enhanced device-centric services include augmented reality with information fetched from sensors, smart phones, wirelessly connected cameras, and other sources. The main features of this use case follow:

- High traffic loads
- Low mobility
- High data rate
- Outdoor coverage
- Outdoor-to-indoor coverage
- Support for both low and high frequency
- Limited interference
- High user density

This use case presents two unique challenges:

- Users expect the same QoE in any context, including at their workplace, enjoying leisure activities such as shopping or being on the move, on foot, or in a vehicle.
- Users in urban environments tend to dynamically cluster. Examples include people waiting at a traffic light or bus stop and conference room meetings at the workplace. These clusters lead to sudden peaks of geographically concentrated mobile broadband demand.

ITU-R Report M.2412^[7] defines *Rural-eMBB* as “...a rural environment with larger and continuous wide area coverage, supporting pedestrian, vehicular and high-speed vehicular users.” The rural deployment scenario focuses on larger and continuous coverage. The key characteristics of this scenario are continuous wide area coverage supporting high-speed vehicles. This scenario uses macro TRxPs, and is noise- and/or interference-limited.

The rural deployment also supports last-mile service to residences and other subscribers to provide telephone and Internet access. Many homes may be near a fiber connection, but the deployment of the last mile of the cabling can be very expensive and not necessarily cost-effective. The addition of new subscribers, or households, may be very expensive if new cables need to be installed. It may also require the operator to support two distinct systems, each with its own subscription management, for wired and wireless subscribers. To address this problem, delivering the last mile wirelessly may be a viable option. Such solutions are known as *Wireless Local Loop* (WLL), where the last mile is delivered wirelessly.

Massive Machine Type Communications

Massive Machine Type Communications (mMTC) is characterized by a very large number of connected devices typically transmitting a relatively low volume of non-delay-sensitive data. However, the machine-to-machine communications involves a range of performance and operational requirements. Devices are required to be low-cost and have a very long battery life, such as five years or longer.

The mMTC usage scenario defined by ITU-R represents a subset of the total IoT universe. A white paper from Ericsson^[10] lists four segments that comprise IoT:

- *Massive IoT*: Massive IoT is characterized by huge volumes of constrained devices that send and/or receive messages infrequently. The traffic is often tolerant of delay. Examples of use cases include low-cost sensors, meters, wearables, and trackers. Such devices are often deployed in challenging radio conditions such as in the basement of a building. Therefore, they require extended coverage and may rely solely on a battery power supply that puts extreme requirements on the life of the battery.
- *Broadband IoT*: Broadband IoT is an application of eMBB to the IoT environment, providing high data rates and relatively low latencies. Examples of use cases are in the areas of automotive, drones, *Augmented Reality/Virtual Reality* (AR/VR), utilities, manufacturing, and wearables.
- *Critical IoT*: Critical IoT is an application of *Ultra Reliable and Low Latency Communications* (URLCC) to the IoT environment, providing extremely low latencies and ultra-high reliability at a variety of data rates. In contrast to Broadband IoT, which achieves low latency on best effort, critical IoT is intended to deliver data within strict latency bounds with required guarantee levels, even in heavily loaded networks. Examples of use cases are in the areas of intelligent transportation systems, smart utilities, remote healthcare, smart manufacturing, and fully immersive AR/VR.
- *Industrial Automation IoT*: This segment supports seamless integration of cellular connectivity into the wired industrial infrastructure used for real-time advanced automation. These applications have extremely demanding requirements such as very accurate indoor positioning and time synchronization across devices and networks.

Massive IoT, as defined by Ericsson, is equivalent to mMTC defined by ITU-R. In terms of the number of connections, mMTC is the most rapidly growing segment of IoT^[11]. Table 2, based on a 2020 Ericsson white paper^[12], indicates likely mMTC use cases that 5G supports.

Table 2: Industry and Society Applications Enabled by Massive IoT

Application Area	Use Cases
Transport and Logistics	Fleet Management Goods Tracking
Agriculture	Climate / Agriculture Monitoring Livestock Tracking
Environment	Process Monitoring and Control Maintenance Monitoring
Industrial	Process Monitoring and Control Maintenance Monitoring
Utilities	Smart Metering Smart Grid Management
Smart Cities	Parking Sensors Smart Bicycles Waste Management Smart Lighting
Smart Buildings	Smoke Detectors Alarm Systems Home Automation
Consumers	Wearables Children/Elderly Tracking Medical Monitoring

An important group of mMTC use cases are in the general category of smart cities. ITU-T 4900^[13] defines a *Smart Sustainable City*, or simply *Smart City*, as follows: A smart sustainable city is an innovative city that uses *Information and Communication Technologies* (ICTs) and other means to improve quality of life, efficiency of urban operation and services, and competitiveness, while ensuring that it meets the needs of present and future generations with respect to economic, social, environmental, and cultural aspects.

The sustainability of a smart city is based on four main aspects:

- *Economic*: The ability to generate income and employment for the livelihood of the inhabitants.
- *Social*: The ability to ensure that the welfare (safety, health, education) of the citizens can be equally delivered despite differences in class, race, or gender.
- *Environmental*: The ability to protect future quality and reproducibility of natural resources.
- *Governance*: The ability to maintain social conditions of stability, democracy, participation, and justice.

Some smart-city use cases, such as *Public Protection and Disaster Relief* (PPDR), fit into the URLCC usage scenario, but many others fall into the mMTC usage scenario category.

As examples, *ITU-T Series Y Supplement 56*^[14] lists the following mMTC demonstration examples of smart-city use cases:

- *Pedestrian Monitoring for Decisive Disaster Response*: Involves the installation of surveillance cameras throughout a city that can monitor crowd size and behavior and transmit this information to a central monitoring/management source. The cameras monitor locations that are likely to draw large groups of people, such as near a railroad or subway station or near a school. If a disaster occurs near one of these sites, the system provides real-time information about the size of the crowd at risk. In addition, the pedestrian monitoring system facilitates the understanding of the behavior of crowds and the detection of abnormal situations by analyzing images captured by surveillance cameras. If it detects any abnormality, the system automatically provides information or instructions for evacuation from the disaster site or for prevention of accidents.
- *Citizens' Safety Services*: Involves interworking between smart-city operation centers and fire and police stations for the citizens' safety services. Surveillance cameras are deployed throughout the city to provide extensive coverage with a minimum number of cameras. IoT-enabled traffic sensors deployed throughout the city can measure rate and volume of traffic. The operations center connects wirelessly to the cameras and sensors to provide a central source of information. The operations center provides traffic information to the first responders to enable them to take the best route to the scene of an emergency.
- *Lift Monitoring Services*: Involves monitoring lifts, or elevators, throughout a city. One such system developed by Surbana Jurong is deployed in Singapore and other Asian cities. The system consists of a central *Lift Monitoring System* (LMS) and IoT-enabled sensors installed in lifts throughout the city. The installation in Singapore monitors more than 26,000 lifts across 10,000 housing units. The system enables rapid response to elevator malfunction. In addition, the sensor devices capture data on an ongoing basis. This data, using machine-learning algorithms, is used to predict future failures, allowing for optimized maintenance and reduced downtime.
- *Infrastructure Monitoring*: Involves using IoT sensor devices to monitor aging infrastructure elements to support automated inspection, diagnosis, confirmation of repair effort, and subsequent status check. The scheme can be applied to bridges, tunnels, and paved roads.
- *Citizen Identification System Using Biometric*: Has objective to provide a digital identity to the entire population to serve as the basis for accessing social services and interacting with the government at various levels. One such system, called *Aadhaar*, is deployed nationwide in India and currently has over one billion people registered. In any large Indian city, there are tens or even hundreds of thousands of Aadhaar devices in use for registration and service access. These devices form a massive IoT network connected to a central server.

Ultra-Reliable and Low-Latency Communications

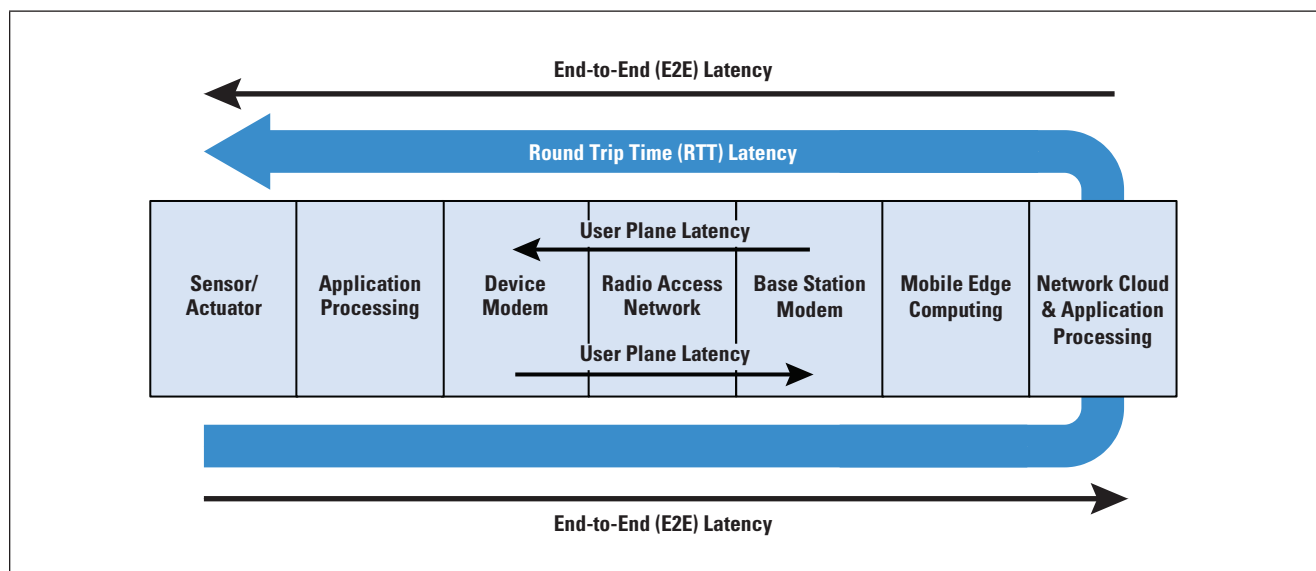
URLLC is a form of machine-to-machine communications that enables delay-sensitive and mission-critical services that require very low end-to-end delay, such as tactile Internet, remote control of medical or industrial robots, driverless cars, and real-time traffic control.

Figure 7 (earlier in this article) indicates that two parameters are of high importance for URLLC: *latency* and *mobility*. ITU-R Report M.2410 breaks the latency requirement into two parts:

- *User-Plane Latency*: is the contribution by the radio network to the time from when the source sends a packet to when the destination receives it (in ms). It is defined as the one-way time it takes to successfully deliver an application layer packet/message from the radio protocol Layer 2/3 SDU ingress point to the radio protocol Layer 2/3 SDU egress point of the radio interface in either uplink or downlink in the network for a given service in unloaded conditions, assuming the mobile station is in the active state. The minimum requirement (that is, the maximum allowable value) is 1 ms assuming unloaded conditions (that is, a single user) for small IP packets (for example, 0-byte payload + IP header), for both downlink and uplink.
- *Control-Plane Latency*: refers to the transition time from a most battery efficient state (for example, Idle state) to the start of continuous data transfer (for example, Active state). The minimum requirement is 20 ms.

User-Plane latency, however, is only one component that UE experiences overall, as illustrated in Figure 8. The *End-to-End (E2E)* latency is generally defined as the time it takes from when a data packet is sent from the transmitting end to when it is received at the receiving entity; for example, Internet server or other device. The measurement reference is the interface between Layers 2 and 3. It is also referred to as *One-Trip Time (OTT)*. It includes the user-plane latency in one direction, transport network delays, and application processing time.

Figure 8: E2E Latency and Round-Trip Time Latency



A related measure is *Round-Trip Time* (RTT), which is the time from when a data packet is sent from a source device until an acknowledgment or response is received from the destination device. Unfortunately, E2E latency is sometimes equal to RTT latency in the literature, even in some 3GPP documents. However, the implication in most standards and specification documents is that E2E latency refers to one-way latency, not round-trip.

Mobility is the maximum UE speed (in km/h) at which a QoS can be achieved. Mobility assumes a seamless transfer between radio nodes that may belong to different layers and/or radio access technologies (multi-layer/RAT) can be achieved. The following classes of mobility are defined:

- *Stationary*: 0 km/h
- *Pedestrian*: 0 to 10 km/h
- *Vehicular*: 10 to 120 km/h
- *High-speed vehicular*: 120 to 500 km/h

M.2410 does not provide a specific measure of QoS. Report ITU-R M.2412 defines QoS as successful delivery of 99% of messages within 10s.

Another aspect of mobility addressed in M.2410 is *Mobility Interruption Time*, which is the shortest time duration supported by the system during which UE cannot exchange user-plane packets with any base station during transitions. This number includes the time required to execute any RAN procedure, radio resource control signaling protocol, or other message exchanges between the mobile station and the RAN. The minimum requirement for mobility interruption time is 0 ms. Thus, there should be no interruption of service when moving UE switches from one base station to another.

URLLC Use Cases

A URLLC white paper from 5G Americas^[15], one of the 3GPP market representation partners, provides a useful way of understanding the wide variety of URLLC use cases by focusing on emerging mission-critical applications that have demanding reliability and latency requirements. These use cases include:

- Smart Factory
- Ground Vehicles, Drones, and Robots
- Tactile Interaction
- Augmented Reality and Virtual Reality
- Emergency, Disasters, and Public Safety
- Urgent Health Care
- Intelligent Transportation

The area that has perhaps received the most attention as an application area that requires URLLC support is that of the *Smart Factory* or *Industrial Automation*. This application area is typified by extremely demanding reliability and latency requirements for 5G communication links between sensors, actuators, and controllers.

Traditionally, Ethernet has been used to provide network connectivity. For smart factories, wireless networks provide many advantages over Ethernet:

- Reduced cost of manufacturing, installation, and maintenance
- Higher long-term reliability as wired connections suffer from wear and tear in motion applications
- Inherent deployment flexibility

With 5G, dispersed IoT sensors, actuators, controllers, and robots driven by software command and control can expand the ability to more fully automate an industrial process.

The application area that encompasses *Ground Vehicles, Drones, and Robots* refers to remotely controlled mobile devices and robots. Such devices are in common use in factory applications, but are also deployed in other contexts, such as smart agriculture. One area of particular interest is unmanned aircraft traffic management.

Tactile Interaction refers to a level of responsiveness that works at a human scale. For example, remote health care or gaming applications may require very low round-trip times to convince human senses that the perceived touch, sight, and sound are lifelike. These use cases involve interaction between humans and systems, where humans wirelessly control real and virtual objects, and the interaction requires a tactile control signal with audio or visual feedback. Robotic controls and interaction include several scenarios with many applications in manufacturing, remote medical care, and autonomous cars. The tactile interaction requires real-time reactions on the order of a few milliseconds. Remote surgery, discussed later in this article, is perhaps the most demanding use case. Table 3 gives typical values of *Key Performance Indicators* (KPIs) for tactile Internet applications.

Table 3: Key Performance Indicators for Tactile Internet

KPI	Value
Traffic Volume Density	0.03–1 Mbps/m ² / (cell radius 100 m ²)
Experienced User Throughput	0.3–1 Mbps (UL)
Latency	User-plane latency less than 2 ms
Availability	>99.999%
Reliability	>99.999 % for healthcare or remote driving/manipulation 95 % for remote gaming or remote augmented reality

AR and VR tend to have relatively high data-rate requirements. Some specific use cases also have URLLC requirements.

A paper from the *Next Generation Mobile Networks Alliance* (NGMN)^[16] lists three AR/VR examples with URLLC requirements:

- *Augmented Worker*: Augmented work is work that integrates digital technologies into the industrial environment to improve how work is done. Augmented work is appropriate for situations when it is not cost-effective or even possible to fully automate tasks, but it is desirable to augment the capabilities of the human worker. A good example is a task such as equipment repair where the access is difficult (for example, a hazardous environment) or the expert is at a remote place. The remote worker can be equipped with an AR headset and some sort of tactile interface for remote control. Sensor information from the remote target location in terms of audio, video, and haptic (tactile) enables the remote operator to control actuators at the target location to achieve the required work.
- *360 Panoramic VR View Video Broadcasting*: 360-degree videos are video recordings where a view in every direction is recorded at the same time, shot using an omnidirectional camera or a collection of cameras. With 360 panoramic VR view video broadcasting, the video is broadcast in real time. Remote users with VR headsets can view the live video feed, and by turning their head, see the point of view change in real time.
- *AR and MR Cloud Gaming*: A good example of an application in the AR/VR area that requires URLLC performance is AR and *Mixed Reality* (MR) cloud gaming, which is real-time game playing using a thin client with the bulk of the software on edge servers. This online gaming service provides on-demand streaming of games onto computers, consoles, and mobile devices. Thus, the user does not have to upgrade frequently and to deal with compatibility issues. Highly interactive games with tight QoS requirements generate the need for low-latency network performance.

Use cases in the category of *Emergency, Disasters, and Public Safety* generally require high reliability to enable response to natural disasters and emergencies. Accurate position location and very low latency to enable rapid response are also often critical requirements.

The *Urgent Health Care* category refers to applications involving remote diagnosis and treatment. A white paper from 5G Americas^[17] lists the following examples in this category:

- *Remote Patient Monitoring*: This use case involves remote patient monitoring via communication with devices that measure certain health indicators, such as pulse, blood glucose, blood pressure, and temperature. On an individual basis, the data rate and latency requirements are modest. However, for this use case to become pervasive, 5G is needed to support the massive increase in the number of connections per square meter while still maintaining the requisite QoS.

- *Remote Health Care:* This use case provides for individualized consultation, treatment, and patient monitoring built on a video linkup capability. The video conferencing can be augmented with remote transfer of health-related data in real time. Treatment could also be offered using smart pharmaceutical devices that correctly administer approved dosages of a drug on a schedule specified by the physician or practitioner.
- *Remote Surgery:* More demanding is remote surgery via control of robotic devices. This application area may be appropriate in ambulances, disaster sites, and remote areas. Important requirements are precise control and very low latency, very high reliability, and tight security.

The *Next Generation Mobile Networks Alliance*^[16] lists the following examples in the Intelligent Transportation category:

- *Assisted Driving:* 5G enables the delivery of advanced driver-assistance features that reduce fatal accidents and traffic congestion. These features include real-time maps for navigation, speed warnings, road hazards, vulnerabilities, heads-up display systems, and sensor data sharing. These features will enable the vehicle to dynamically change its course on the road under certain scenarios and conditions. *Vehicle-to-Network* (V2N) communication is necessary for this use case. Information from the vehicle enables the remote application to perform short-range modelling and recognition of surrounding objects and vehicles plus mid- to long-range modelling of the surroundings using information on the latest digital maps, traffic signs, traffic-signal locations, road construction, and traffic congestion.
- *Autonomous Driving:* Fully autonomous driving involves the capability of a vehicle to sense its environment and navigate without human input under all scenarios and conditions. A 5G network with URLLC capability enables numerous necessary features, including the use of complex algorithms to distinguish between different cars on the road and identify appropriate navigation paths given obstacles and considering the rules of the road, and the exchange of information in real time between thousands of cars connected in the same area.
- *Tele-operated Driving:* This use case refers to the use of remote driver assistance in areas where automatic driving is not possible. This assistance can provide enhanced safety for disabled people, elderly populations, and drivers in complex traffic situations. Typical application scenarios include disaster areas and unexpected and difficult terrains for manual driving such as in mining and construction. Tele-operated driving requires the wireless network to support V2N communication of video, sound feed information, and diagnostics from the vehicle, along with environmental information, to the remote driver. The network must support transmitting control commands from the remote driver to the vehicle to maneuver the vehicle in real time.

References and Further Reading

- [1] ITU-R, “IMT Vision—Framework and overall objectives of the future development of IMT for 2020 and beyond,” ITU-R Recommendation M.2083, September 2015.
- [2] ITU-R, “Minimum requirements related to technical performance for IMT-2020 radio interface(s),” ITU-R Report M.2410, November 2017.
- [3] ITU-T, “Requirements of the IMT-2020 network,” ITU-T Recommendation Y.3101, April 2018.
- [4] ITU-T, “High-level technical characteristics of network softwarization for IMT-2020,” ITU-T Recommendation Y.3150, January 2018.
- [5] ITU-T, “Framework of the IMT-2020 network,” ITU-T Recommendation Y.3102, May 2018.
- [6] 3GPP TS 22.261, “Technical Specification Group Services and System Aspects; Service requirements for the 5G system; Stage 1 (Release 18),” January 2021.
- [7] ITU-R, “Guidelines for evaluation of radio interface technologies for IMT-2020,” ITU-R Report M.2412, October 2017.
- [8] Mandar N. Kulkarni, Sarabjot Singh, and Jeffrey G. Andrew, “Coverage and Rate Trends in Dense Urban mmWave Cellular Networks,” 2014 *IEEE Global Communications Conference*, December 2014.
- [9] EU Project METIS, “Updated scenarios, requirements and KPIs for 5G mobile and wireless system with recommendations for future investigations,” Deliverable D1.5, April 2015.
- [10] Ericsson, “Cellular IoT Evolution for Industry Digitalization,” Ericsson White Paper, January 2019.
- [11] Ericsson, *Ericsson Mobility Report*, November 2020. (Issued annually).
- [12] Ericsson, “Cellular networks for Massive IoT,” Ericsson White Paper, January 2020.
- [13] ITU-T, “Overview of key performance indicators in smart sustainable cities,” ITU-T Recommendation Y.4900, June 2016.
- [14] ITU-T, “Supplement on use cases of smart cities and communities,” IT-T Series Y Supplement 56, December 2019.
- [15] 5G Americas, “New Services & Applications with 5G Ultra-Reliable Low-Latency Communications,” November 2018.
- [16] Next Generation Mobile Networks Alliance, “Verticals URLLC Use Cases and Requirements,” February 2020.
- [17] 5G Americas, “5G Services & Use Cases,” November 2017.

- [18] William Stallings, “Network Functions Virtualization,” *The Internet Protocol Journal*, Volume 24, No. 2, July 2021.
- [19] William Stallings, *5G Wireless: A Comprehensive Introduction*, ISBN-13: 9780136767299, Pearson, 2021.

WILLIAM STALLINGS is a consultant, lecturer, and author of over a dozen books on data communications and computer networking. He has a PhD in computer science from M.I.T. He has written numerous books on computer networking and computer architecture. His home in cyberspace is **WilliamStallings.com** and he can be reached at **ws@shore.net**

Our Privacy Policy

The *General Data Protection Regulation* (GDPR) is a regulation for data protection and privacy for all individual citizens of the *European Union* (EU) and the *European Economic Area* (EEA). Its implementation in May 2018 led many organizations worldwide to post or update privacy statements regarding how they handle information collected in the course of business. Such statements tend to be long and include carefully crafted legal language. We realize that we may need to provide similar language on our website and in the printed edition, but until such a statement has been developed here is an explanation of how we use any information you have supplied relating to your subscription:

- The mailing list for *The Internet Protocol Journal* (IPJ) is entirely “opt in.” We never have and never will use mailing lists from other organizations for any purpose.
- You may unsubscribe at any time using our online subscription system or by contacting us via e-mail. We will honor any request to remove your name and contact information from our database.
- We will use your contact information only to communicate with you about your subscription; for example, to inform you that a new issue is available, that your subscription needs to be renewed, or that your printed copy has been returned to us as undeliverable by the postal authorities.
- We will never use your contact information for any other purpose or provide the subscription list to any third party other than for the purpose of distributing IPJ by post or by electronic means.
- If you make a donation in support of the journal, your name will be listed on our website and in print unless you tell us otherwise.

Low Earth Orbit Satellite Systems for Internet Access

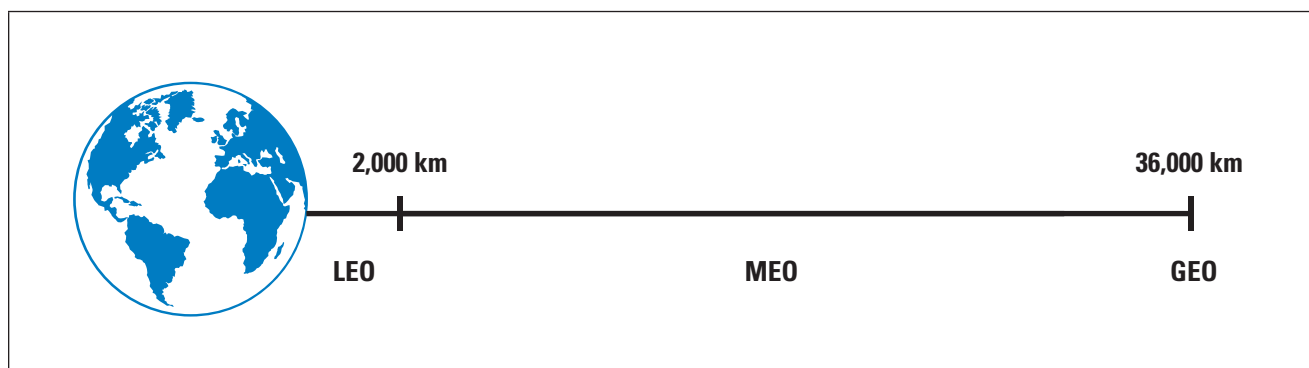
by Dan York, Internet Society, and Geoff Huston, APNIC

Satellites have been providing Internet connectivity for a few decades now, but new technologies using satellites in *Low Earth Orbit* (LEO) have created a race to offer space-based Internet access that provides ubiquitous, high-speed, and low-latency connectivity. This article explores what is happening, what is new, and why our planet may soon be circled with 50,000 to 90,000 new satellites!

Understanding Orbits

To understand the excitement over LEOs, you need to first understand the orbits of satellites (Figure 1).

Figure 1: Orbital Altitudes



Geostationary Orbit

If you fire a projectile at a speed greater than the Earth’s “escape velocity,” (11.2 km/s) it will head away from the Earth. If you reduce the speed slightly, the projectile will be caught by the Earth’s gravity and try to fall back to Earth. If you also incline the aiming trajectory, then instead of falling back to Earth, it will settle into an orbit around the Earth. The orbital speed relative to the Earth is a function of the altitude of the object. At very high altitudes, such as the moon, the orbital period is slower than the Earth’s rotation, while at very low altitudes the orbital period is down to a small number of hours, implying that there is a mid-spot where the orbital period is the same as the Earth’s rotation. If you launch a satellite to an altitude of 35,768 km above the equator, orbiting in the same direction as the Earth on the equatorial plane, then from the Earth the spacecraft appears to sit in a stationary position when observed from the Earth’s surface; this orbit is a *Geostationary Orbit* (GEO)^[0]. This geometry is what allows people to set up a satellite send/receive dish and point it at a specific location in the sky where the satellite is positioned—and never change the orientation of the dish.

GEO satellites are sufficiently distant that they can cover an entire hemisphere of the Earth's surface. However, they are normally equipped with a collection of transponders, most of which are focused on smaller areas, allowing the satellite to service multiple specific target regions at once, with greater total capacity as a result. A satellite operator can achieve global coverage with as few as three satellites. In addition to global service platforms, many nations have launched GEO satellites that are stationed over their country to provide communication services across their region.

These satellites are typically the size of a large bus and are expensive in terms of both construction costs and launch cost. Both the Moon and the Sun exert gravitational effects on the satellite, and, to a lesser extent, solar radiation pressure, all causing the satellite to drift away from its geosynchronous position. To counter this drift, the satellite is equipped with thrusters and some form of propellant. The total amount of onboard fuel defines an upper limit to the time that station position can be maintained, and these satellites typically have an operational lifespan of around 15–20 years.

In order to keep the level of radio interference between adjacent satellites to an acceptably low level, there are a limited number of geostationary orbit locations. Typically, geosynchronous satellite stations are separated by 2 degrees of angle as seen from the Earth, or 1,471 km apart in orbit. Disputes between nations over the deployment of satellites in this orbit are addressed through the coordination work hosted by the *International Telecommunications Union Radiocommunications Sector* (ITU-R).

A challenge with using GEO satellites for Internet access is that they are so far away from Earth. It takes a minimum of 238 ms for a signal to travel from the surface of the Earth to a satellite positioned 35,768 km away and back again. The *Round Trip Time* (RTT) to propagate an outbound packet via a GEO satellite and receive a reply is a minimum of 476 ms. The distance to the satellite increases as you move away from the position directly underneath the satellite on the Earth's equator, and the propagation time for the round-trip approaches 560 ms as you approach the limit of clear signal access near the polar areas. When you add delays for signal encoding, switching, and other terrestrial elements, the delivered performance of a service based on geosynchronous satellites is a typical RTT of around 660 ms, or two-thirds of a second. For many applications that are tuned to operate efficiently on faster terrestrial paths, this extended delay often causes the application to be sluggish and unresponsive in terms of its performance.

It is also the case that at this altitude the Earth's magnetic field provides far less shielding from solar radiation via the *Van Allen Belt*, so the electronics for GEO satellites need to have appropriate shielding, and the onboard electronics must tolerate a certain amount of radiation exposure.

Low Earth Orbit

An orbiting spacecraft needs to be positioned at least 160 km above the surface of the Earth, or it will encounter significant drag from the top of the Earth's atmosphere and its orbit will quickly decay, with an inevitable result. Above this altitude, it is viable to position orbiting spacecraft without needing to provide large quantities of continual propulsion (although some residual drag is experienced in orbital altitudes up to 500 km or so). For example, the International Space Station orbits at an altitude 400 km, with an orbital period of some 90 minutes. This region of space, where the orbits are higher than around 160 km and below 2,000 km, is termed the *Low Earth Orbit* (LEO) region. This is the region where we've positioned most of our satellites, as they are more accessible in terms of launch cost.

LEO satellites are close enough to the Earth's surface that signal propagation time to the satellite and back can be between 4 and 8 ms, which gives a range of RTT measurements for packet transmission via LEO services in the range 10–50 ms, a range comparable to that of terrestrial systems. With per-access service capacities of between 10 and 200 Mbps, LEO services can support most forms of modern real-time communication and online interaction^[1].

However, LEO-based satellite services require more complexity. At an altitude of 550 km, for example, a satellite will be visible from the Earth's surface in a circular area with a radius of some 900 km. Its orbital path is such that its velocity will be some 27,000 km/h, and each spacecraft will be visible from a fixed point on the Earth's surface for 5 minutes. In other words, to provide a continuous service over a fixed point, an evenly distributed collection of a minimum of 21 spacecraft would be needed in an orbital plane to ensure that as one satellite falls below the horizon, another is rising from the opposite horizon. Higher numbers of satellites in the orbital plane ensure a more reliable service and allow the Earth stations to avoid using spacecraft that are low in the horizon. To cover the entire Earth's surface, you need a minimum of 21 such orbital planes if you are using a 550-km altitude. The result is that, instead of just three satellites to provide a GEO service to anywhere on the planet, you would need hundreds or even thousands of satellites to provide the same comprehensive coverage. However, this scenario has some benefits in that the total capacity available from the system would also be many thousands of times greater in aggregate than the sparse GEO arrangement (refer to Figure 2).

LEO satellite systems may also be in multiple “shells” at different altitudes. For example, at the time of this article in July 2023, SpaceX's Starlink has launched 3,982 satellites into shells from 550 to 570 km, and another 750 satellites into shells from 525 to 535 km. Another provider, OneWeb, has launched 634 satellites into shells at 1,200 km in altitude^[2].

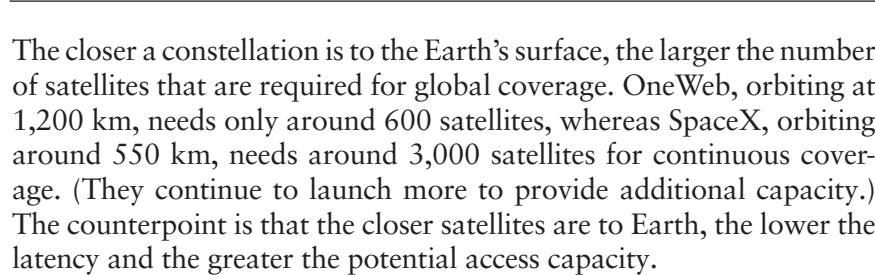
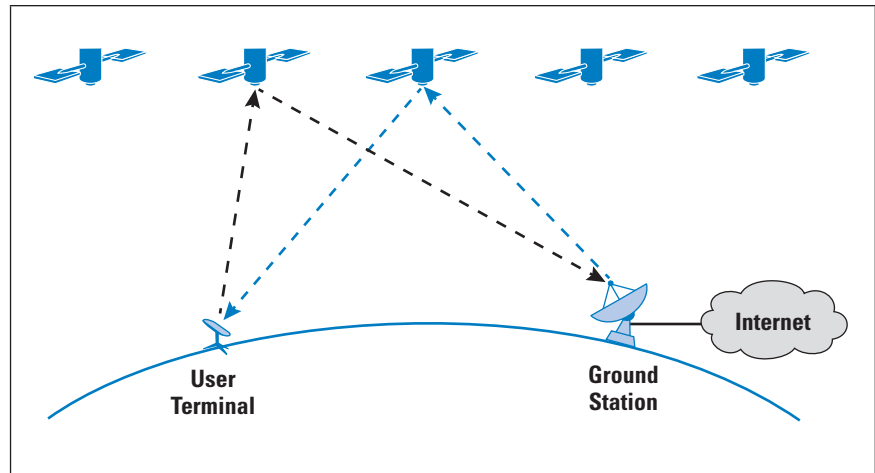


Figure 3: Satellite Handoff



Medium Earth Orbit

The region of space between 2,000 and 36,000 km above the Earth is referred to as *Medium Earth Orbit* (MEO). At the low end of this range, 2,000 km, the orbital period is around 2 hours, and the time to pass a signal up to the spacecraft and back is around 13 ms if you are located directly under the spacecraft. The higher the altitude the longer the orbital period, which, in turn reduces the number of spacecraft required to support a continuously available service. The transmission time to bounce a signal off the spacecraft at an altitude of 20,000 km is around 140 ms.

This zone encompasses the inner and outer Van Allen belts, which are belts of energetic charged particles that are trapped into an Earth orbit because of the Earth's magnetic field. The good news is that these belts protect the Earth's atmosphere from being blown away by solar radiation (as appears to have happened to Mars when its inner core solidified). The not-so-good news for satellites is that orbiting in this belt is like wandering through a firing range—there is always the possibility the sensitive electronics are damaged by a strike from one of these energetic particles. The outer belt is less dense, but the particles can have significantly higher energy levels. Beyond the Van Allen belts spacecraft encounter far higher levels of risk of damage from cosmic rays and solar radiation.

A region between the inner and outer Van Allen belts lies approximately between 12,000 and 24,000 km in altitude, which has a lower incidence of such energetic particles. The belts fluctuate in size and shape due to changes in the levels of solar radiation. The Earth itself acts as a shield, so that the belt is more compressed facing towards the Sun and extends further out on the “dark” side of the Earth.

The major satellites in this region are the satellite systems that support navigation, such as the *Global Positioning System* (GPS) using 31 spacecraft orbiting at some 20,200 km, *Galileo* with 24 active spacecraft at 23,222 km, *GLONASS* using 24 orbiting spacecraft at an altitude of 19,100 km, and *BeiDou* with 30 MEO satellites at 21,150 km.

For Internet access, the only major provider currently operating in MEO is the O3b network of around 20 satellites, operating at an altitude of approximately 8,000 km. O3b (which originally stood for “other 3 billion,” referring to the number of people still offline) began offering Internet connectivity in 2014 and was acquired by SES in 2016. SES is continuing to expand the service and is in the process of launching a new generation of 11 “O3b mPOWER” satellites into MEO, with promises to offer speeds up to multiple gigabits per second to its commercial customers.

Signal latency is higher than for LEOs, but significantly less than that of GEOs. Presumably feeling the competitive pressure from the LEO industry, SES has been working with many of the GEO providers to provide “multi-orbit” connectivity options that combine both MEO and GEO systems.

MEOs are a compromise in many ways. The Earth equipment still needs to perform tracking of the satellite, and that limits the power and sensitivity of the MEO antennae, yet the increased distance limits the performance and capacity of the system. The higher altitude provides greater coverage per satellite, allowing for broad coverage of the Earth’s surface with fewer spacecraft in the constellation, but the smaller number of satellites limits the overall capacity of the system.

If launch costs had remained high, then MEO systems made more sense in terms of minimizing the initial cost of the operation and maximizing the potential user base for the MEO satellite service, but the dramatic change in launch costs for LEO systems coupled with the use of phased array low-power steerable antennae has shifted the position quite dramatically in favor of LEO services. While GEO and MEO services tend to operate as wholesale services to a limited set of commercial customers using a conventional leased circuit service model, LEO systems have entered the consumer market, operating a direct access service as a retail service.

User Equipment

Consumers who want to connect to a LEO service for Internet access need to purchase what the satellite industry calls a *User Terminal*. This equipment includes the antenna and some access terminal, such as a small Ethernet switch or Wi-Fi access point. The phased-array antennas are compact, lightweight, and user-installable. Amazon has demonstrated some prototype user terminals that are small enough to fit in a backpack, and some mobile carriers have entered into agreements with Starlink to provide mobile handset access services directly from the handset to the satellite system (presumably with a significantly lower service capacity because of the limitations of the radio antenna on the handset).

There are more-complex user terminals that use parabolic dish antennae. These systems can achieve higher capacity and superior performance, but they require some form of mechanical steering to track the LEO satellite.

You can achieve the satellite-to-satellite handover function using dual antennae terminals, with one tracking the satellite for the active connection while the second pivots to focus on the next satellite in sequence.

Earth Stations and Inter-Satellite Connections

Satellite communications systems have conventionally operated as a “mirror in the sky.” The satellite receives a signal sent up from a user terminal and switches to a sending transponder that beams the signal back to an Earth station, which then passes the signal into the terrestrial network. With communications systems based on a GEO configuration it was possible to use just three Earth stations to service the entire system, given the hemispherical visibility of GEO satellites.

MEO and LEO satellites have more limited visibility, and in a “mirror” mode of operation, the density of Earth stations depends on the satellite altitude. For Starlink satellites at a 550-km altitude, Earth stations would need to be configured in a grid with around a 1,000-km spacing. This setup may be feasible in some locales, but if the intended service coverage includes more remote areas and coverage across oceans, then a different approach to Earth stations would be needed for LEO systems. An interesting technical development with LEO constellations to respond to this situation is the development of *Inter-Satellite Lasers* (ISLs) to send data between satellites within a constellation.

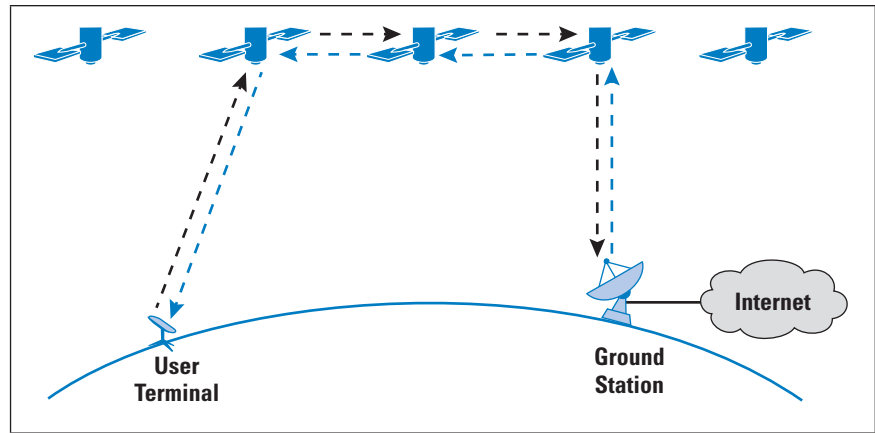
With the use of ISLs in the Starlink constellation, connections can now be passed from one satellite to another, and they do so in a relay form until they can drop down to an Earth station or user terminal. For example, SpaceX used this solution in early 2023 to connect remote Antarctic research stations, without having any ground stations on Antarctica^[4]. Details of Starlink’s ISL system, including capacity and configuration, are not yet published, but there have been recent service announcements that show service availability in regions where no Earth stations exist.

A simple approach to ISLs would be to link adjacent satellites in a common orbital plane (Figure 4). A more flexible approach might be to link adjacent satellites across different orbital paths, but the relative closing speeds of satellites on different orbital planes may well be beyond the capabilities of steerable laser systems. For example, a spacecraft travelling north/south at 27,000 km/h attempting to track a satellite passing on an east-west plane at the same speed would have a closing speed of some 700 km/h and an angular velocity of one angle of degree every three seconds, which may present some challenges to an onboard laser steering system.

Spectrum and Regulatory Approvals

Satellite communication requires allocation of certain frequencies in the radio spectrum for communication from the consumer equipment (such as the antenna) to the satellites, and from the satellites to a ground station where connections are made to the Internet.

Figure 4: Inter-Satellite Lasers



Spectrum management could easily consume an entire article by itself, but at a high level, spectrum “allocations” are coordinated through the *International Telecommunication Union (ITU)*, specifically the *ITU Radiocommunication Sector (ITU-R)*. Most of the LEO systems use frequencies within the Ku (10.7–14.5 GHz) and Ka (17.3–30 GHz) bands, which are the bands that are intended for use by Broadband Satellite Services. There is also the more recently allocated Q/V band (37.5–51.4 GHz), which is available for use in this context and has been deployed already in some systems.

Initial Allocations

The world’s nations have agreed that above some altitude “outer space” began. The implication of this agreement is that the sovereign rights that apply to the defined surface parts of the Earth extend only up to the point of “outer space.” Oddly enough, the world’s nations did not agree as to where “outer space” begins, and some nations claim sovereignty up to an altitude of only 100 km, while others extend that further to 160 km or more. In any case, the result is that there is no national regime that must approve or otherwise say what a spacecraft may do in outer space in the form of “over flying” its territory. However, some conventions apply to assist various folk to coordinate their actions in space and assist in resolving any disputes that may arise. The use of GEO station slots by various nations is one such area where conventions apply, and the ITU-R assists in this coordination activity.

However, when the topic shifts to that of communications between Earth and satellites, there is a requirement to get various forms of national regulatory approval, based on the location of the Earth stations. The failure to gain such approvals for the Iridium service was the major cause of the early business failure of this venture in 2000. The approval is not quite as simple as just approval for the operation of Earth stations. When a company wants to launch a LEO satellite constellation, it conventionally obtains approval from its national regulator. For example, SpaceX and Amazon are both US companies, so they filed their requests with the US *Federal Communication Commission (FCC)*. The filings include the radio frequencies they want to use and the number, altitudes, and orbital planes at which their satellites will operate. These filings are also forwarded to the ITU.

The LEO and MEO space allocations are generally operated on a first-come, first-served basis, but to prevent people from “squatting” on spectrum and altitude allocations, the ITU requires LEO satellite constellations to have 10 percent of their constellation in orbit within the first two years after the start of deployment, 50 percent in five years, and 100 percent in seven years. This factor is part the reason why there is a great amount of heightened activity to launch LEO constellations. The various companies who have lodged applications need to meet these deployment milestones or they risk losing the exclusivity of their spectrum allocations.

For GEO satellites the ITU-R coordinates the spectrum allocations and the orbital slots. For LEO or MEO satellites, the ITU-R coordinates only the spectrum allocations, and does not coordinate the orbital planes and altitudes for these satellites. That aspect is handled entirely by national regulators.

When a company obtains the necessary spectrum and altitude allocations from its national regulator, it then can launch its satellite platforms into orbit. This launch activity requires completely separate approvals and involves processes whose descriptions fall outside the scope of this article.

National Approvals

As part of obtaining the initial spectrum allocations, a LEO system provider receives the approval to operate within its home nation. Then the provider has to go to the spectrum regulators in every country in which it wishes to operate to service and receive regulatory approval to use the spectrum in that country.

In some cases, the requested spectrum may be already in use, and the country faces the difficult decision of whether to re-configure its local use situation or be denied use of the LEO system. An example is Armenia, where a national regulator representative informed the audience of the *Armenian Internet Governance Forum 2022* that using Starlink was not possible anytime soon because the frequencies were in use by the Armenian military and government. Given that SpaceX will not change its frequencies, there probably will not be an option until the Armenian government changes its own systems to use different frequencies, which could take time and some amount of unplanned costs.

While it is technically possible for a LEO provider to offer service in a country for which it does not have permission to operate (SpaceX activated its service during the protests in Iran in late 2022^[5]), it is not legally permitted to simply bring an antenna into a country and start using it with a LEO system. The LEO providers must obtain permission to operate the service in each country before they can make their service available to customers. Additionally, the LEO providers may also need to get approval to distribute the consumer equipment, and approval to interconnect with local terrestrial infrastructure.

Standards

We have very little visibility into how the internal networks operate, but at the Internet Protocol and application layers, the LEO constellations so far seem to support all the conventional standards for Internet Protocol forwarding and Internet operations created by the *Internet Engineering Task Force* (IETF).

It seems at the moment that the major LEO operators (Starlink, OneWeb, Amazon Project Kuiper) are all pursuing their own proprietary systems for communication between their user terminals (antennas) and their satellites, and between their satellites and their ground stations. This process will require consumers to purchase completely separate user terminals in order to use each different system. Perhaps at some point this process will become standardized, but not in this initial period of deployment.

Still Many Questions

That point is perhaps a critical one. Regardless of any marketing hype, the reality is that the LEO Internet access industry is very much still in its infancy. Only SpaceX's Starlink has global coverage. OneWeb has launched sufficient satellites to attain global coverage and is in the process of getting all its satellites in position. It hopes to offer global connectivity by the end of 2023, concentrating its service on the government and enterprise sector. Amazon's Project Kuiper has been manufacturing its satellites and equipment, but is still waiting for rocket availability to get its satellites into space. Many other companies are in various stages of getting their systems underway.

There are still many open questions, many of which the Internet Society explored in a recent document about LEO satellites^[6]. What will the capacity of these LEO systems truly be? Will they be able to support all the many devices we want to connect to them? Will the systems be affordable by those who need the connectivity the most? Using space-based platforms to provide global coverage to the billions of unconnected people probably would require some significant changes in the service model, because the challenges, particularly in terms of affordability, are still significant.^[8]

Will these constellations all be able to operate without interfering with each other? Will consumers tolerate the costs of proprietary equipment and the high cost of switching? What about the problems of "space debris" resulting from collisions or inactive satellites? Do we understand the potential environmental impact of having so many satellites burning up in our upper atmosphere when they reach their five-year end-of-life? Or the impact of all the regular rocket launches needed to resupply the constellations with new satellites? Questions abound, and many of them we may not be able to answer until we have the experience with getting more LEO constellations online.

Looking Ahead

Without a doubt, the next few years will be extremely active:

- SpaceX plans to complete its “Gen1” constellation of 4,408 satellites. In addition, SpaceX has received US FCC approval for 7,500 satellites in its “Gen2” constellation, which the company hopes to grow to almost 30,000 satellites. The company has also announced numerous “direct-to-phone” services with a collection of mobile network operators.
- OneWeb aims to have its global connectivity service available by the end of 2023.
- Over the next two years, Amazon is seeking to launch its Project Kuiper, a direct competitor to Starlink in the consumer market, and have it operational in 2025.
- The Chinese government is seeking to launch its own LEO constellation, called “Guowang,” which will have almost 13,000 satellites.

Filings with the ITU show that there is a path where as many as 90,000 satellites could be launching into LEO over the next several years. Here are a few examples:

SpaceX Starlink Gen 1	4,408
SpaceX Starlink Gen 2	29,988
OneWeb, Phase 1	718
OneWeb, Phase 2	6,372
Amazon Project Kuiper	7,774
China Guowang	12,992
Astra	13,620
Boeing	5,842
Globalstar	3,080
Lynk	2,000
Telesat Lightspeed	1,969
Spin Launch	1,190
TOTAL	89,953

To put this information in perspective, LEO, MEO, and GEO orbits currently have only about 8,000 active satellites—and SpaceX operates over 4,500 of them!

Even more satellite systems are in planning stages; this article has covered only LEO satellites used for Internet access. In addition, LEO constellations are being launched for sensor networks (as in the Internet of Things), imaging/photography, and much more.

How many of these constellations will actually be launched into orbit and be used as a service is an open question. Navigating all the required regulatory approvals across all the various national regulatory regimes is very challenging. Unless you are SpaceX and own your own rockets, launching satellites into space is extremely difficult right now, as many current rocket programs are delayed or simply unavailable.

We are also seeing that the rise of the competition of space-based Internet systems is causing ground-based Internet service providers to accelerate their plans for deploying terrestrial networks. As good as LEO systems may be, fiber networks can still provide even higher speeds. There is still much in the way of potential to use fiber-based trunk networks with last-mile access provided by mobile radio technologies, and these networks are well-understood technologies with an already-dominant user base. If the terrestrial access service providers are successful in making even faster connectivity more widely available and affordable, then the competition for the user between space and terrestrial-based systems would increase in intensity, and we can anticipate that such competition will result in lower costs, wider coverage, and improved performance for consumers.

In any case, there are tremendous opportunities for LEO satellite systems to help us connect the unconnected, create more resilient networks, help coordinate disaster-relief efforts, and generally bring high-speed, low-latency Internet connectivity to everyone, everywhere. The next few years will show us whether we can make these goals a reality.

References and Further Reading

- [0] Note that in most satellite-related policies from the ITU and other regulators, geostationary orbit is abbreviated as “GSO,” and both LEO and MEO satellites are grouped as “Non-Geosynchronous Orbit” or “NGSO” satellites.
- [1] Josh Fomon, “New Speedtest Data Shows Starlink Users Love Their Provider,” *Ookla Insights Articles*, May 8, 2023.
- [2] Jonathan’s Space Pages:
<https://planet4589.org/space/con/conlist.html>
- [3] *Our World in Data*, “Cost of space launches to low Earth orbit.”
- [4] Kevin Hurler, “Starlink Is Now Connecting Remote Antarctic Research Camps to the Internet,” *Gizmodo*, January 23, 2023.
- [5] Karl Vick, “Inside the Clandestine Efforts to Smuggle Starlink Internet Into Iran,” *TIME*, January 25, 2023.
- [6] “Perspectives on LEO Satellites: Using Low Earth Orbit Satellites for Internet Access,” Internet Society, 2022.
- [7] “sat-int” IETF mailing list for discussion of IETF technologies for satellite networking:
<https://www.ietf.org/mailman/listinfo/Sat-int>

- [8] “Satellite broadband for the masses: Are we there yet?” Presentation by Mike Puchol to the APRICOT 2023/APNIC 55 Conference, March 2023.
- [9] Ulrich Speidel, “Getting hands-on experience with Starlink,” *APNIC Blog*, March 16, 2023.
- [10] Ulrich Speidel, “Everything you wanted to know about LEO satellites, Part 1: The basics,” *APNIC Blog*, May 20, 2021.
- [11] Ulrich Speidel, “Everything you wanted to know about LEO satellites, Part 2: Constellations, Gateways and Antennas,” *APNIC Blog*, May 27, 2021.
- [12] Ulrich Speidel, “Everything you wanted to know about LEO satellites, Part 3: Bandwidth, System Capacity and Inter-satellite Routing,” *APNIC Blog*, June 2, 2021.
- [13] Ulrich Speidel, “Everything you wanted to know about LEO satellites, Part 4: Why direct to site?,” *APNIC Blog*, June 11, 2021.
- [14] Ulrich Speidel, “Satellite still a necessity for many Pacific Islands,” *APNIC Blog*, September 18, 2018.

DAN YORK serves the *Internet Society* (ISOC) as the Director, Internet Technology, with his focus on helping explain the changes happening to the Internet and how we use it. He led ISOC’s 2022 *Low Earth Orbit* (LEO) satellite project and continues to monitor LEO systems as part of ISOC’s work to connect the unconnected. Dan is also active within the DNS and real-time communications areas of the *Internet Engineering Task Force* (IETF). Since the mid-1980s Dan has been working with online communication technologies and helping businesses and organizations understand how to use and participate in these new media. An author of multiple books on networking, security, IPv6, and Linux, he frequently presents at industry conferences and events. He has been blogging, writing online, and podcasting since the early 2000s. He lives in Vermont, USA. E-mail: york@isoc.org

GEOFF HUSTON AM, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

Check your Subscription Details!

Make sure that both your postal and e-mail addresses are up-to-date since these are the only methods by which we can contact you. If you see the words “Invalid E-mail” on your printed copy this means that we have been unable to contact you through the e-mail address on file. If this is the case, please contact us at ipj@protocoljournal.org with your new information. The subscription portal is located here: <https://www.ipjsubscription.org/>

Letter to the Editor

On 25 years of The Internet Protocol Journal

Geoff Huston's magnum opus covering the last 25 years of the history of the Internet inspires me to reflect on those 25 years from several perspectives^[1]. The year 1998 saw the launch of two extraordinary companies resting on and extending the Internet infrastructure: Google and Akamai. Each, in its own way, survived the "dot-bust" and went on to make significant contributions to the utility of the Internet and the World Wide Web. Many other companies have come, some have gone, and some have continued to enjoy robust growth. For many of us, 2023 feels like a new inflection point in communications. Low earth-orbiting satellites promise to provide access to the Internet from every square inch of the planet. Inter-continental fiber networks continue to expand in number and capacity, and there is energizing interest in developing mesh-like infrastructure to respond to cable cuts through optical re-routing in switching units on the ocean floor.

Ironically, 1998 was also the year in which the *Interplanetary Special Interest Group* of the Internet Society was formed, and it still continues as the Interplanetary chapter (ipnsig.org). A new suite of interplanetary "Bundle Protocols" has been standardized by the IETF and by the *Consultative Committee on Space Data Systems* (CCSDS), whose standards are also part of the ISO standards library. As the *Artemis*, *LunaNet* (NASA), and *MoonLight* (ESA) missions unfold, bringing us back to the Moon, the prospect of commercial operations looms large and immediate. The governance challenges of competitive and cooperative public/private engagements will be topics of urgent discussion well before this decade is out. There will be lessons from the multi-stakeholder policy-making practices derived from the building of the Internet and neo-institutional energy as the need for new governance mechanisms emerge.

Wireless technologies continue to evolve, and edge computing and more disciplined forms of Wi-Fi are emerging. Social media and recently emerging *Large Language Model* neural networks are confounding technologists and policy makers as they seek to make the Internet a safer place while preserving its historical openness to new ideas, new applications, and new ways to share and discover knowledge. These new neural transformer networks are living up to their name by transforming our awareness of the importance of reliable information sources in a sea of misinformation and disinformation. Preservation of human rights has become an ever more urgent priority in the face of scaled abuse of online resources. We can hope that the solutions to the problems of artificial intelligence will be found in the technology itself.

[1] Geoff Huston, "Twenty-Five Years Later," *The Internet Protocol Journal*, Volume 25, No. 1, June 2023.

And for all those years, *The Internet Protocol Journal* has consistently shed light on the conundrums that confound Internet engineers, scientists, operators, and policy makers. Hats off to its long-time editor, Ole Jacobsen, for persistent and quality reporting and sharing of timely and technically sound information.

—Vint Cerf, Woodhurst, May 2023

Internet Society Launches NetLoss Calculator

The Internet Society recently launched the *NetLoss* calculator, a revolutionary tool that measures the economic impact of Internet shutdowns around the world. Hosted on the Internet Society's *Pulse* Platform that tracks and analyzes shutdowns, NetLoss uses a groundbreaking econometric framework to understand the impacts of shutdowns and provides an unprecedented level of rigor and precision in estimating their economic damage.

Internet shutdowns globally reached a record high in 2022, with governments around the world ordering Internet access and services to be restricted or blocked during civil unrest, school exams, and during elections, which resulted in major economic consequences.

According to the NetLoss calculator:

- The shutdown in Sudan in April 2023 is estimated to have cost the country more than \$3 million USD, as well as the loss of 560 jobs.
- The shutdown in Pakistan in May 2023 is estimated to have cost more than \$13 million USD, as well as increased unemployment.
- The shutdown in Guinea in June 2023 is estimated to have cost the country nearly half a million USD and job losses.

Governments often mistakenly believe that Internet shutdowns will quell unrest, stop the spread of misinformation, or reduce harm from cybersecurity threats. But shutdowns are extremely disruptive to economic activity: they halt e-commerce, generate losses in time-sensitive transactions, increase unemployment, interrupt business-customer communications, and create financial and reputational risks for companies. They also hurt a country's growth as research shows Internet adoption positively impacts *Gross Domestic Product* (GDP).

The Internet Society has long opposed the practice of Internet shutdowns, and urges all governments to refrain from implementing them due to the damage they inflict on a nation's economy, civil society, and Internet infrastructure. With NetLoss, organizations and advocates can demonstrate to governments and regulators how a shutdown will negatively impact their nation's economy.

In addition to the estimated cost of an Internet shutdown (that is, the loss in GDP), the Internet Society NetLoss calculator also estimates:

- The change in the unemployment rate due to a shutdown.
- The amount of *Foreign Direct Investment* (FDI) lost due to a shutdown.
- Risk of a shutdown: the probability that a country will experience a shutdown.

“The global rise in Internet shutdowns shows that governments continue to ignore the negative consequences of undermining the open, accessible, and secure nature of the global Internet. The calculator is a major step forward for the community of journalists, policymakers, technologists, and other stakeholders who are pushing back against the damaging practice of Internet shutdowns. Its groundbreaking and fully transparent methodology will help show governments around the world that shutting down the Internet is never a solution,” said Andrew Sullivan, President and CEO, of The Internet Society.

The calculator considers a wide range of economic impacts beyond traditional measures of economic output, such as GDP, to demonstrate the financial impact of an Internet shutdown. It also includes the change in the unemployment rate, the amount of FDI lost, and the risk of a future shutdowns.

In addition to its primary indicators, the NetLoss calculator’s methodology also takes into account other factors that can impact country-specific economic outcomes, including the age dependency ratio (percentage of working 18–65 years old to total population), the fraction of the population residing in urban areas, and the percentage of the labor force with basic education.

By using the following open data sets, the NetLoss calculator’s methodology is reproducible and transparent:

- *Shutdown Data*: Includes detailed event-level data on government-mandated shutdown events.
- *Protests and Civil Unrest*: Includes detailed event-level data on various events, their start and end dates, involved parties, and associated fatalities.
- *Elections*: The Constituency-Level Elections Archive maintained by Yale University provides elections data from 150 countries since 1960.
- *Socioeconomic Indicators*: The World Bank provides data on economic indicators including GDP per capita, employment, inflation, and foreign investment.

The framework used in the NetLoss calculator builds on the Internet Society’s longstanding research and advocacy on this issue via the Pulse Platform. Launched in December 2020, Internet Society Pulse consolidates trusted third-party Internet measurement data from various sources into a single platform to examine Internet trends and tell data-driven stories so that policymakers, researchers, journalists, network operators, civil society groups and others can better understand the health, availability, and evolution of the Internet. The source of data for NetLoss is the World Bank’s *World Development Indicators*, which typically corrects for minor statistical changes. Data used in the calculator is updated quarterly.

The NetLoss calculator can be found on the Pulse platform:
<https://pulse.internetsociety.org/netloss>

Reflections on Ten Years Past the Snowden Revelations

Authored by Stephen Farrell, Farzaneh Badii, Bruce Schneier, and Steven M. Bellovin, RFC 9446 contains the thoughts and recountings of events that transpired during and after the release of information about the United States *National Security Agency* (NSA) by Edward Snowden in 2013. There are four perspectives: that of someone who was involved with sifting through the information to responsibly inform the public, that of a security area director of the IETF, that of a human rights expert, and that of a computer science and affiliate law professor. The purpose of this memo is to provide some historical perspective, while at the same time offering a view as to what security and privacy challenges the technical community should consider. These essays do not represent a consensus view, but that of the individual authors. The RFC can be found here: <https://www.rfc-editor.org/info/rfc9446>

APNIC Celebrates 30 Years

This year, the *Asia Pacific Network Information Centre* (APNIC) enters its third decade. Starting from a tiny office in Tokyo in 1992—with three people and a spreadsheet serving less than 600 delegated entities—APNIC has grown to a community of nearly 24,000 organizations across 56 economies. The APNIC of today serves 2.6 billion Internet users, more than half the global Internet. APNIC economies also comprise more than half the global IPv6 capability.

Despite many changes in technology and policy worldwide, APNIC has remained committed to: “A global, open, stable, and secure Internet.” The Asia Pacific is home to nine of the world’s 46 Least Developed Countries (as defined by the United Nations). Fifteen states and seven affiliated economies of APNIC Members are *Small Island Developing States* (SIDS), characterized by low population, distance and remoteness. However, the region also includes global economic superpowers and some of the most populated economies on Earth.

The combination of language, culture, distance, isolation and the different scale of the region’s communities magnifies the importance of consensus policy making on an equal basis. The APNIC community has developed policy that reflects and enables Internet growth across our region, and has ensured an Asia Pacific voice has been heard at a global level.

To mark this 30-year milestone, the *APNIC Blog* will run a series looking back at the past and into the future. The intention is to share stories, anecdotes, milestones and insights that capture some of the essence of the last 30 years. For more information visit:

<https://blog.apnic.net/2023/08/08/apnic-celebrates-30-years/>

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Kjetil Aas	Ilia Bromberg	David Dillow	Serge Van Ginderachter	Richard Johnson
Fabrizio Accatino	Lukasz Bromirski	Richard Dodsworth	Greg Goddard	Jim Johnston
Michael Achola	Václav Brožík	Ernesto Doelling	Tiago Goncalves	Jonatan Jonasson
Martin Adkins	Christophe Brun	Michael Dolan	Ron Goodheart	Daniel Jones
Melchior Aelmans	Gareth Bryan	Eugene Doroniuk	Octavio Alfageme	Gary Jones
Christopher Affleck	Ron Buchalski	Michael Dragone	Gorostiaga	Jerry Jones
Scott Aitken	Paul Buchanan	Joshua Dreier	Barry Greene	Michael Jones
Jacobus Akkerhuis	Stefan Buckmann	Lutz Drink	Jeffrey Greene	Amar Joshi
Antonio Cuñat Alario	Caner Budakoglu	Aaron Dudek	Richard Gregor	Javier Juan
William Allaire	Darrell Budic	Dmitriy Dudko	Martijn Groenleer	David Jump
Nicola Altan	BugWorks	Andrew Dul	Geert Jan de Groot	Anders Marius Jørgensen
Shane Amante	Scott Burleigh	Joan Marc Riera	Ólafur Guðmundsson	Merike Kao
Marcelo do Amaral	Chad Burnham	Duocastella	Christopher Gumez	Andrew Kaiser
Matteo D'Ambrosio	Randy Bush	Pedro Duque	Gulf Coast Shots	Naoki Kambe
Selva Anandavel	Colin Butcher	Holger Durer	Sheryll de Guzman	Christos Karayiannis
Jens Andersson	Jon Harald Bøvre	Karlheinz Dölger	Rex Hale	Daniel Karrenberg
Danish Ansari	Olivier Cahagne	Mark Eanes	Jason Hall	David Kekar
Finn Arildsen	Antoine Camerlo	Andrew Edwards	James Hamilton	Stuart Kendrick
Tim Armstrong	Tracy Camp	Peter Robert Egli	Darow Han	Robert Kent
Richard Artes	Brian Candler	George Ehlers	Handy Networks LLC	Jithin Kesavan
Michael Aschwanden	Fabio Caneparo	Peter Eisses	Stephen Hanna	Jubal Kessler
David Atkins	Roberto Canonico	Torbjörn Eklöv	Martin Hannigan	Shan Ali Khan
Jac Backus	David Cardwell	Y Ertur	John Hardin	Nabeel Khatri
Jaime Badua	Richard Carrara	ERNW GmbH	David Harper	Dae Young Kim
Bent Bagger	John Cavanaugh	ESdatCo	Edward Hauser	William W. H. Kimandu
Eric Baker	Lj Cemerar	Steve Esquivel	David Hauweele	John King
Fred Baker	Dave Chapman	Jay Etchings	Marilyn Hay	Russell Kirk
Santosh Balagopalan	Stefanos Charchalakakis	Mikhail Evstiounin	Headcrafts SRLS	Gary Klesk
William Baltas	Molly Cheam	Bill Fenner	Hidde van der Heide	Anthony Klopp
David Bandinelli	Greg Chisholm	Paul Ferguson	Johan Helsingius	Henry Kluge
A C Barber	David Chosrova	Ricardo Ferreira	Robert Hinden	Michael Kluk
Benjamin Barkin-Wilkins	Marcin Cieslak	Kent Fichtner	Damien Holloway	Andrew Koch
Feras Batainah	Lauris Cikovskis	Armin Fisslthaler	Alain Van Hoof	Ia Kochiashvili
Michael Bazarewsky	Brad Clark	Michael Fiumano	Edward Hotard	Carsten Koempe
David Belson	Narelle Clark	The Flirble Organisation	Bill Huber	Richard Koene
Richard Bennett	Horst Clausen	Jean-Pierre Forcioli	Hagen Hultzs	Alexader Kogan
Matthew Best	James Cliver	Gary Ford	Kauto Huopio	Matthijs Koot
Hidde Beumer	Guido Coenders	Susan Forney	Asbjørn Højmark	Antonin Kral
Pier Paolo Biagi	Joseph Connolly	Christopher Forsyth	Kevin Iddles	Robert Krejčí
Arturo Bianchi	Steve Corbató	Andrew Fox	Mika Ilvesmaki	John Kristoff
John Bigrow	Brian Courtney	Craig Fox	Karsten Iwen	Terje Krogdahl
Orvar Ari Bjarnason	Beth and Steve Crocker	Fausto Franceschini	Joseph Jackson	Bobby Krupczak
Tyson Blanchard	Dave Crocker	Valerie Fronczak	David Jaffe	Murray Kucherawy
Axel Boeger	Kevin Croes	Tomislav Futivic	Ashford Jaggernaut	Warren Kumari
Keith Bogart	John Curran	Laurence Gagliani	Thomas Jalkanen	George Kuo
Mirko Bonadei	André Danthine	Edward Gallagher	Jozef Janitor	Dirk Kurfuerst
Roberto Bonalumi	Morgan Davis	Andrew Gallo	Martijn Jansen	Mathias Körber
Lolke Boonstra	Jeff Day	Chris Gamboni	John Jarvis	Darrell Lack
Julie Bottorff	Fernando Saldana Del Castillo	Xosé Bravo Garcia	Dennis Jennings	Andrew Lamb
Photography		Oswaldo Gazzaniga	Edward Jennings	Richard Lamb
Gerry Boudreaux	Rodolfo Delgado-Bueno	Kevin Gee	Aart Jochem	Yan Landriault
Leen de Braal	Julien Dhallenne	Rodney Gehrke	Nils Johansson	Edwin Lang
Kevin Breit	Freek Dijkstra	Greg Giessow	Brian Johnson	Sig Lange
Thomas Bridge	Geert Van Dijk	John Gilbert	Curtis Johnson	Markus Langenmair

Fred Langham	David Millsom	Derrell Piper	Philip Schneck	Douglas Thompson
Tracy LaQuey Parker	Desiree Miloshevic	Rob Pirnie	James Schneider	Kerry Thompson
Alex Latzko	Joost van der Minnen	Jorge Ivan Pincay	Peter Schoo	Lorin J Thompson
Jose Antonio Lazaro	Thomas Mino	Ponce	Dan Schrenk	Fabrizio Tivano
Lazaro	Rob Minshall	Marc Vives Piza	Richard Schultz	Peter Tomsu Fine Art
Antonio Leding	Wijnand Modderman-	Victoria Poncini	Timothy Schwab	Photography
Rick van Leeuwen	Lenstra	Blahoslav Popela	Roger Schwartz	Joseph Toste
Simon Leinen	Mohammad Moghaddas	Andrew Potter	SeenThere	Rey Tucker
Robert Lewis	Charles Monson	Ian Potts	Scott Seifel	Sandro Tumini
Christian Liberale	Andrea Montefusco	Eduard Llull Pou	Paul Selkirk	Angelo Turetta
Martin Lillepuu	Fernando Montenegro	Tim Pozar	Andre Serralheiro	Michael Turzanski
Roger Lindholm	Roberto Montoya	David Raistrick	Yury Shefer	Phil Tweedie
Link Light Networks	Joel Moore	Priyan R Rajeevan	Yaron Sheffer	Steve Ulrich
Mike Lochocki	John More	Balaji Rajendran	Doron Shikmoni	Unitek Engineering AG
Chris and Janet Lonvick	Maurizio Moroni	Paul Rathbone	Tj Shumway	John Urbanek
Sergio Loreti	Brian Mort	William Rawlings	Jeffrey Sicuranza	Martin Urwaleck
Eric Louie	Soenke Mumm	Mujtiba Raza Rizvi	Thorsten Sideboard	Betsy Vanderpool
Adam Loveless	Tariq Mustafa	Bill Reid	Greipur Sigurdsson	Surendran Vangadasalam
Josh Lowe	Stuart Nadin	Petr Rejhon	Fillipe Cajaiba da Silva	Ramnath Vasudha
Guillermo a Loyola	Michel Nakhla	Robert Remenyi	Andrew Simmons	Randy Veasley
Hannes Lubich	Mazdak Rajabi Nasab	Rodrigo Ribeiro	Pradeep Singh	Philip Venables
Dan Lynch	Krishna Natarajan	Glenn Ricart	Henry Sinnreich	Buddy Venne
David MacDuffie	Naveen Nathan	Justin Richards	Geoff Sisson	Alejandro Vennera
Sanya Madan	Darryl Newman	Rafael Riera	John Sisson	Luca Ventura
Miroslav Madić	Thomas Nikolajsen	Mark Risinger	Helge Skrivervik	Scott Vermillion
Alexis Madriz	Paul Nikolich	Fernando Robayo	Terry Slattery	Tom Vest
Carl Malamud	Travis Northrup	Michael Roberts	Darren Sleeth	Peter Villemoes
Jonathan Maldonado	Marijana Novakovic	Gregory Robinson	Richard Smit	Vista Global Coaching
Michael Malik	David Oates	Ron Rockrohr	Bob Smith	& Consulting
Tarmo Mammers	Ovidiu Obersterescu	Carlos Rodrigues	Courtney Smith	Dario Vitali
Yogesh Mangar	Jim Oplotnik	Magnus Romedahl	Eric Smith	Rüdiger Volk
John Mann	Tim O'Brien	Lex Van Roon	Mark Smith	Jeffrey Wagner
Bill Manning	Mike O'Connor	Marshall Rose	Tim Sneddon	Don Wahl
Harold March	Mike O'Dell	Alessandra Rosi	Craig Snell	Michael L Wahrman
Vincent Marchand	John O'Neill	David Ross	Job Snijders	Lakhinder Walia
Normando Marcolongo	Carl Örne	William Ross	Ronald Solano	Laurence Walker
Gabriel Marroquin	Packet Consulting	Boudhayan	Asit Som	Randy Watts
David Martin	Limited	Roychowdhury	Ignacio Soto Campos	Andrew Webster
Jim Martin	Carlos Astor Araujo	Carlos Rubio	Evandro Sousa	Jd Wegner
Ruben Tripiana Martin	Palmeira	Rainer Rudigier	Peter Spekrijse	Tim Weil
Timothy Martin	Gordon Palmer	Timo Ruiters	Thayumanavan Sridhar	Westmoreland
Carles Mateu	Alexis Panagopoulos	RustedMusic	Paul Stancik	Engineering Inc.
Juan Jose Marin Martinez	Gaurav Panwar	Babak Saberi	Ralf Stempfer	Rick Wesson
Ioan Maxim	Chris Parker	George Sadowsky	Matthew Stenberg	Peter Whimp
David Mazel	Alex Parkinson	Scott Sandefur	Martin Štěpánek	Russ White
Miles McCredie	Craig Partridge	Sachin Sapkal	Adrian Stevens	Jurrien Wijlhuizen
Brian McCullough	Manuel Uruena Pascual	Arturas Satkovskis	Clinton Stevens	Derick Winkworth
Joe McEachern	Ricardo Patara	PS Saunders	John Streck	Pindar Wong
Alexander McKenzie	Dipesh Patel	Richard Savoy	Martin Streule	Makarand Yerawadekar
Jay McMaster	Dan Paynter	John Sayer	David Strom	Phillip Yialeloglou
Mark Mc Nicholas	Leif Eric Pedersen	Phil Scarr	Colin Strutt	Janko Zavernik
Olaf Mehlberg	Rui Sao Pedro	Gianpaolo Scassellati	Viktor Sudakov	Bernd Zeimet
Carsten Melberg	Juan Pena	Elizabeth Scheid	Edward-W. Suor	Muhammad Ziad
Kevin Menezes	Chris Perkins	Jeroen Van	Vincent Surillo	Ziayuddin
Bart Jan Menkveld	Michael Petry	IngenSchenau	Terence Charles Sweetser	Tom Zingale
Sean Mentzer	Alexander Peuchert	Carsten Scherb	T2Group	Jose Zumalave
Eduard Metz	David Phelan	Ernest Schirmer	Roman Tarasov	Romeo Zwart
William Mills	Harald Pilz	Benson Schliesser	David Theese	廖明沂.

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at **ole@protocoljournal.org** or **olejacobsen@me.com**

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Internet
Society



Diamond Sponsors

Your logo here!

Ruby Sponsors



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

John Crain, Senior Vice President and Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Principal – Internet Technology, Policy, and Advocacy
The Internet Society

Dr. Jun Murai, Founder, WIDE Project
Distinguished Professor, Keio University
Co-Director, Keio University Cyber Civilization Research Center, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

December 2023

Volume 26, Number 3

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
Introduction to 5G	2
Why ATM Failed.....	22
20 Years of Cellular and Wi-Fi Integration	30
RDRS	37
Fragments	39
Thank You!	40
Call for Papers	42
Supporters and Sponsors	43

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

In our previous issue, we published Part One of “Introduction to 5G” by William Stallings. Part One introduced the standards and specifications that define 5G and described the usage scenarios that 5G supports. Part Two, included in this issue, provides an overview of the structure and function of 5G networks. A third article, on *Network Slicing*, which is closely related to 5G, will be published in a future edition of this journal.

This journal, as well as its predecessor *ConneXions—The Interoperability Report*, has covered numerous networking technologies over the last 35 years. Some of these technologies have become important building blocks for all networks, for example, *Ethernet*, which for more than 50 years has seen further improvements and standardization. We will publish an article on the history of Ethernet in a future issue. Other technologies have emerged, only to later fade into oblivion—an example being *Asynchronous Transfer Mode* (ATM). Our second article, by Craig Partridge, explores the reasons why ATM failed.

Modern smartphones and other mobile devices rely heavily on the use of numerous radio-based technologies such as *Near Field Communication* (NFC), Bluetooth, *Global Positioning System* (GPS), Wi-Fi, and cellular data. Our third article, by Mark Grayson, examines efforts to integrate cellular and Wi-Fi services into a single architecture.

The WHOIS protocol and its associated server were first introduced in 1982 in RFC 812. Described as “...a server ... that delivers the full name, U.S. mailing address, telephone number, and network mailbox for ARPANET users,” the protocol specification was revised and finalized in RFC 3912 in 2004. WHOIS is an essential tool for anyone seeking information about a particular domain registration. Because of personal data protection laws, many ICANN-accredited registrars are now required to redact personal data from WHOIS lookups, yet certain parties may still have a legitimate need to access non-public information. ICANN has recently launched the *Registration Data Request Service* (RDRS) to address this need. Adiel Akplogan describes RDRS in our final article.

Publication of this journal is made possible by the generous support of our donors, supporters, and sponsors. In 2023 we were pleased to welcome *.au Domain Administration Limited* (auDA) and *Flexoptix* as our newest sponsors. If you would like to donate to or sponsor IPJ, please contact us at ipj@protocoljournal.org

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

Introduction to 5G

Part Two: Core Network, Radio Access Network, and Air Interface

by William Stallings

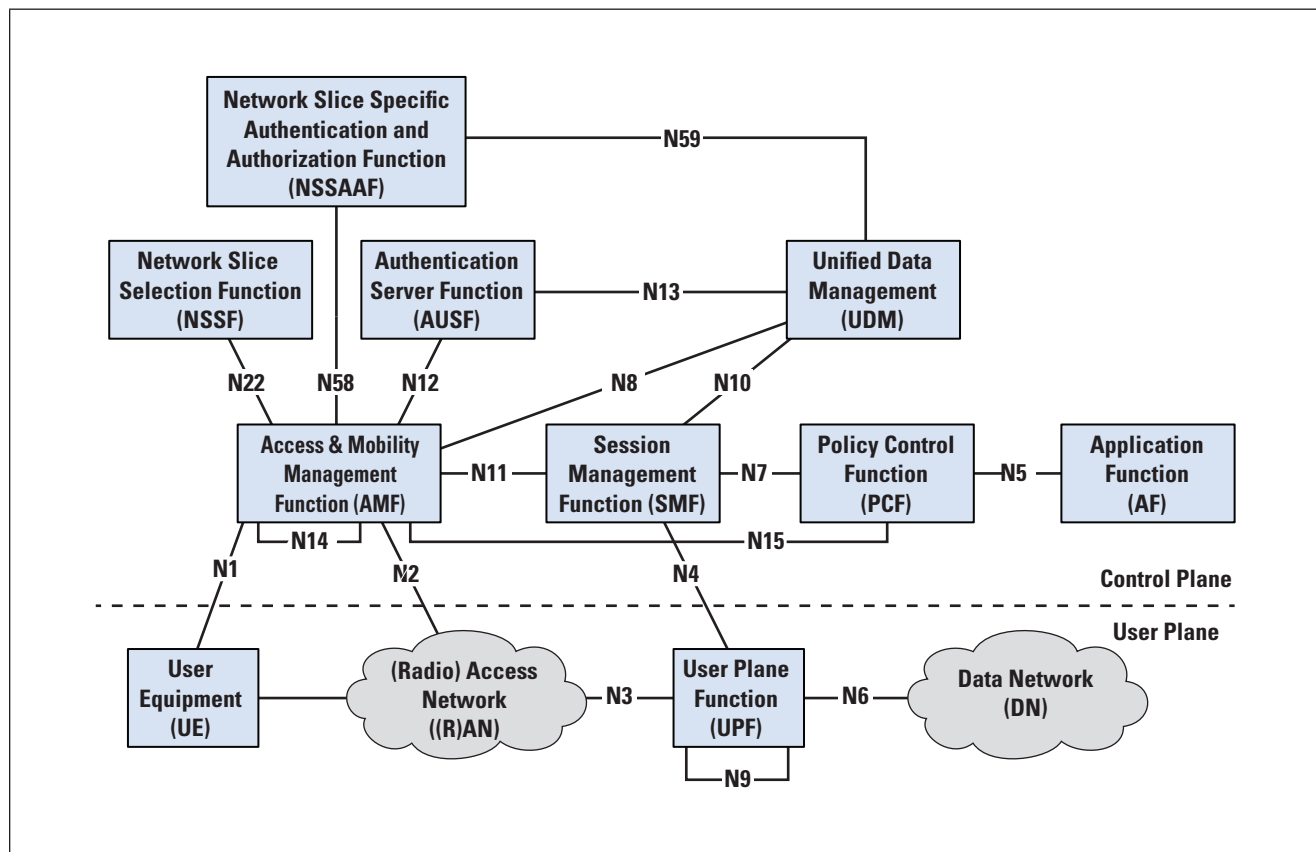
Part One of this 5G introduction^[0] addressed requirements, standards, and applications for 5G. This part provides an overview of the structure and function of 5G networks.

Core Network

Dozens of 3GPP specifications are related to the 5G core network that together describe a system of significant complexity. A key document in this collection is *3GPP Technical Specification TS 23.501*^[1]. This document, 450 pages long in its current release, provides a detailed technical overview of the core network architecture, procedures, services, and interfaces, and is the primary basis for the presentation in this article.

The core network architecture can be viewed as a set of interconnected *Network Functions* (NFs). An NF is a processing function in a network, which has defined functional behavior and interfaces. An NF can be implemented as a network element on dedicated hardware, a software instance running on dedicated hardware, or a virtualized function instantiated on an appropriate platform.

Figure 1: 5G Core Network Functional Architecture



TS 23.501 contains numerous architecture diagrams from several different points of view and at varying levels of detail. Figure 1 depicts the basic 5G architecture using a reference-point representation, showing how the NFs interact with each other.

The figure includes the following NFs and other modules:

- *Authentication Server Function* (AUSF): Performs authentication between *User Equipment* (UE) and the network
- *Access and Mobility Management Function* (AMF): Receives all connection- and session-related information from the UE (N1/N2) but is responsible only for handling connection, registration, reachability, and mobility management tasks. All messages related to session management are forwarded to the *Session Management Function* (SMF).
- *Network Exposure Function* (NEF): Provides an interface for outside applications to communicate with the 5G network to obtain network-related information about the capabilities of the network.
- *Network Repository Function* (NRF): Allows NFs to register their functionality and to discover the services offered by other NFs present in the network.
- *Network Slice Selection Function* (NSSF): Selects the set of network slice instances to accommodate the service request from a UE. When a UE requests registration with the network, AMF sends a network slice selection request to NSSF with preferred network slice selection information. The NSSF responds with a message including the list of appropriate network slice instances for the UE.
- *Network Slice-Specific Authentication and Authorization* (NSSAAF): Performs authentication and authorization specific to a slice.
- *Policy Control Function* (PCF): Provides functionalities for the control and management of policy rules including rules for *Quality of Service* (QoS) enforcement, charging, and traffic routing. PCF enables end-to-end QoS enforcement with QoS parameters (for example, maximum bit rate, guaranteed bit rate, and priority level) at the appropriate granularity (for example, per UE, per flow, and per PDU session).
- *Session Management Function* (SMF): Responsible for *Protocol Data Unit* (PDU) session establishment, modification, and release between a UE and a data network. A PDU session, or simply session, is an association between the UE and a data network that provides a PDU connectivity service. A PDU connectivity service is a service that provides for the exchange of PDUs between a UE and a data network.
- *Unified Data Management* (UDM): Responsible for access authorization and subscription management. UDM works with AMF and AUSF as follows: The AMF provides UE authentication, authorization, and mobility management services. The AUSF stores data for authentication of UEs, and the UDM stores UE subscription data.

- *User Plane Function (UPF)*: Handles the user plane path of PDU sessions. UPF functions include packet routing and forwarding, QoS handling, traffic usage reporting, and policy rule enforcement.
- *Application Function (AF)*: Provides session-related information to PCF so that SMF can ultimately use this information for session management. AF interacts with application services that require dynamic policy control. AF extracts session-related information (for example, QoS requirements) from application signaling and provides it to PCF in support of its rule generation.
- *User Equipment (UE)*: Gives users access to network services. An example is a mobile phone. For the purpose of 3GPP specifications, the interface between the UE and the network is the radio interface.
- *(Radio) Access Network [(R)AN]*: A network that provides access to a 5G core network. It includes the 5G RAN and other wireless and wired access networks.
- *Data Network (DN)*: A network to which UE is logically connected by a session. It may be the Internet, a corporate intranet, or an internal services function within the mobile network operator's core (including content-distribution networks).
- *Service Communication Proxy (SCP)*: NFs and NF services can communicate directly or indirectly via the SCP. The SCP enables multiple NFs to communicate with each other and with user plane entities in a highly distributed multi-access edge compute cloud environment. These services provide routing control, resiliency, and observability to the core network.

In Figure 1, two reference points loop back to the same function: N9 and N14. The N9 reference point is an interface between two distinct UPFs used for forwarding packets. The N14 reference point is between two AMFs, one acting as a source AMF for a data transfer and the other acting as a destination AMF.

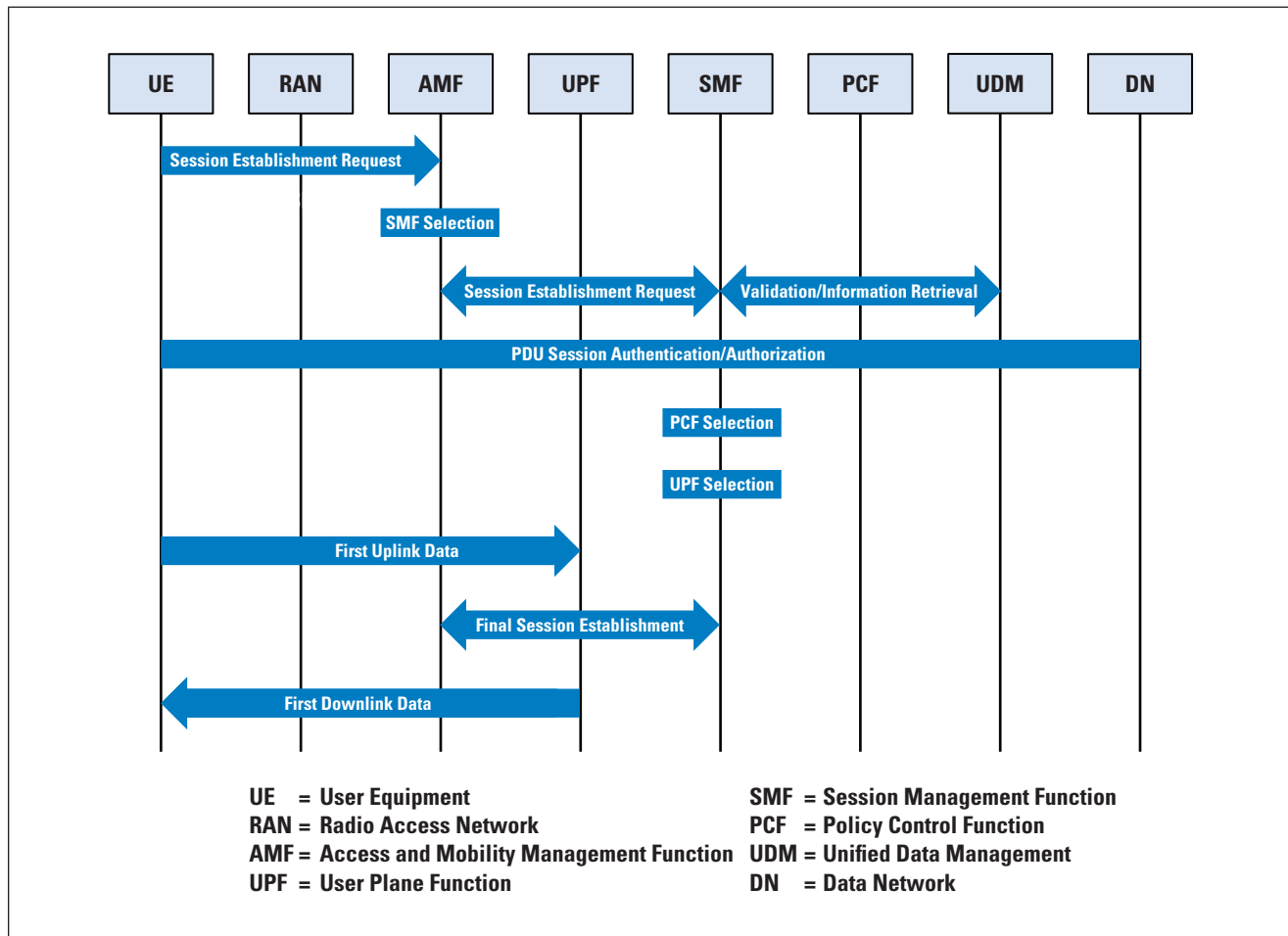
An example of the interaction of the various NFs is the session-establishment procedure, which is defined in TS 23.502^[2]. Figure 2 provides a much-simplified view of the interaction between the various network components during session establishment. Session establishment begins with a request from the UE over the RAN, which is directed to the AMF. An SMF is selected to manage the PDU session. SMF uses UDM in the process of creating a session and performing authentication and authorization. SMF selects a PCF for the session and a UPF to handle data plane PDU forwarding in both directions. SMF establishes a session with the DN. After a few more exchanges, the UE is able to communicate over a session with the DN.

SDN and NFV

Two essential enablers of 5G services provided by core networks are *Software-Defined Networking (SDN)* and *Network Functions Virtualization (NFV)*.^[15]

ITU-T Y.3300^[3] defines SDN as a set of techniques that enables users to directly program, orchestrate, control, and manage network resources, thereby facilitating the design, delivery, and operation of network services in a dynamic and scalable manner.

Figure 2: UE-Requested PDU Session Establishment

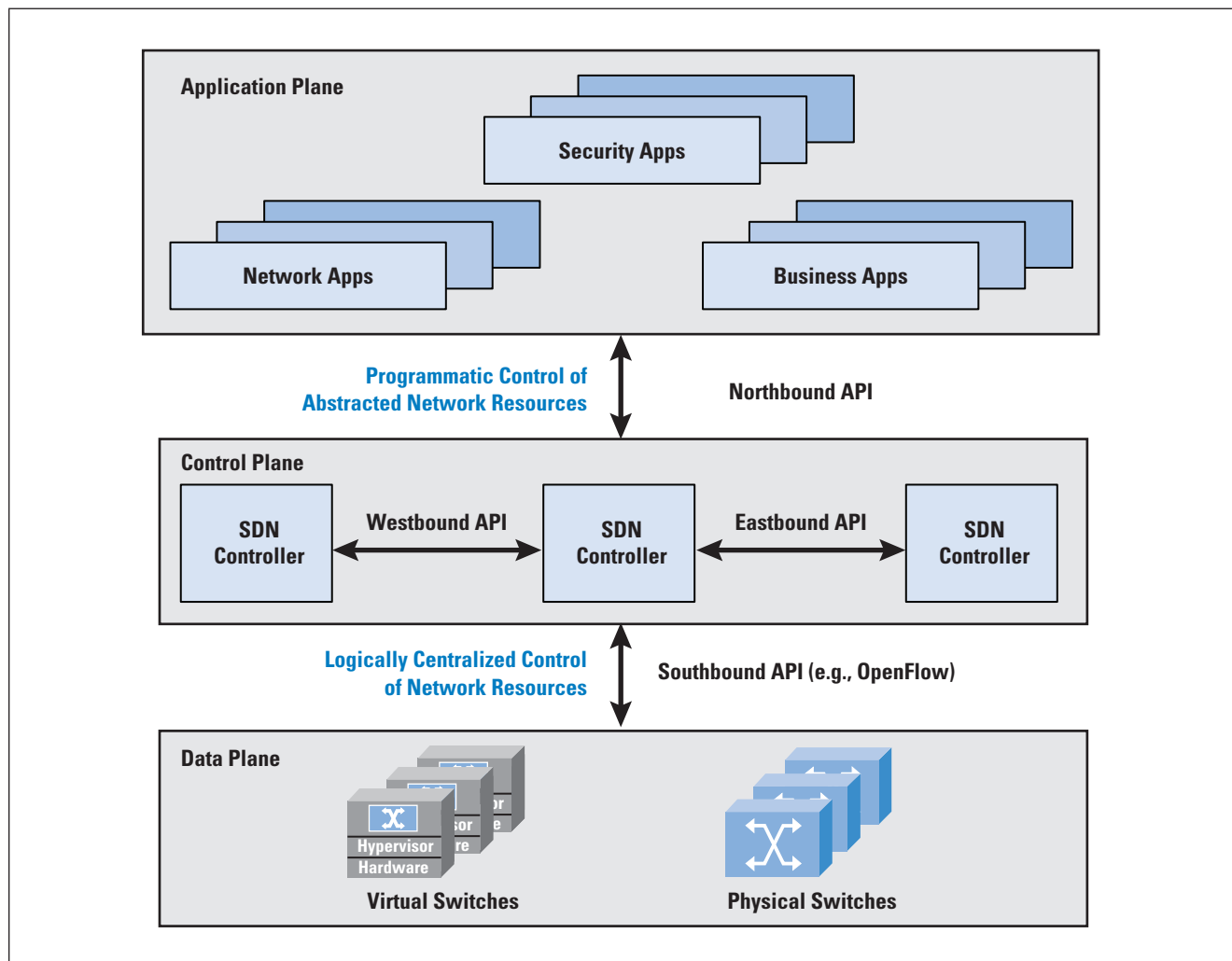


The two elements involved in forwarding packets through routers are a *control function*, which decides the route the traffic takes and the relative priority of traffic, and a *data function*, which forwards data based on control-function policy. Prior to SDN, these functions were performed in an integrated fashion at each network device (router, bridge, packet switch, etc.). Control in such a traditional network is exercised with a routing and control network protocol that is implemented in each network node. This approach is relatively inflexible and requires all of the network nodes to implement the same protocols. With SDN, a central controller performs all complex functionality, including routing, naming, policy declaration, and security checks. This central controller constitutes the SDN control plane, and consists of one or more SDN controllers. The SDN controller defines the data flows that occur in the SDN data plane. Each flow through the network is configured by the controller, which verifies that the communication is permissible by the network policy.

If the controller allows a flow requested by an end system, it computes a route for the flow to take, and adds an entry for that flow in each of the switches along the path. With all complex functions subsumed by the controller, switches simply manage flow tables whose entries can be populated only by the controller. The switches constitute the data plane. Communication between the controller and the switches uses a standardized protocol.

Figure 3 illustrates the SDN architecture. The data plane consists of physical switches and virtual switches, both of which are responsible for forwarding packets. The internal implementation of buffers, priority parameters, and other data structures related to forwarding can be vendor-dependent. However, each switch must implement a model, or an abstraction, of packet forwarding that is uniform and open to the SDN controllers. This model is defined in terms of an open *Application Programming Interface* (API) between the control plane and the data plane (that is, the southbound API).

Figure 3: SDN Architecture

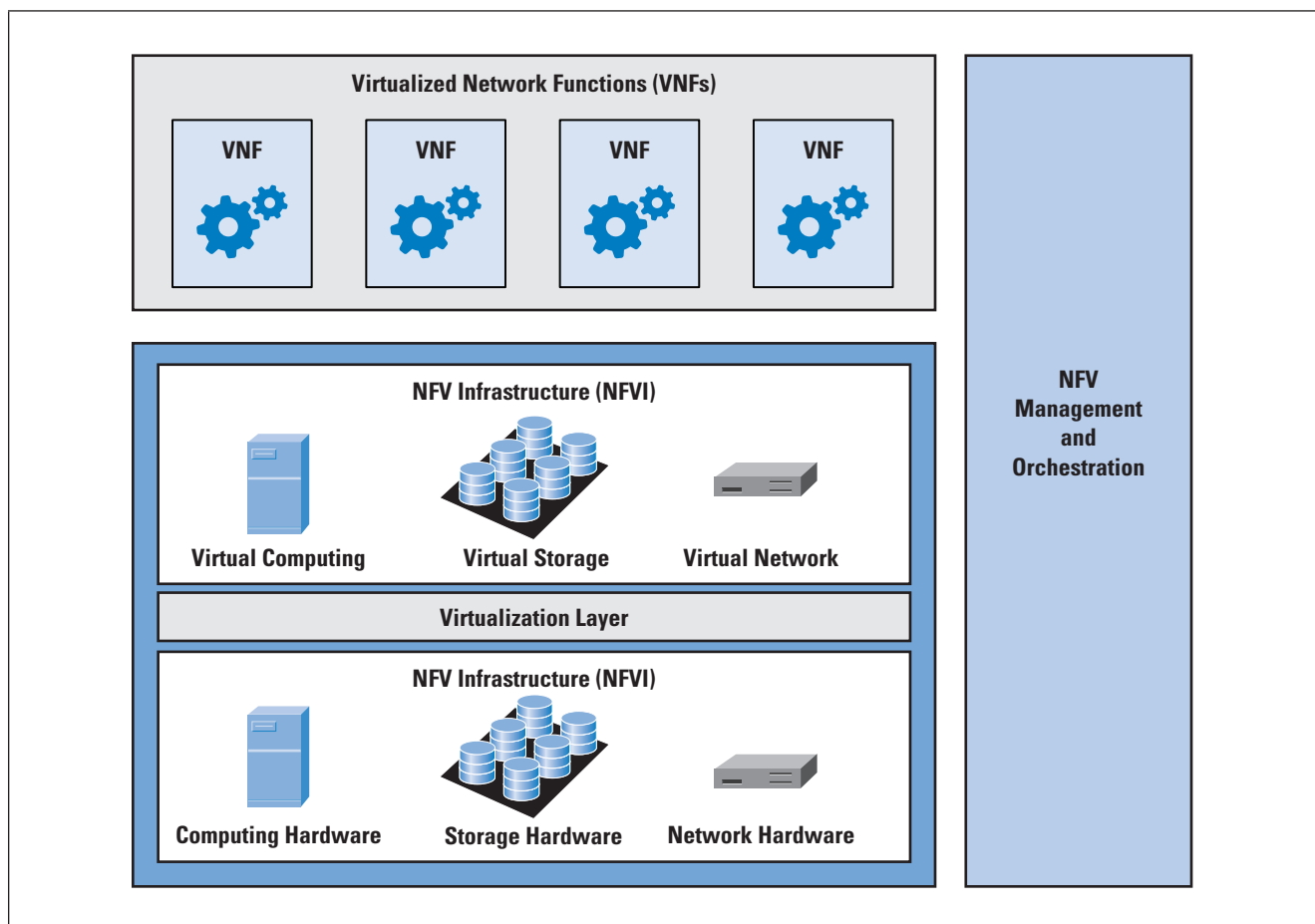


Similarly, SDN controllers can be implemented directly on a server or on a virtual server. An API is used to control the switches in the data plane. In addition, controllers use information about capacity and demand obtained from the networking equipment through which the traffic flows. SDN controllers also expose northbound APIs, meaning that developers and network managers can deploy a wide range of off-the-shelf and custom-built network applications, many of which were not feasible prior to the advent of SDN.

NFV decouples network functions, such as routing, firewalling, intrusion detection, and network address translation, from proprietary hardware platforms and implements these functions in software. It uses standard virtualization technologies that run on high-performance hardware to virtualize network functions. It is applicable to any data plane processing or control plane function in both wired and wireless network infrastructures.

Figure 4 shows a high-level view of the NFV framework, which supports the implementation of network functions as software-only VNFs.

Figure 4: High-Level NFV Framework



The NFV framework consists of three domains of operation:

- *Virtualized Network Functions* (VNFs): These functions consist of a collection of VNFs, implemented in software, run over the *NFV Infrastructure* (NFVI)
- *NFV Infrastructure* (NFVI): The NFVI performs a virtualization function on the three main categories of devices in the network service environment: computer devices, storage devices, and network devices.
- *NFV Management and Orchestration*: This domain encompasses the orchestration and lifecycle management of physical and/or software resources that support the infrastructure virtualization and the lifecycle management of VNFs. NFV management and orchestration focuses on all virtualization-specific management tasks necessary in the NFV framework.

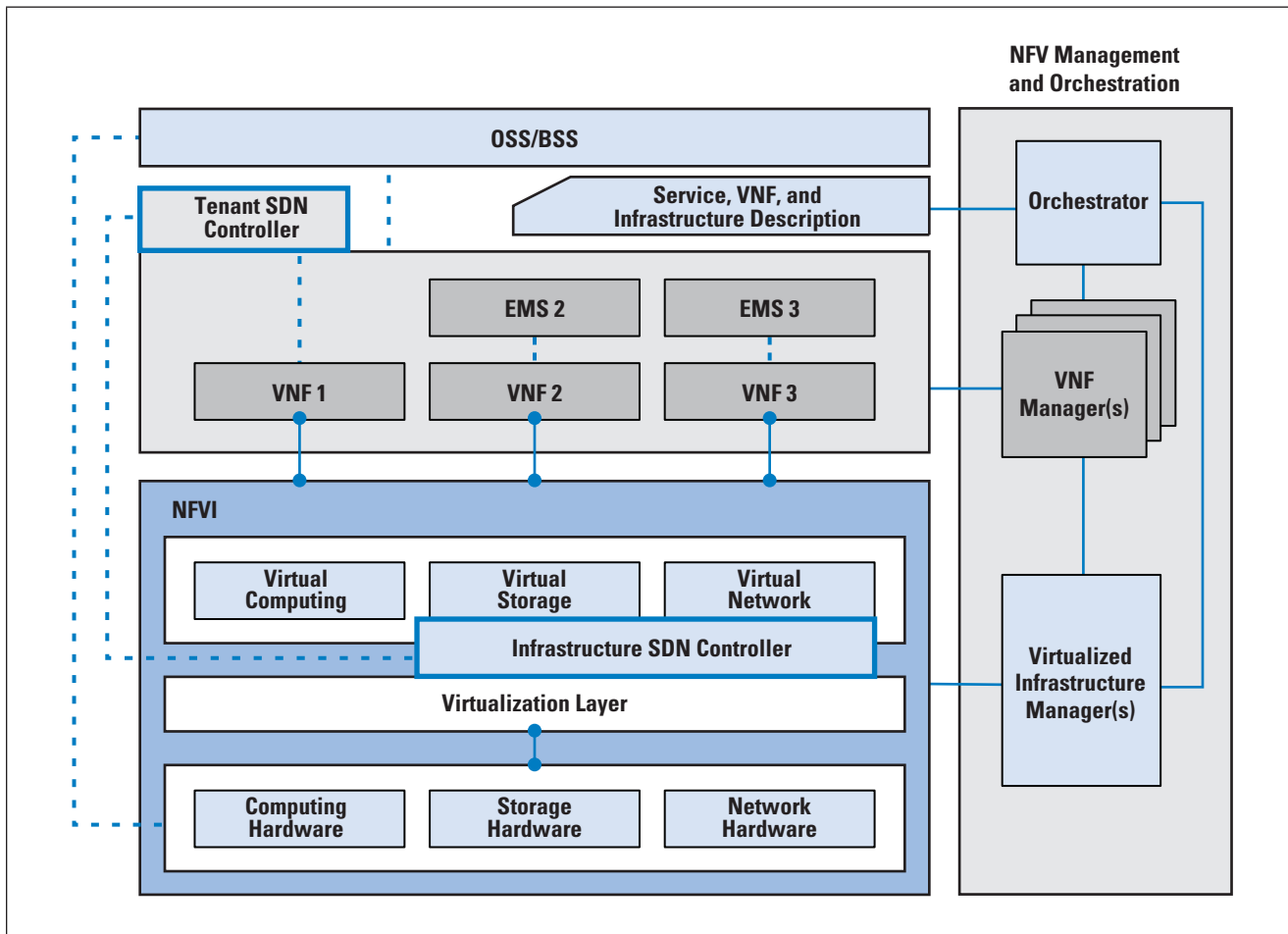
NFV and SDN are independent but complementary schemes. SDN decouples the data and control planes of network traffic control, making the control and routing of network traffic more flexible and efficient. NFV decouples network functions from specific hardware platforms via virtualization to make the provision of these functions more efficient and flexible. Virtualization can be applied to the data plane functions of the routers and other network functions, including SDN controller functions. Thus, either can be used alone, but the two can be combined to reap greater benefits.

The European Telecommunications Standards Institute (ETSI) document GS NFV-EVE 005^[4] examines the manner in which SDN can be incorporated in the NFVI to provide connectivity services.

The framework incorporates two controllers: one logically placed at the tenant level and another at the NFVI level. Each controller centralizes the control plane functionalities and provides an abstract view of all the connectivity-related components it manages. The controllers are as follows:

- *Infrastructure SDN Controller* (IC): This controller enables communication among VNFs and among their components, including the cases when those VNFs are instantiated in separated *Points of Presence* (PoPs), reachable through a WAN connection. Managed by the *Virtualized Infrastructure Manager* (VIM), this controller may change infrastructure behavior on demand according to VIM specifications adapted from tenant requests.
- *Tenant SDN Controller* (TC): Instantiated in the tenant domain as one of the VNFs or as part of the *Network Management System* (NMS), this second controller dynamically manages the pertinent VNFs used to realize the tenant's network service(s). These VNFs are the underlying forwarding plane resources of the TC. The operation and management tasks that the TC carries out are triggered by the applications running on top of it (for example, the *Operations Support System* [OSS]).

Figure 5: Integrating SDN Controllers into the Reference NFV Architectural Framework



Network Slicing

One of the most important features of 5G is *Network Slicing*. Indeed, network slicing is essential to the exploitation of the capabilities defined for 5G. Network slicing enables a 5G network operator to provide customized networks by creating multiple virtual and end-to-end networks, referred to as network slices. Each network slice can be defined according to different requirements on functionality, QoS, and specific users.

“Network Slicing for 5G: Challenges and Opportunities”^[5] lists the following advantages of slicebased networking compared with traditional networking:

- Network slicing can provide logical networks with better performance than one-size-fits-all networks.
- A network slice can scale up or down as service requirements and the number of users change.
- Network slices can isolate the network resources of one service from the others; the configurations among various slices don’t affect each other. Therefore, the reliability and security of each slice can be enhanced.
- A network slice is customized according to QoS requirements, which can optimize the allocation and use of physical network resources.

Network slicing is made possible by NFV and SDN. NFV implements the *Network Functions* (NFs) in a network slice, enabling the isolation of each network slice from all other network slices. Isolation can be achieved by one or more of the following: (1) using a different physical resource; (2) separation by virtualization, which may allow sharing of physical resources; or (3) through sharing a resource with the guidance of a respective policy that defines the access rights for each tenant. Isolation assures QoS and security requirements for that slice, independent of other slices operating on the network from the same or different users. After a network slice is defined, SDN operates to monitor and enforce QoS requirements by controlling the behavior of the QoS flow for each slice.

Figure 6, based on concepts in a *Next Generation Mobile Networks* report^[6], illustrates network slicing concepts.

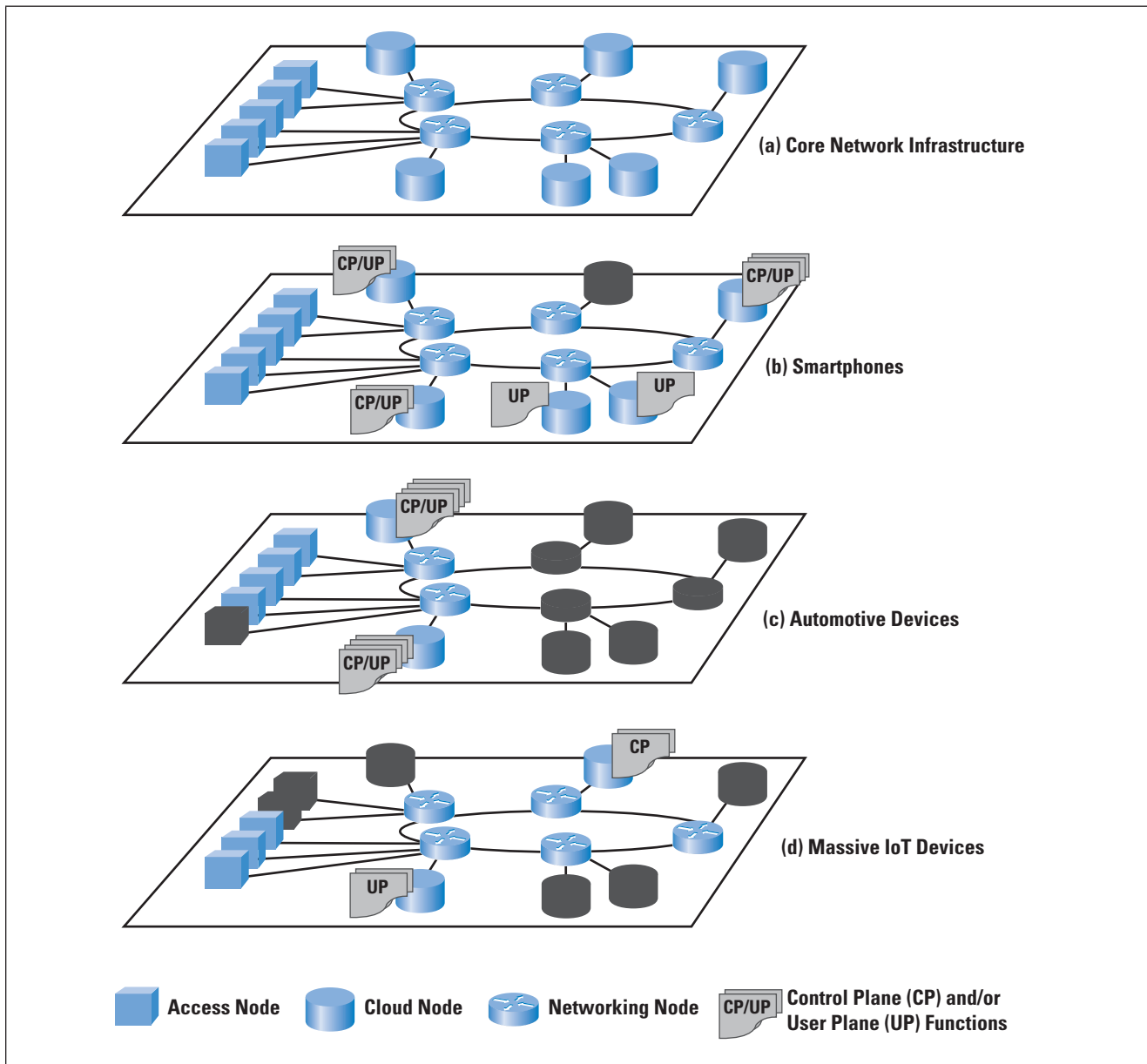
The figure shows a simple core network configuration composed of three types of devices:

1. *Cloud Nodes*: These nodes provide cloud services, software, and storage resources. There are likely to be one or more central cloud nodes that provide traditional cloud computing service. In addition, cloud-edge nodes provide low-latency and higher-security access to client devices at the edge of the network. All of these nodes include virtualization system software to support virtual machines and containers. NFV enables effective deployment of cloud resources to the appropriate edge node for a given application and given fixed or mobile user. The combination of SDN and NFV enables the movement of edge resources and services to dynamically accommodate mobile users.
2. *Networking Nodes*: These nodes are IP routers and other types of switches for implementing a physical path through the network for a 5G connection. SDN provides for flexible and dynamic creation and management of these paths.
3. *Access Nodes*: These nodes provide an interface to radio access networks, which in turn provide access to mobile UE. SDN creates paths that use an access node for one or both ends of a connection involving a wireless device.

The remainder of Figure 6 illustrates three use cases. The blacked-out core network resources represent resources not used to create the network slice. Cloud nodes that are part of the slice may include the following:

- Control plane functions associated with one or more user plane functions (for example, a reusable or common framework of control)
- Service- or service category-specific control plane and user plane function pairs (for example, a user-specific multimedia application session)

Figure 6: 5G Network Slices Implemented on the Same Infrastructure



The first network slice depicted in Figure 6 is for a typical smartphone use case. Such a slice might have fully fledged functions distributed across the network. The second network slice in the figure indicates the type of support that may be allocated for automobiles in motion. This use case emphasizes the need for security, reliability, and low latency. A configuration to achieve these needs would limit core network resources to nearby cloud edge nodes, plus the recruitment of sufficient access nodes to support the use case. The final use case illustrated in Figure 6 is for a massive *Internet of Things* (IoT) deployment, such as a huge number of sensors. The slice can contain just some specific *Control Plane* (CP) and *User Plane* (UP) functions with, for example, no mobility functions. The CP and UP functions might include filtering and preliminary data analysis at the edge and big data types of analysis at a more central node. This slice would need to engage only access nodes nearest to the IoT device deployment.

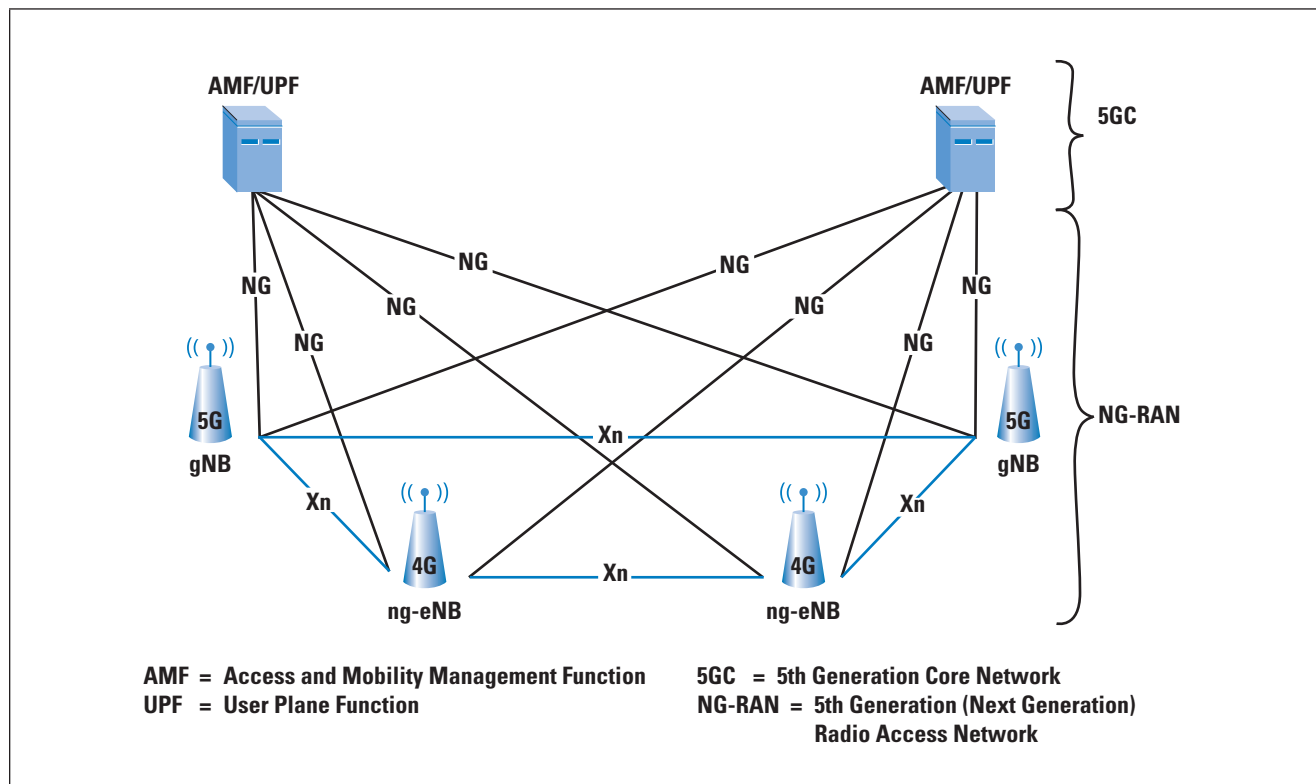
Radio Access Network

As for the 5G core network, there are dozens of 3GPP specifications related to the 5G RAN. A key document in this collection is *3GPP Technical Specification TS 38.300*^[7]. This document provides a detailed technical overview of the RAN architecture, protocols, functionality, and interfaces, and is the primary basis for the presentation in this article.

The overall RAN architecture, in terms of the deployment of base-station RAN nodes, is dictated by the need to coexist with 4G UE and 4G core networks for an extended period. In 2019, 4G became the dominant mobile technology across the world, with more than 4 billion connections, accounting for 52% of total connections (excluding licensed cellular IoT). 4G connections will continue to grow for the next few years, peaking at just under 60% of global connections by 2023^[8]. It is clear that 4G UE will form a major portion of the cellular demand for quite a few years to come.

The most important requirement for 5G carriers is to provide full support for both 4G and 5G UE. Figure 7, from TS 38.300, is a simplified view of the overall RAN architecture and its interface to the 5G core network for providing that support. The figure depicts two types of base stations. The gNB node provides 5G *New Radio* (NR) user plane and control plane protocol terminations toward the UE and connects via the NG interface to the 5GC (5G core). The ng-eNB node provides 4G, referred to as *Evolved Universal Terrestrial Radio Access* (E-UTRA) user plane and control plane protocol terminations toward the UE and connects via the NG interface to the 5GC.

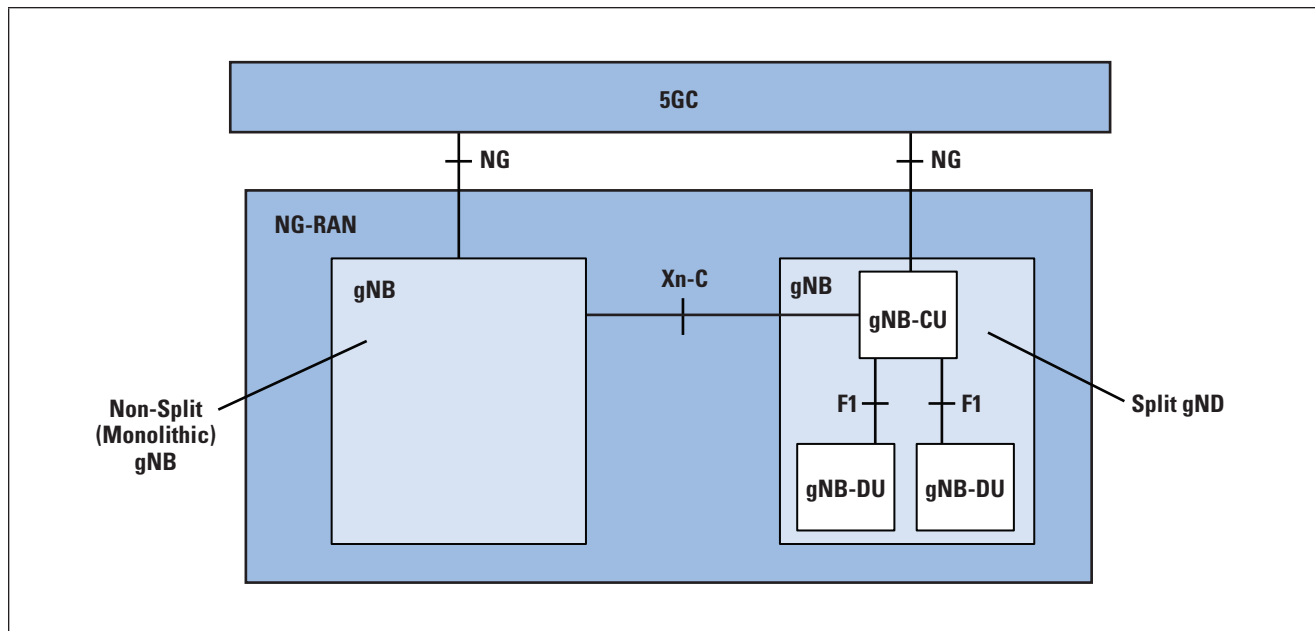
Figure 7: Overall Radio Access Network Architecture



The base stations interconnect with each other by means of the Xn interface. Base stations connect to the core network through the NG interfaces, more specifically to the AMF (access and mobility management function) by means of the NG-C interface and to the UPF by means of the NG-U interface.

Figure 8, from TS 38.401^[9], provides a different perspective on key 5G RAN interfaces. A gNB may be a single integrated system, referred to as a *monolithic* or *non-split* node. Or a gNB may be organized as a split node, consisting of a *gNB-Central Unit* (gNB-CU) and one or more *gNB-Distributed Units* (gNB-DUs). The CU processes non-real-time protocols and services, and the DU processes physical layer protocol and real-time services. One gNB-DU supports one or multiple cells. One cell is supported by only one gNB-DU. A gNB-CU and the gNB-DU units are connected via the F1 logical interface. One gNB-DU is connected to only one gNB-CU. For resiliency, a gNB-DU may also be able to connect to another gNB-CU (if the primary gNB-CU fails) through appropriate implementation. NG, Xn, and F1 are logical interfaces.

Figure 8: RAN Interfaces



Air Interface

As with other aspects of 5G, dozens of 3GPP specifications are related to the 5G RAN. However, the definitive document is *ITU-R Recommendation M.2150*^[10], issued in February 2021. The current version of the Recommendation adopts three radio interface technologies: 3GPP 5G-SRIT, 3GPP 5GRIT, and 5Gi (India/TSDSI). However, the 5Gi specification is unlikely to achieve widespread adoption outside of India^[11]. Accordingly, the coverage in this article of the air interface standards is based on the 3GPP specifications. This article summarizes three key aspects of the air interface: antennas, physical layer, and channel coding.

Antennas

5G systems use *Multiple-Input/Multiple-Output* (MIMO) antenna systems extensively. Key features are base-station antennas consisting of large arrays of antennas and the use of *Beamforming*, and *Beam Management*.

In a MIMO scheme, the transmitter and receiver employ multiple antennas. The source data stream is divided into n substreams, one for each of the n transmitting antennas. The individual substreams are the input to the transmitting antennas (multiple input). At the receiving end, m antennas receive the transmissions from the n source antennas via a combination of line-of-sight transmission and multipath. The output signals from all of the m receiving antennas (multiple output) are combined. With a lot of complex math, the result is a much better receive signal than can be achieved with either a single antenna or multiple frequency channels. Note that the terms *input* and *output* refer to the input to the transmission channel and the output from the transmission channel, respectively.

MIMO systems are characterized by the number of antennas at each end of the wireless channel. Thus, an 8×4 MIMO system has eight antennas at one end of the channel and four at the other end. In configurations with a base station, the first number typically refers to the number of antennas at the base station. There are two types of MIMO transmission schemes:

1. *Spatial Diversity*: The same data is coded and transmitted through multiple antennas, effectively increasing the power in the channel proportionally to the number of transmitting antennas. This mechanism improves the *Signal-to-Noise Ratio* (SNR) for cell edge performance. Further, diverse multipath fading offers multiple “views” of the transmitted data at the receiver, thus increasing robustness. In a multipath scenario where each receiving antenna would experience a different interference environment, there is a high probability that if one antenna suffers a high level of fading, another antenna will have sufficient signal level.
2. *Spatial Multiplexing*: A source data stream is divided among the transmitting antennas. The gain in channel capacity is proportional to the available number of antennas at the transmitter or receiver, whichever is less. Spatial multiplexing can be used when transmitting conditions are favorable and for relatively short distances compared to spatial diversity. The receiver must do considerable signal processing to sort out the incoming substreams, all of which are transmitting in the same frequency channel, and to recover the individual data streams.

Beamforming is one of the essential technologies in developing advanced cellular antenna systems. Beamforming is a technique by which an array of antennas can be steered to transmit radio signals in a specific direction. Rather than simply broadcasting energy/signals in all directions, the antenna arrays that use beamforming determine the direction of interest and send/receive a stronger beam of signals in that specific direction.

In this technique, each antenna element is fed separately with the signal to be transmitted. The phase and amplitude of each signal are then added constructively and destructively in such a way that they concentrate the energy into a narrow beam or lobe. The various transmitted signals merge in the air by normal coherence of the electromagnetic waves, thereby forming a virtual beam in a predetermined direction. To understand how this procedure works, consider a signal that is fed to different antenna elements shifted in phase different amounts for each element. Now picture the transmitted energy from each element at an angle of 45° . At any point along that 45° line, the distance traveled by electromagnetic waves from different antenna elements is not equal. If the phase shifting is such that at 45° signals from all antenna elements arrive at the same phase, then the beam is strongest in that direction.

Beam Management refers to techniques and processes used to achieve the transmission and reception of data over relatively narrow beams. Beamforming and beam management are essential for using the *millimeter-wave* (mmWave) region over the 5G air interface. Narrow beams are needed to compensate for high path loss and blockage. With the use of narrow beams, and especially if the UE is mobile, beam management provides the means for both the base-station antenna and the onboard UE antenna to “lock on” to a beam that provides an optimal path from transmitter to receiver.

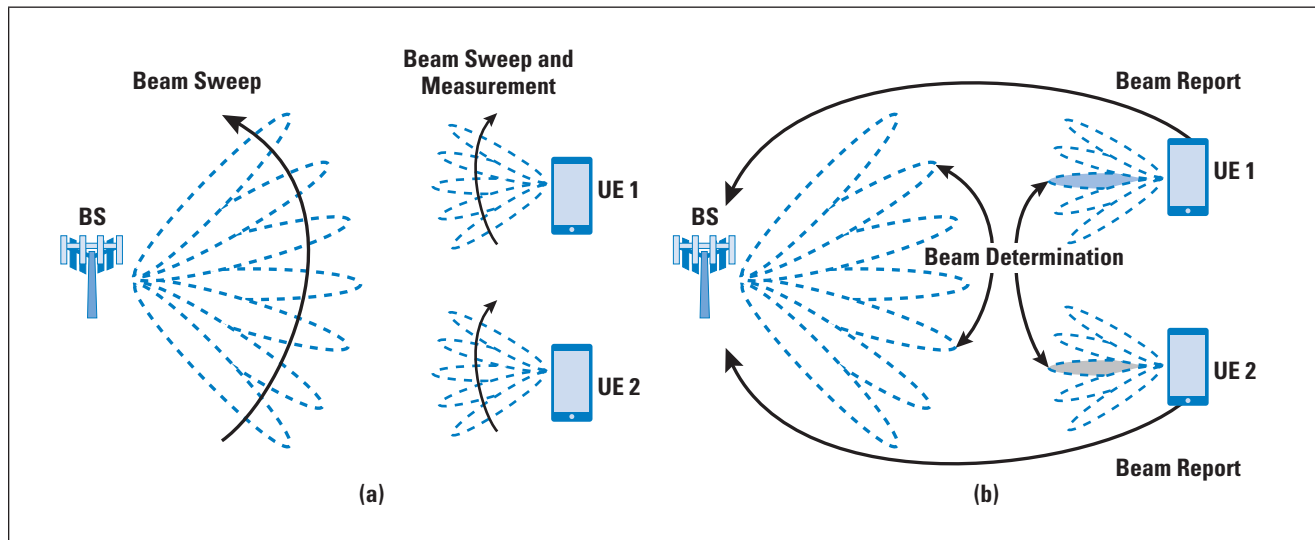
By adjusting the phase and amplitude parameters, a MIMO antenna can generate multiple beams, with each beam covering part of the cell area. For downlink transmission, the objective of beam management is to select a transmit beam to a UE so that the UE can receive the signal with the highest power and best SNR. For uplink transmission, the base station attempts to choose the receive beam for a UE with the best receive beamforming gain. Similarly, if the UE antenna system is capable of beamforming, the UE can use the beams to improve link quality.

The beam management procedure involves beamforming, beam sweeping, beam measurement, beam determination, and beam reporting, as shown in Figure 9.

This figure is taken from [12], which indicates the following elements in the context of downlink transmission:

- *Beam Sweeping*: The base-station antenna (that is, the 5G radio access network node gNB) transmits beams in a predetermined sequence for beam measurement at the UE side.
- *Beam Measurement*: The UE measures the characteristics of received beamformed signals.
- *Beam Determination*: The UE selects the optimal beam. In essence, the UE isolates the receive beam, which affords the best reception. The best results are obtained when the transmitting and receiving beam pair is optimal for the location of the UE at the time.
- *Beam Reporting*: The UE reports back to the gNB the information based on beam measurement.

Figure 9: Beam Management Procedures with Downlink Transmissions



Beam management is an ongoing dynamic process that involves selecting an initial beam pair and then modifying the selection as transmit/receive conditions change.

The term *full-dimension MIMO* (FD-MIMO), or *3D-MIMO*, refers to a MIMO antenna system that is capable of varying the direction of a beam in both horizontal (azimuth) and vertical (elevation) dimensions. Thus, FD-MIMO can project a beam in any direction in three-dimensional space. This capability has advantages, especially in dense urban environments. The ability to adjust transmitted beams in the vertical dimension can improve the received signal power of terminals deep inside high-rise buildings and help overcome some of the building penetration loss. FD beamforming is also advantageous in indoor deployments in high-rise buildings, where a single base station may be able to optimize coverage over more than one floor. Such techniques directly increase spectral efficiency.

OFDMA and SC-FDMA

An important access of the air interface is the way multiplexing and multiple access is achieved over a physical transmission channel. Two techniques are employed for the 5G air interface: *Orthogonal Frequency Division Multiple Access* (OFDMA) and *Single-Carrier Frequency Division Multiple Access* (SC-FDMA). These two schemes use the following foundational techniques:

- *Frequency-Division Multiplexing* (FDM): A physical-layer technique in which multiple baseband signals are modulated on different frequency carrier waves and added together to create a composite signal. The effect of FDM is to divide a transmission bandwidth into multiple subchannels, each of which is dedicated to a particular baseband signal.

- *Frequency-Division Multiple Access (FDMA)*: An access method at the data-link layer based on FDM principles, in which different frequency bands are allocated to different data streams, possibly from different users. The data-link layer in each station tells its physical layer to make a bandpass signal from the data passed to it. The signal must be created in the allocated band. There is no multiplexer at the physical layer. The signals created at each station are automatically bandpass-filtered. They are mixed when they are sent to the common channel. FDMA supports demand assignment, in which the assignment of frequency bands to users changes over time.

Orthogonal Frequency Division Multiplexing (OFDM), also called *Multicarrier Modulation*, is dedicated to a single data source. It uses multiple carrier signals at different frequencies, sending some of the bits on each channel. It differs from ordinary FDM in that the individual subcarriers are orthogonal to one another. In essence, signals are orthogonal to one another if the peaks of the power spectral density of each subcarrier occur at a point at which the power of other subcarriers is zero. A result of this property is that adjacent subcarriers can be packed closely together, making efficient use of the bandwidth.

OFDM has several advantages. First, frequency-selective fading affects only some subcarriers and not the whole signal. If the data stream is protected by a forward error-correcting code, this type of fading is easily handled. More importantly, OFDM overcomes *Intersymbol Interference (ISI)* in a multipath environment. ISI has a greater impact at higher bit rates because the distance between bits, or symbols, is smaller. With OFDM, the data rate is reduced by a factor of N , and this reduction increases the symbol time by a factor of N . This increase dramatically reduces the effect of ISI. As a result of these considerations, with the use of OFDM it may not be necessary to deploy equalizers, which are complex devices whose complexity increases with the number of symbols over which ISI is present.

Like OFDM, OFDMA employs multiple closely spaced subcarriers, but for OFDMA, the subcarriers are divided into groups of subcarriers. Each group is referred to as a *subchannel*. The subcarriers that form a subchannel need not be adjacent. In the downlink, a subchannel may be intended for different receivers. In the uplink, a transmitter may be assigned one or more subchannels.

Subchannelization defines subchannels that can be allocated to *Subscriber Stations (SSs)* depending on their channel conditions and data requirements. Using subchannelization within the same time slot, a *Base Station (BS)* can allocate more transmit power to user devices (SSs) with lower SNR, and less power to user devices with higher SNR. Subchannelization also enables the BS to allocate higher power to subchannels assigned to indoor SSs, resulting in better in-building coverage.

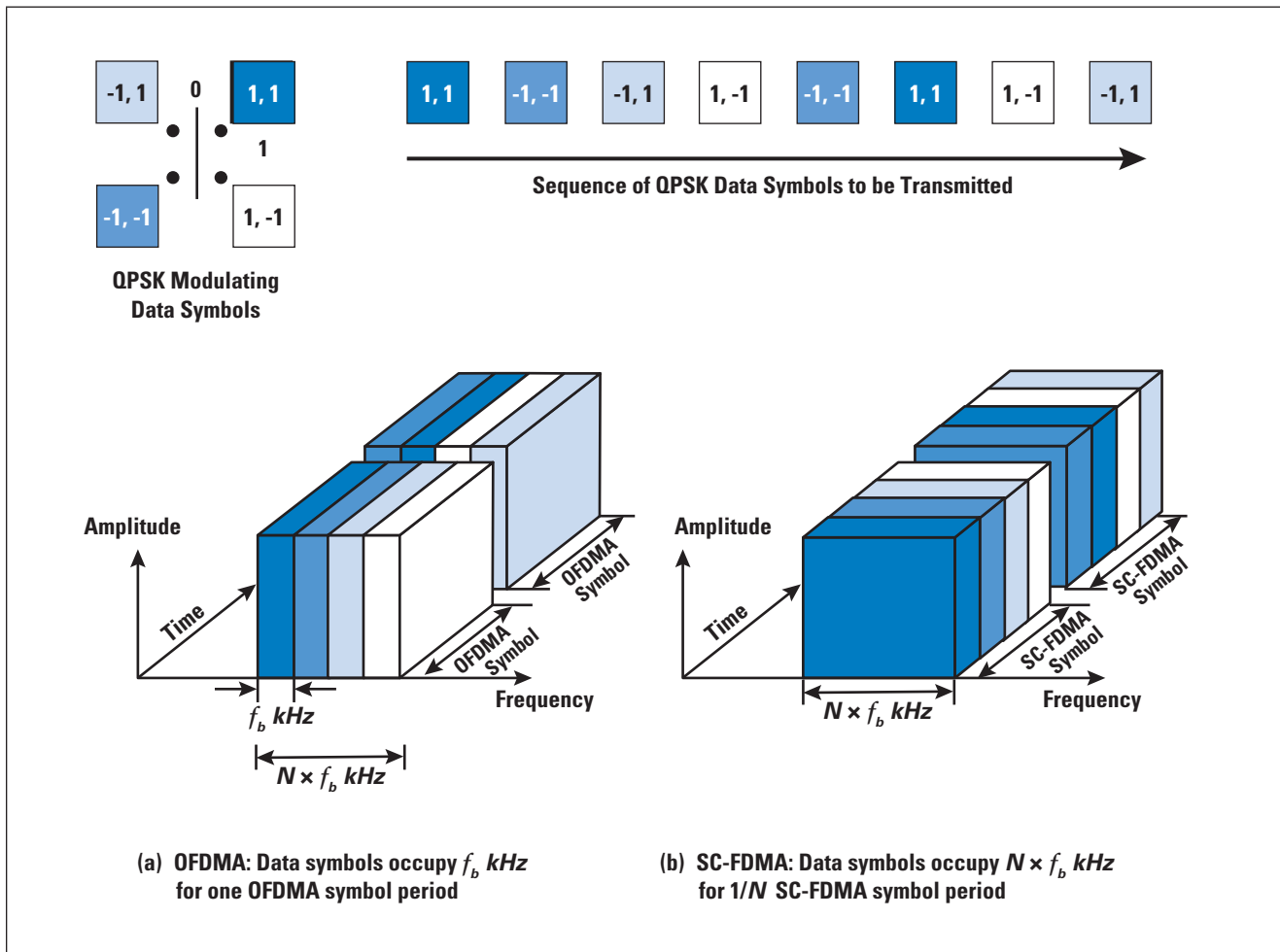
Subchannels are further grouped into *bursts*, which can be allocated to wireless users. Each burst allocation can be changed from frame to frame as well as within the modulation order. This capability allows the base station to dynamically adjust the bandwidth usage according to the current system requirements. Subchannelization in the uplink can save user-device transmit power because it can concentrate power on only certain subchannels allocated to it. This power-saving feature is particularly useful for battery-powered user devices, the likely case in mobile 4G and 5G.

SC-FDMA is a more recently developed multiple-access technique that is similar in structure and performance to OFDMA. One prominent advantage of SC-FDMA over OFDMA is the lower *Peak-to-Average Power Ratio* (PAPR) of the transmit waveform, which benefits the mobile user in terms of battery life and power efficiency. OFDM signals have a higher PAPR because, in the time domain, a multicarrier signal is the sum of many narrowband signals. At some time instances this sum is large, and at other times it is small, meaning the peak value of the signal is substantially larger than the average value. Thus, SC-FDMA is superior to OFDMA. However, it is restricted to uplink use because the increased time-domain processing of SC-FDMA would entail considerable burden on the base station.

Figure 10 provides an example of how the OFDM and SC-FDMA signals appear. As the figure illustrates, with OFDM, a source data stream is divided into N separate data streams, and these streams are modulated and transmitted in parallel on N separate subcarriers, each with bandwidth f_b . The source data stream has a data rate of R bps, and the data rate on each subcarrier is R/N bps. For SCFDMA, it appears from Figure 10 that the source data stream is modulated on a single carrier (hence the SC prefix to the name) of bandwidth $N \times f_b$ and transmitted at a data rate of R bps. The data is transmitted at a higher rate but over a wider bandwidth compared to the data rate on a single subcarrier of OFDM. However, because of the complex signal processing of SC-FDMA, the preceding description is not accurate. In effect, the source data stream is replicated N times, and each copy of the data stream is independently modulated and transmitted on a subcarrier, with a data rate on each subcarrier of R bps.

Compared with OFDM, SC-FDMA transmits at a much higher data rate on each subcarrier, but because the same data stream is on each subcarrier, it is still possible to reliably recover the original data stream at the receiver.

Figure 10: Example of OFDMA and SC-FDMA



Channel Coding

3GPP TS 38.212^[13] specifies two *Forward Error Correction* (FEC) techniques for the air interface: *Low-Density Parity-Check Coding* and *Polar Coding*. TR 38.802^[14] contains the results of a study into NR physical-layer aspects, and is useful for understanding the reasoning behind the concepts. An adequate overview of these two techniques is beyond the scope of this article, but a brief overview of the concepts is provided.

An (n, k) *parity-check code* encodes k data bits into an n -bit codeword. Typically, and without loss of generality, the convention used is that the leftmost k bits of the codeword reproduce the original k data bits, and the rightmost $(n - k)$ bits are the check bits. Such a code is defined by a set of $m = (n - k)$ simultaneous linear equations. If there are m linearly independent equations, there will be some set of k of the variables that can be arbitrarily specified such that one can solve for the other $(n - k)$ variables.

A parity-check code that produces n -bit codewords is the set of solutions to the following equations:

$$\begin{aligned}
 h_{11}c_1 \oplus h_{12}c_2 \oplus \dots \oplus h_{1n}c_n &= 0 \\
 h_{21}c_1 \oplus h_{22}c_2 \oplus \dots \oplus h_{2n}c_n &= 0 \\
 &\vdots \\
 h_{m1}c_1 \oplus h_{m2}c_2 \oplus \dots \oplus h_{mn}c_n &= 0
 \end{aligned}$$

...where the coefficients h_{ij} take on the binary values 0 or 1. The specific set of values of h_{ij} define a specific code.

The $m \times n$ matrix $H = [h_{ij}]$ is called the *Parity Check Matrix*. Each of the m rows of H corresponds to one of the individual equations. Each of the n columns of H corresponds to one bit of the codeword. If we represent the codeword by the row vector $c = [c_i]$, then the equation set can be represented as:

$$Hc^T = cH^T = 0$$

A *Low-Density Parity-Check* (LDPC) code is one in which H has a small density of 1s. That is, the elements of H are almost all equal to 0. Hence the designation *low density*. LDPC codes are enjoying increasing use in high-speed wireless specifications, including Wi-Fi, satellite, and cellular. LDPC codes exhibit performance in terms of bit error probability that is very close to the Shannon limit and can be efficiently implemented for high-speed use.

LDPC codes are suitable for larger blocks of data and are used for 5G data channels. Greater efficiency for small blocks of data is achievable with polar codes, and thus they are used for control channels. Polar codes involve a relatively complex mathematical transformation that involves splitting a communication channel into numerous synthetic bit channels, some of which have extremely low bit error rates and the remainder have high bit error rates, with the data bits being sent over the reliable bit channels. The mathematics behind this transformation is fairly complex and is not pursued here.

References and Further Reading

- [0] William Stallings, "Introduction to 5G Part One: Standards, Specifications, and Usage Scenarios," *The Internet Protocol Journal*, Volume 26, No. 2, September 2023.
- [1] 3GPP TS 23.501, "Technical Specification Group Services and System Aspects; System Architecture for the 5G System (5GS); Stage 2 (Release 16)," December 2020.
- [2] 3GPP TS 23.502, "Technical Specification Group Services and System Aspects; Procedures for the 5G System (5GS); Stage 2 (Release 16)," December 2020.

- [3] ITU-T, “Framework of software-defined networking,” ITU-T Recommendation Y.3300, June 2014.
- [4] ETSI, “Network Functions Virtualisation (NFV); Ecosystem; Report on SDN Usage in NFV Architectural Framework,” ETSI GS NFV-EVE 005, December 2015.
- [5] Xin Li, Mohammed Samaka, H. Anthony Chan, Deval Bhamare, Lav Gupta, Chengcheng Guo, and Raj Jain, “Network Slicing for 5G: Challenges and Opportunities,” *IEEE Internet Computing*, September/October 2017.
- [6] Next Generation Mobile Networks Alliance (NGMN Alliance), “5G End-to-End Architecture Framework,” August 2019.
- [7] GPP TS 38.300, “Technical Specification Group Radio Access Network; NR; NR and NG-RAN Overall Description; Stage 2 (Release 16),” January 2021.
- [8] GSM Association, *The Mobile Economy*, published annually, <https://www.gsma.com/mobileeconomy/>
- [9] 3GPP TS 38.401, “Technical Specification Group Radio Access Network; NG-RAN; Architecture Description (Release 16),” September 2020.
- [10] ITU-R, “Detailed specifications of the terrestrial radio interfaces of International Mobile Telecommunications-2020,” ITU-R M.2150, February 2021.
- [11] Manu Kaushik, “India’s own 5G standard could delay its 5G launch,” *Business Today*, February 20, 2021.
- [12] Guosen Yue, Lingjia Liu, Yongxing Zhou, and Jianzhong Zhang, “MIMO Technologies in 5G New Radio,” *GetMobile: Mobile Computing and Communications*, March 2017.
- [13] 3GPP TS 38.212, “Technical Specification Group Radio Access Network; NR; Multiplexing and channel coding (Release 16),” December 2020.
- [14] 3GPP TR 38.802, “Technical Specification Group Radio Access Network; Study on New Radio Access Technology Physical Layer Aspects (Release 14),” September 2017.
- [15] William Stallings, “Network Functions Virtualization,” *The Internet Protocol Journal*, Volume 24, No. 2, July 2021.
- [16] William Stallings, *5G Wireless: A Comprehensive Introduction*, ISBN-13: 9780136767299, Pearson, 2021.

WILLIAM STALLINGS is a consultant, lecturer, and author of over a dozen books on data communications and computer networking. He has a PhD in computer science from M.I.T. He has written numerous books on computer networking and computer architecture. His home in cyberspace is **WilliamStallings.com** and he can be reached at **ws@shore.net**

Why ATM Failed

by Craig Partridge, Colorado State University

In the late 1980s and early 1990s, *Asynchronous Transfer Mode* (ATM) was widely viewed as the new Internet architecture poised to take the place of the nascent Internet and to inaugurate a world-wide high-speed communications infrastructure. It didn't happen. Instead, after several years of uncertainty, the Internet swept ATM to the side and grew into the global infrastructure we know today.

Even today, 30+ years later, there are different views about how and why ATM “failed.” This essay, while acknowledging ATM had to overcome some technical hurdles, argues that the central problem was a fast-moving window of opportunity that was squandered, largely because of poor standards leadership.

Origins of ATM

Jonathan Turner's forward-looking essay “New Directions in Communications (or Which Way to the Information Age?),” published in 1986, is widely viewed as the paper that launched ATM^[1]. Turner examined the growing diversity of applications using data networks—in particular, advent of many-channel cable television. He looked at the rapidly diminishing error rates in transmission networks, thanks to the advent of fiber-optic cables.

Turner predicted that to meet future needs, our communications networks should be designed around high-performance parallel switches optimized for short packets of information sent over low-error links. Building on prior work at Bell Labs, Turner anticipated that to keep overhead (notably headers) in the short packets small, packets would contain small labels that associated the packets with established connections through the network (both point-to-point and point-to-multipoint).

That, simply, was the idea behind Asynchronous Transfer Mode: implement a futuristic network in which data was reliably transported at high speed in small, fixed-size packets called *cells* over connections.

Turner's central insights were right. Today's data communications networks are built around highly reliable transmission networks and make heavy use of high-performance parallel switches. Furthermore, those high-performance parallel switches internally move data in ways akin to Turner's proposed small packets.

When Turner made his predictions, the Internet was entering one of most difficult stages in its explosive growth. In late 1986, the Internet became plagued by congestion and routing collapses and struggled with inadequate network-management capabilities.

The Internet standards process was not up to the task of handling these concurrent challenges and, as a result, it was necessary to create the *Internet Engineering Task Force* (IETF) to coordinate the substantial efforts to convert research into standardizations. It would take until 1989 for the critical issues to get resolved and their solutions standardized.

Thus, in 1986, future-focused people in data networking saw a lot of wisdom in Turner's paper and a level of worry when they looked at the state of the Internet. Planning for a better data-communications future based on Turner's ideas was appealing.

Nonetheless, ATM, the realization of Turner's vision, failed for three reasons. First, ATM suffered from exceptionally poor standards leadership. Second, and due in part to the poor standards leadership, ATM missed the window(s) of opportunity to capture the local-area network market. Finally, ATM had difficulties matching the needs of the emerging wide-area Internet market.

The ATM Choices in 1986

Before delving into where ATM failed, it is useful to look at the futures one could envision for ATM in 1986 after reading Turner's paper. Turner suggested at least three choices, and the choices were not exclusive:

- ATM as the future universal data-networking protocol. This choice was the most intellectually popular one and, furthermore, the one the telephony industry wanted to see. In this plan, ATM, with its high performance, would sweep away the various networking alternatives such as TCP/IP, *Open Systems Interconnection* (OSI), Apple's *AppleTalk*, IBM's *Systems Network Architecture* (SNA), and Novell *Netware*. ATM would be end-to-end, from the wall jack in your home or office to the data center or business. Furthermore, ATM could also support voice and video—in ways the Internet could not yet—thus positioning the telephone industry to retain their existing voice business and take customers from the cable TV industry.
- ATM as a *Local-Area Network* (LAN) technology^[2]. In 1986, the state of local-area networks was poor. The major technology was original Ethernet, which required heavy coax cables. Interconnecting Ethernets was tedious and required hand configuration to avoid routing loops. It was just in 1986 that Digital Equipment Corporation introduced the *Spanning Tree Protocol*, which prevented loops and made it much easier to connect Ethernet segments. Thus, in 1986, LANs looked clunky, and ATM looked like a way to make it easier to build large corporate networks. This approach was particularly appealing to Silicon Valley startups, because it was relatively simple to build a 4- or 8-port ATM switch, and they could envision that if their small switches sold well, they would be positioned to move into the wide-area ATM market when that market matured.

- Finally, you could envision ATM as just-another-link layer over which you ran TCP/IP. At its simplest, the idea was that ATM would be the next-generation *Wide-Area Networking* (WAN) protocol and TCP/IP would run on top of it. This idea was popular primarily within the burgeoning TCP/IP community, but it was also acceptable in Silicon Valley, where the ATM-as-a-LAN product vendors were happy to have customers for their products, even if those customers used TCP/IP. It was, of course, anathema to the telephony community, which sought to use ATM to take control of the growing data-communications market.

Exceptionally Poor Standards Leadership

Turner's vision of the future was particularly appealing to the international telephony community, which had missed the early stages of the data-communications revolution and the cable-television revolution. Telephone companies, many of which were quasi-government owned, were (correctly) concerned that their business model centered on voice communication was going to be destroyed and their business would shrink to managing the fiber and cables over which other companies would make money selling data services.

So, the telephone companies tasked a committee of their standards organization—the *International Telecommunication Union* (ITU)—to make Turner's vision a reality. The *Consultative Committee for International Telegraphy and Telephony* (CCITT) began creating ATM. (In 1992, CCITT became the *International Telecommunication Union Telecommunication Standardization Sector* [ITU-T]).

CCITT was utterly unsuited to create data-communications standards. It was a standards group that drew its expertise largely from telephony laboratories. A senior US data-communications researcher from *Bell Communications Research* ["Bellcore"] (which had spun off from Bell Labs) attended one of the early ATM standards meetings, assessed the data-communications expertise in the room, and promptly announced to his friends that ATM was an acronym for "Another Telephony Mistake."

CCITT launched straight into standards making without doing an architectural review. Turner had assumed a world of high-speed, low-error networks, and those assumptions had architectural implications. For example, in 1990, Julio Escobar and I did a detailed study of how to do error detection and recovery for ATM in fiber-optic networks. When Julio took our results to a CCITT meeting to figure out what we had missed, he was stunned to be welcomed as the first person who had done an analysis.

CCITT assumed that the world in which the telephone, the cable TV, and the computer each had its own distinctive wall plug would persist in ATM. This assumption was at odds with Turner's vision. In his paper, Turner had pointed out that application-specific networks were a mistake.

It was widely understood in the data-communications community that all data was *bits* and, for the most part, those bits did not care about the format or the wire they were transmitted over. But CCITT's member telephone companies wanted to charge different tariffs for different services, so having a uniform protocol for all data was undesirable. Rather the notion was they would charge distinct tariffs for voice, video, and data, and enforce this differentiation by formatting the different types of data in distinct formats in ATM. Accordingly, CCITT went ahead and defined distinct cell formats, called *ATM Adaptation Layers* (AALs), for voice, video, and data. Each AAL was mapped into a standard 53-byte cell format.

The development of the AAL for data was a particular disaster. The standards group initially planned for two data AALs, AAL 3 for digital video and AAL 4 for computer data. They realized the two could be combined, and they created AAL 3/4, which was swiftly standardized, but was so badly designed that data-communications-savvy members of the standards committee consulted with members of the IETF to propose a new AAL for data, which became AAL 5^[3].

Finally, that 53-byte cell size is worth attention. It came about because of a disagreement about data rates. Telephone companies with less-developed networks were planning to offer ATM over 1.5- or 2-Mbps links (US T1, European E0), and they wanted to support voice calls over ATM. To avoid jitter, that support dictated a small cell size. In contrast, companies in the US expected to offer ATM on 155-Mbps (OC-3) or greater high-speed links and wanted larger cell sizes to reduce the cost of fragmenting data into cells. Competing proposals of 16- and 128-byte cells moved to 32 and 64 bytes, respectively, and the compromise was 48 bytes plus a 5-byte header. The result was to meet neither party's goals.

By late 1991, the failings of the ATM standards process were so severe that the emerging ATM vendor community had to step in. The vendors announced an industry-driven standards group, *The ATM Forum*. The Forum effectively took over the ATM standards process.

A Window of Opportunity in the Local Area

In the grand telephony vision, ATM was designed to be a new networking technology, delivered by the telephone company to your office or home. To achieve that dream, ATM needed to win the competition for the home and office network. In 1986, when ATM was first envisioned, its adoption looked eminently possible, because there was no home network (people dialed into their office computer [not the network, but a specific computer] using a modem) and the primary office network technology was the cumbersome original Ethernet.

But change was also coming. The first thin-wire Ethernet standard, 802.3e, was standardized in 1986. The standard for 10BASE-T, which worked over twisted pair, would come in 1990.

With the advent in 1989 of the first multi-port Ethernet switches (from Kalpana), setting up an office Ethernet became a matter of running 4-wire cables to a closet that held an Ethernet switch. Physically that was the same service ATM was planning to offer, and thanks to Ethernet bridging standards, it was easier to install and operate Ethernet than the nascent ATM LANs.

Consumers could see that more Internet-compatible technology was coming soon. Plans for *100BASE-T* (100-Mbps Ethernet) were soon well-known. It came out in 1995. Wireless networking appeared in 1990 (*WaveLAN*), and standardization efforts leading to *Wi-Fi*, which was intentionally made easy to integrate with wired Ethernet, soon followed.

Concurrently, the Internet was booming. Its major technical problems having been resolved, its user population grew 16X between 1990 and 1995. The World Wide Web appeared in 1991. The Internet was built around Ethernet at the edges and long-haul leased lines in the core. If you had joined the Internet revolution, ATM meant changing your working installed technology. If you could convince yourself that the Ethernet growth curve was good enough, then you didn't need ATM LANs.

Thus, by late 1990, ATM had serious market competition at the edge. It was clear the competition was going to increase. There was a narrow opportunity (perhaps already lost) to capture market share and make ATM the networking service at the network edge. The realization the market window was closing helps explain the vendors' desire to fix the ATM standards process and create The ATM Forum in 1991.

In retrospect, ATM never managed to grab that market window, and the advent of 100-Mbps Ethernet slammed the window shut. Corporate customers had bought a handful of ATM switches to see if they might work if Ethernet and Wi-Fi did not evolve fast enough. But Ethernet and Wi-Fi did evolve fast enough. Furthermore, experience with ATM LANs did not make a compelling case for change.

A Last Chance in the Wide Area

While ATM was rapidly losing credibility for the office and home LAN market at the start of the 1990s, it still had a viable potential role as the technology for wide-area networks. The Internet relied on telecommunications companies for its long-haul links. It was entirely possible that those long-haul links could be running ATM.

At the start of the 1990s, no one knew how to build a multigigabit Internet router. The only working devices that could move data at line speed between multiple gigabit links were ATM switches. Admittedly the switches were prototypes and vendors were waiting for ATM standards to finalize the products, but realized versions of those devices existed, whereas a multigigabit router did not.

So, in many ways, a version of Turner's vision was still alive. The difference from 1986 was that rather than seeing the whole network as running something like ATM, the vision was that there would be islands of Ethernets running IP interconnected via ATM (with no routers in the network center because the routers were too slow).

As in the LAN market, emerging technologies and market forces meant ATM needed to move swiftly to grab this opportunity. In defense of the folks working on ATM for the wide area, unlike the LAN market (where the competitive situation was obvious to anyone who wished to look), the competition for the wide area was not widely recognized. It looked as if there was plenty of time to make this vision happen. As a result, just as in the LAN market, the wide-area ATM effort did not move fast enough.

Critical to this portion of the ATM story is the rise of businesses that made their money installing and selling long-haul runs of fiber-optic links. A British company, *Cable and Wireless*, jumped into this business when AT&T was broken up in 1984, and other companies, including what would become Level 3, followed. These companies started by renting fiber-optic links to the telephone carriers, but in the early 1990s they realized there was a possible market selling links to *Internet Service Providers* (ISPs).

The links ran a point-to-point protocol using the *Synchronous Optical Network* (SONET) or *Synchronous Digital Hierarchy* (SDH) protocols (which differed in minor details). SONET/SDH deliver blocks of bytes at a range of speeds from 155 Mbps up to many gigabits.

By mid-1993, there was a draft of a standard for running IP (and other protocols) over ATM. This standard was recognized to be at least a year later than the market needed it. In this case, arguments among vendors had caused the IETF to take too long. Nonetheless, when the standard was issued, the telephone industry could have sold high-speed (155 Mbps and faster) static ATM circuits running to ISPs. There was early adopter interest in such a product. But, for reasons unclear, the telephone industry did not offer a product.

Enterprising technologists realized there was an opening to develop a competing product to meet ISPs' need to be able to transmit IP data-grams over SONET links. In mid-1994, a draft of a standard for the *Point-to-Point Protocol* (PPP) operating over SDH/SONET appeared. PPP over SONET/SDH was an appealing product because ISPs could simply lease a fiber between two points of presence and run PPP. That approach allowed the ISPs to bypass the telephone companies by renting fiber directly and running PPP. An added benefit was that PPP made better use of the link (less overhead) than ATM and was simpler to manage.

At this point, the ATM dream was nearing its end. The only remaining hope was that the TCP/IP protocols would fail to scale cleanly to gigabit speeds.

By 1994, the only challenge remaining for TCP/IP was the development of multi-gigabit (10+ Gbps and faster) Internet routers. But by early 1995, multiple router vendors were indicating to their customers that, while challenges remained, the customers should expect multi-gigabit IP routers to appear.

Conclusion

In 1996, ten years after Turner's paper, ATM as a forward-looking networking protocol was effectively dead. The IETF's IP over ATM working group held its last meeting in March 1996.^[4,5,6] Perhaps more vividly illustrating the situation, a startup named Ipsilon was marketing a product that sought to repurpose ATM switches as IP routers.

ATM had its origins in Jon Turner's clear vision of the future, and it had a 2- to 5-year head start on the protocols it would compete with. This essay suggests that the ATM community squandered that lead. A valid alternative explanation is that ATM was too complex. Creating a high-speed, circuit-switched, data-networking protocol was a difficult problem in the 1980s, and ATM struggled with challenges such as address resolution and congestion control. But those challenges were ultimately solved. I suggest they would have been solved earlier had the ATM community felt more urgency, and thus this essay focuses on missed chances rather than a technical reason for the ATM failure.

References and Further Reading

- [1] Jonathan Turner, "New Directions in Communications (or Which Way to the Information Age?)," in *IEEE Communications Magazine*, Volume 24, No. 10, pp. 8–15, October 1986.
- [2] J. Bryan Lyles and Daniel C. Swinehart, "The emerging gigabit environment and the role of local ATM," in *IEEE Communications Magazine*, Volume 30, No. 4, pp. 52–58, April 1992.
- [3] Daniel H. Greene and J. Bryan Lyles, "Reliability of Adaptation Layers," in *Proceedings of the IFIP WG6.1/WG6.4 Third International Workshop on Protocols for High-Speed Networks*, Stockholm, Sweden, May 13–15, 1992. IFIP Transactions C-9, North-Holland 1993, ISBN 0-444-89925-1.
- [4] Mark Laubach, "Classical IP and ARP over ATM," RFC 1577, January 1994.
- [5] Mark Laubach and Drew Perkins, "IP over ATM Working Group's Recommendations for the ATM Forum's Multiprotocol BOF Version 1," RFC 1754, January 1995.
- [6] Robert G. Cole, David H. Shur, and Curtis Villamizar, "IP over ATM: A Framework Document," RFC 1932, April 1996.
- [7] George Clapp and Mike Zeug, "Components of OSI: Asynchronous Transfer Mode (ATM) and ATM Adaptation Layers," *ConneXions—The Interoperability Report*, Volume 6, No. 4, April 1992.

<https://archive.org/details/ConneXions.06.04>

- [8] Tom Lyon, “Simple and Efficient Adaptation Layer (SEAL),” ANSI Standards Project T1S1.5 AAL, August 1991.
- [9] John T. Lewis, Raymond Russell, Fergal Toomey, Brian McGurk, Simon Crosby, and Ian Leslie, “Practical connection admission control for ATM networks based on on-line measurements,” *Computer Communications*, Volume 21, Issue 17, pp. 1585–1596, November 25, 1998.
- [10] Alexander G. (Sandy) Fraser, “Towards a Universal Data Transport System,” in *IEEE Journal on Selected Areas in Communications*, Volume 1, No. 5, pp. 803–816, November 1983.

CRAIG PARTRIDGE is a professor at Colorado State University (CSU). Prior to coming to CSU, he was Chief Scientist at BBN Technologies. A member of the *Internet Hall of Fame*, Craig was one of the founding members of the *Internet Engineering Steering Group* (IESG), and he designed MX resource records. Craig can be reached at: **craig.partridge@colostate.edu**

Our Privacy Policy

The *General Data Protection Regulation* (GDPR) is a regulation for data protection and privacy for all individual citizens of the *European Union* (EU) and the *European Economic Area* (EEA). Its implementation in May 2018 led many organizations worldwide to post or update privacy statements regarding how they handle information collected in the course of business. Such statements tend to be long and include carefully crafted legal language. We realize that we may need to provide similar language on our website and in the printed edition, but until such a statement has been developed here is an explanation of how we use any information you have supplied relating to your subscription:

- The mailing list for *The Internet Protocol Journal* (IPJ) is entirely “opt in.” We never have and never will use mailing lists from other organizations for any purpose.
- You may unsubscribe at any time using our online subscription system or by contacting us via e-mail. We will honor any request to remove your name and contact information from our database.
- We will use your contact information only to communicate with you about your subscription; for example, to inform you that a new issue is available, that your subscription needs to be renewed, or that your printed copy has been returned to us as undeliverable by the postal authorities.
- We will never use your contact information for any other purpose or provide the subscription list to any third party other than for the purpose of distributing IPJ by post or by electronic means.
- If you make a donation in support of the journal, your name will be listed on our website and in print unless you tell us otherwise.

Lessons Learned from 20 Years of Cellular and Wi-Fi Integration

by Mark Grayson, Cisco Systems

By any measure, cellular and Wi-Fi wireless technologies can be viewed as having fundamentally transformed the way users consume Internet services. The number of cellular users worldwide now exceeds 7 billion^[1], and the number of Wi-Fi devices shipped now exceeds 39 billion^[2]. These numbers reflect the complementary nature of the wireless systems, with the exclusively licensed cellular systems enabling cellular operators to provide Internet access over wide geographic coverage and the shared unlicensed Wi-Fi systems enabling businesses and individuals to provide targeted local-area coverage.

Whereas the default approach is for these two wireless systems to simply co-exist while supporting their respective value propositions, there have been various attempts over the last 20 years to integrate the two technologies into a single “converged” architecture. Over those 20 years, and the different “Gs,” a myriad of architectural approaches has been specified by the *3rd Generation Partnership Project* (3GPP) to integrate Wi-Fi and cellular architectures.

This article looks at some of the key takeaways from the last 20 years of attempting to converge cellular and Wi-Fi systems.

3GPP Specifications

Efforts started in December 2003, when 3GPP approved its first work item for “UMTS-WLAN Interworking”^[3]. The justification behind the work item highlighted the complementary nature of cellular and Wi-Fi deployments:

“WLAN technology can complement 3GPP based networks in deployment environments with high user density and demand for higher data rates. However, in order to provide flexible use of both technologies in these environments and to provide mobility of services between the two technologies it is sensible that some degree of interworking exists between the two technologies/systems.”

Fast forward 20 years and there now have been over a dozen different approaches specified that look to “converge” cellular. These are listed in Table 1, with some of the architectures illustrated in Figure 1.

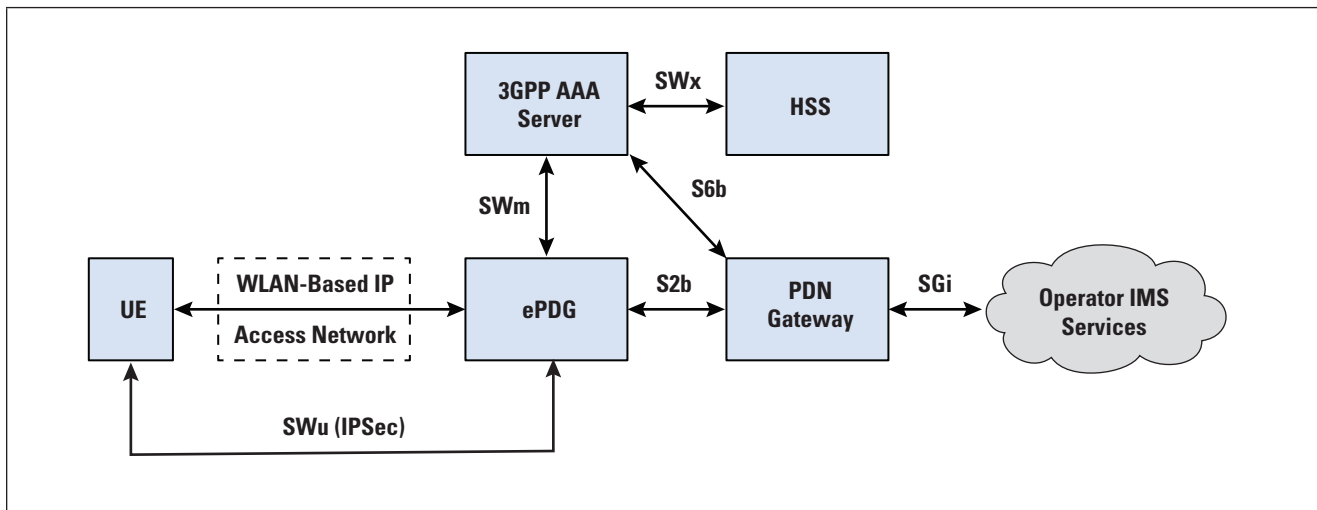
Table 1: A range of 3GPP approaches for converging 3GPP and Wi-Fi systems.

2G-Based	Generic Access Network
3G-Based	Interworking WLAN
4G-Based	Access Network Discovery and Selection Function, evolved Packet Data Gateway, Dual Stack Mobile IPv6, IP Flow Mobility and Seamless WLAN offload, Proxy Mobile IPv6 Trusted WLAN, S2a GPRS Tunnelling Protocol Trusted WLAN, LTE WLAN Aggregation, LTE/WLAN Radio Level Integration with IPsec Tunnel, Network-Based IP Flow Mobility
5G-Based	Non-3GPP Interworking Function, 5G Access Traffic Steering Switching and Splitting

Whereas the majority of these architectures have failed to see wide-scale adoption, the one solution that has seen significant deployment is the evolved *Packet Data Gateway* (ePDG) used to support *Wi-Fi Calling*, illustrated in Figure 1.

Although standardized in 2008^[4], it wasn't until 2014 with the launch of Apple iOS 8 and its native Wi-Fi Calling capability that the functionality became widely adopted, allowing transparent access to *IP Multimedia Subsystem* (IMS)-based rich media communications over both *Long-Term Evolution* (LTE) and Wi-Fi access networks. This type of deployment leverages enhanced *User Equipment* (UE) functionality to use an *IP Security* (IPSec) tunnel between the UE and ePDG to support the IMS-based services. Figure 1 shows the architecture of 3GPP.

Figure 1: 3GPP's ePDG Architecture



GSMA IR.67^[5] and 3GPP TS23.003^[6] have defined a standard realm that mobile operators may use in their Wi-Fi Calling deployments to enable their ePDGs to be discoverable over the public Internet. The *Fully Qualified Domain Name* (FQDN) is of the form:

epdg.epc.mnc<MNC>.mcc<MCC>.pub.3gppnetwork.org

...where <MCC> represents an E.212 *Mobile Country Code* and <MNC> represents the E.212 *Mobile Network Code* allocated to the mobile operator.

Wi-Fi Calling Adoption

In June 2023, the list of allocated MCC and MNC values published at <http://mcc-mnc.com/> was used to determine whether the operator that had been allocated a particular MCC and MNC had configured a *Domain Name System* (DNS) entry to enable its ePDG to be discovered. Table 2 shows the results, which indicate that over 100 countries have deployed Wi-Fi Calling where the ePDG is discoverable using the standard FQDNs defined by 3GPP.

Table 2: Countries where at least one operator has configured a standard DNS entry for ePDG discovery (Source: Cisco Systems)

Albania	Colombia	Iceland	Morocco	Saint Lucia
Anguilla	Croatia	India	Myanmar	Saint Vincent and the Grenadines
Antigua and Barbuda	Cyprus	Indonesia	Namibia	Saudi Arabia
Argentina	Czech Republic	International Networks	Nepal	Singapore
Armenia	Denmark	Ireland	Netherlands	Slovakia
Australia	Dominica	Israel	New Zealand	Slovenia
Austria	Dominican Republic	Italy	Norway	South Africa
Bahamas	Ecuador	Jamaica	Oman	Spain
Bahrain	Egypt	Japan	Pakistan	Sri Lanka
Bangladesh	Estonia	Jordan	Panama	Sudan
Barbados	Faroe Islands	Kazakhstan	Paraguay	Sweden
Belarus	Finland	Kuwait	Peru	Switzerland
Belgium	France	Latvia	Philippines	Taiwan
Brazil	Germany	Liechtenstein	Poland	Thailand
British Virgin Islands	Ghana	Lithuania	Portugal	Türkiye
Brunei	Greece	Luxembourg	Puerto Rico	Turks and Caicos Islands
Bulgaria	Grenada	Malaysia	Qatar	Ukraine
Cambodia	Guadeloupe and Martinique and French Guiana	Maldives	Reunion	United Arab Emirates
Canada	Guam	Monaco	Romania	United Kingdom
Cayman Islands	Hong Kong	Montenegro	Russia	United States of America
Chile	Hungary	Montserrat	Saint Kitts and Nevis	Vietnam

There is clearly a disparity in adoption of the different 3GPP approaches for converging 3GPP and Wi-Fi systems. For instance, in contrast to the over 100 countries that have launched ePDG-based integration, the *Global Mobile Suppliers Association* (www.gsacom.com) reports that only a single operator has invested in LTE WLAN Aggregation.

IMS-based Wi-Fi Calling Observations

Compared to the alternative “trusted” solutions defined by 3GPP for integrating Wi-Fi, the ePDG-based integration can leverage any suitable Wi-Fi network. The result of this leverage may be one of the key reasons that has led to its rapid adoption.

When looking at the Wi-Fi market as a whole, Dell'Oro reports that around 6% of all Wi-Fi equipment revenue is associated with the Service Provider segment^[7]. Only Manufacturing and Logistics segments have lower overall market share, with the Wi-Fi markets for K-12 Education, Higher Education, Finance, Healthcare, Government, Hospitality, and Retail all exceeding the Wi-Fi Service Provider market.

The first lesson learned is to avoid restricting your target market. By enabling all segments deploying Wi-Fi to benefit from ePDG-based integration, the Wi-Fi Calling approach offers the broadest market reach.

The next key observation is that the majority of smartphone data is being sent over connections that use Wi-Fi rather than mobile networks (2G, 3G, 4G, or 5G). The latest data from UK regulator *Ofcom* indicates that nearly three-quarters of all smartphone data is sent over Wi-Fi rather than mobile^[8]. Increasingly these Wi-Fi systems are being dimensioned to deliver gigabit-based services over the fixed network. When comparing data from *Ofcom's* latest *Communications Market Report*^[9], the average volume for fixed broadband, where Wi-Fi dominates, is 453 GB a month, which is 75 times the 6 GB a month for the average cellular subscription.

These figures mean that the smartphone traffic transported over Wi-Fi equates to around 1% of the total fixed broadband traffic that can easily be accommodated. Equally important is the focus of ePDG-based integration on delivering seamless connectivity for IMS-based services, enabling users to receive mobile calls when out of cell tower coverage. With *Ofcom* reporting that the average UK cellphone user calls for 200 minutes of use a month, and a conservative 128 Kbps for the IMS call over Wi-Fi, the impact of Wi-Fi Calling on the cellphone network can be estimated:

$$\begin{aligned} \text{Total 200 minutes (cellular and Wi-Fi)} &\times 75\% = 150 \text{ minutes over Wi-Fi} \\ 150 \times 60 &= 9000 \text{ seconds} \\ 9000 \times 128 \text{ Kb} &= 1.1\text{Gb} = 144 \text{ MB} \end{aligned}$$

Importantly, the 144 MB/month of voice-over-Wi-Fi traffic corresponds to a 2.5% traffic increase compared with the average 6 GB/month used by a cellular subscriber.

This information can be contrasted with other integration approaches that focus on “trusted integration,” where all traffic sent over Wi-Fi is integrated into the cellular provider’s gateway. With 75% of traffic being carried over Wi-Fi, these approaches may result in a 300% increase in traffic across the cellular network. Whether the cellular operators can derive sufficient value from the 300% increase in traffic to cover the additional costs in supporting such is an open issue. However, the increasing adoption of encrypted flows over the Internet has already impacted an operator’s ability to derive value from observing data sent over cellular networks.

The second lesson learned is to avoid thinking of Wi-Fi and cellular as symmetrical services. Wi-Fi is already being dimensioned to support 75 times the traffic load of cellular, and the majority of smartphone traffic continues to be carried on Wi-Fi. Hence, there appears to be advantages to focus integration efforts on systems that avoid transporting the bulk of Wi-Fi data over cellphone networks, such as enabled by IMS-focused ePDG-based integration.

Integrating Native IP-based Services

In the 20 years since 3GPP embarked on the journey to converge 3GPP and Wi-Fi, there has been a significant transition in how Internet services are consumed. Early attempts at convergence were hampered by the binding of sockets to physical interfaces, with applications often stalling as devices made the switch from cellular to Wi-Fi. Hence, initial architectures looked to mask transitions from client-side applications, including the use of Mobile IP client functionality that bound sockets to logical instead of physical interfaces.

In 2023, the Internet is continuing to transition. Not only is over 90% of Internet traffic encrypted, in certain regions of the world we observe that nearly 50% of the traffic has transitioned from regular TCP to HTTP3 transported over *User Datagram Protocol* (UDP)-based QUIC^[10]. Critically, instead of having to mask different paths, the QUIC transport protocol supports native connection migration. Existing connections continue to operate as devices change their endpoint IP addresses when they switch between different networks.

Since 2022, hyperscaler offerings have included native support for HTTP3 and connection migration capability, and the device ecosystem has similarly enabled application developers to benefit from it^[11].

The third lesson learned is to avoid thinking of situations where multiple accesses and multiple paths are available to devices as peculiar. Convergence solutions shouldn't be a "bolt-on" to address specific corner cases. Instead, accept that the Internet has already started its transition to natively support such scenarios.

The Complexities of Path-Selection Policy

Path-selection policy in a heterogeneous environment is a complex issue. Instead of the network-controlled handover approach used in homogeneous cellular networks, the characteristics of Wi-Fi and cellular connections may vary dramatically in terms of costs, quality, and, for moving users, coverage persistence. However, some of the convergence architectures look to expand service providers' cellular network-controlled approach to accommodate Wi-Fi, enabling service providers to define rules that include packet-flow descriptors, access-selection criteria, as well as how to control the steering of flows between Wi-Fi and cellular.

But there is now increasing acceptance that the network provider is but one stakeholder in the complex decision process that is path selection. Identity providers may have preferred relationships that lead them to prioritize the usage of specific paths.

Users are important stakeholders in path selection, including whether the path corresponds to an unmetered private connection or a metered network that may lead to additional charges. Operating System and device vendors can base their preference on near real-time visibility of access networks and associated metrics, and battery levels can be used to drive preference for paths that consume lower energy. Application providers know the metrics that result in the best application experience, whether that is lowest latency for interactive applications or highest throughput for applications that consume significant amounts of data. Applications know in advance the likely duration of application flows and whether it is worth migrating already established flows or waiting for the establishment of new flows over a newly preferred path.

The device ecosystem is looking to meet the needs of their application providers by delivering frameworks that enable applications to configure how multiple paths should be employed, enabling applications developers to easily benefit from the HTTP3 connection migration capability.

The final lesson learned is that a single command-and-control approach to path-selection policy cannot accommodate all stakeholder requirements. And we should recognize that value is continuing to migrate towards the application; application loyalty is the new brand loyalty. So, the goal should be about how to best deliver those application experiences, and what hints and instrumentation can be exchanged between stakeholders to enable better decisions to be made.

Summary

The last two decades have seen significant investment and innovation in the development of cellular and Wi-Fi integration for the delivery of enhanced mobile services. This article has looked at some of the key takeaways from the journey and the lessons learned along the way. In summary: embracing an approach that facilitates integration with the 94% of non-service provider Wi-Fi deployments and leverages the native connection migration support provided by HTTP3, while ensuring that application stakeholders can exchange hints to enable better decisions, is the best route for delivering enhanced services across combined Wi-Fi and cellular networks.

References and Further Reading

- [1] Petroc Taylor, “Forecast number of mobile users worldwide from 2020 to 2025,” *Statista*, January 18, 2023.
<https://www.statista.com/statistics/218984/number-of-global-mobile-users-since-2010/>
- [2] “Value of Wi-Fi,” *Wi-Fi Alliance*,
<https://www.wi-fi.org/discover-wi-fi/value-of-wi-fi>
- [3] “3GPP system - WLAN- Interworking,” 3GPP Technical Specification Group Services and System Aspects, Meeting #22, Maui, Hawaii, USA, 15–18 December 2003,
https://www.3gpp.org/ftp/tsg_sa/TSG_SA/TSGS_22/Docs/PDF/SP-030712.pdf

- [4] “3rd Generation Partnership Project; Technical Specification Group Services and System Aspects; Architecture enhancements for non-3GPP accesses (Release 12),” March 2013,
https://www.3gpp.org/ftp/Specs/archive/23_series/23.402/23402-c00.zip
- [5] “DNS Guidelines for Service Providers and GRX and IPX Providers, Version 21.0,” *GSM Association*, 25 November 2022,
<https://www.gsma.com/newsroom/wp-content/uploads/IR.67-v21.0.pdf>
- [6] “3rd Generation Partnership Project; Technical Specification Group Core Network and Terminals; Numbering, addressing and identification; (Release 18),” March 2023.
https://www.3gpp.org/ftp/Specs/archive/23_series/23.003/23003-i10.zip
- [7] “Market Research Reports on Wireless LAN,” Dell’Oro Group,
<https://www.delloro.com/market-research/enterprise-network-infrastructure/wireless-lan/>
- [8] “Mobile Matters,” *Office of Communications* (Ofcom),
<https://www.ofcom.org.uk/research-and-data/telecoms-research/mobile-smartphones/mobile-matters>
- [9] “Multi-sector Research,” *Office of Communications* (Ofcom),
<https://www.ofcom.org.uk/research-and-data/multi-sector-research/>
- [10] Andreas Enotiadis and Bart Van de Velde, “The New, Encrypted Protocol Stack Taking over the Internet and How to Deal with It,” *Cisco Live*, Las Vegas, June 2023.
<https://www.ciscolive.com/on-demand/on-demand-details.html?#/session/1686177808340001V6Dc>
- [11] Channy Yun, “New – HTTP/3 Support for Amazon CloudFront,” *AWS Blog*, August 15, 2022.
<https://aws.amazon.com/blogs/aws/new-http-3-support-for-amazon-cloudfront/>
- [12] Geoff Huston, “A Quick Look at QUIC,” *The Internet Protocol Journal*, Volume 22, No. 1, March 2019.
- [13] Geoff Huston, “Comparing TCP and QUIC,” *The Internet Protocol Journal*, Volume 25, No. 3, December 2022.

MARK GRAYSON holds an B.Eng. in Electronic and Communication Engineering from the University of Birmingham (England), and a Ph.D. in Wireless Communications from University of Hull. Since 2000, he has worked for Cisco Systems on a variety of wireless-related projects, including early EAP-SIM implementations, 3GPP standardization of Inter-working WLAN and Cisco’s Service Provider Wi-Fi Solutions. He has served on the board of the *Small Cell Forum* where he was responsible for the small cell virtualization efforts that led to the specification of the nFAPI-based split 6. He now serves as co-chair of *O-RAN Alliance’s Fronthaul Working Group*, and is rapporteur for the deliverables dealing with the YANG-based management of the O-RAN Radio Units. He is also the chair of the *Wireless Broadband Alliance’s OpenRoaming Wi-Fi Federation* aimed at lowering barriers to roaming onto private networks and is a Fellow of the *Institute of Engineering and Technology* (IET). He can be reached at: mg@cisco.com

A New and Simplified Way to Request Nonpublic gTLD Registration Data

By Adiel Akplogan, ICANN

The Internet Corporation for Assigned Names and Numbers (ICANN) has recently launched the *Registration Data Request Service* (RDRS). This new service handles requests for access to nonpublic registration data related to *generic Top-Level Domains* (gTLDs). The RDRS is a free and global service that can be an important resource for ICANN-accredited registrars and those who have a legitimate interest in nonpublic data like law enforcement, intellectual property professionals, cybersecurity professionals, consumer-protection advocates, and government officials. The service introduces a more consistent and standardized format for handling these unique requests.

Because of personal data protection laws, many ICANN-accredited registrars are now required to redact personal data from public records such as WHOIS^[1] lookups. With no one way to request or access such data, it can be difficult for interested parties to get the information they need. The RDRS will help ease this problem by providing a simple and standardized process to make these types of requests with benefits for both registrars and requestors.

The RDRS is not only an important tool for the Internet community at large but for the ICANN Board as well. The service was implemented at the direction of the ICANN Board to gather relevant usage data to help inform policy decisions related to a *System for Standardized Access/Disclosure*^[2]. The more registrars and requestors that use the RDRS, the more accurate and valuable the data collected will be toward making that decision. ICANN-accredited registrars are encouraged to opt-in to the service. More information is available at the end of this article.

What Is the RDRS?

The RDRS is a free, global, one-stop-shop ticketing system that handles nonpublic gTLD registration data requests. The RDRS connects requestors of nonpublic data with the relevant ICANN-accredited registrars for gTLD domain names that participate in the service. The service streamlines and standardizes the process for submitting and receiving requests through a single platform. It is important to note that the RDRS does not guarantee access to requested registration data. All communication and data disclosure between the registrars and requestors takes place outside of the system.

Who Can Use the RDRS?

The service is intended for use by ICANN-accredited registrars and individuals and entities with a legitimate interest for access to nonpublic gTLD registration data.

Requestors include but are not limited to: law enforcement, intellectual property professionals, cybersecurity professionals, consumer protection advocates, and government officials. Use by ICANN-accredited registrars is voluntary. More information on how to opt-in to the service is available at the *Naming Services Portal for Registrars*: <https://www.icann.org/resources/pages/nsp-registrars-2018-03-26-en>

Benefits of the Service

One of the key benefits is the simplification of the request process, making it easier to identify the right registrars and provide the necessary information for efficient and timely submission and consideration of disclosure requests. Instead of filling out multiple forms with varying sets of required information, each managed by different registrars, requestors need only to complete a single, standardized form through the service.

Requestors also no longer need to look up the appropriate registrar to contact—the service will do that for them. The service also provides a centralized platform where requestors can conveniently access pending and past requests. They can create new requests, develop request templates for future use, and cancel requests when needed. Registrars can benefit from using the service as it provides a mechanism to manage and track all nonpublic data requests in a single location. Registrars can receive automated alerts anytime they receive a request. The use of a standardized submission form also makes it easier for the correct information and supporting documents to be provided to evaluate a request. For more information on the RDRS, including a flyer for requestors, visit: <https://www.icann.org/rdrs-en>

References

- [1] Leslie Daigle, “WHOIS Protocol Specification,” RFC 3912, September 2004.
- [2] ICANN Generic Name Supporting Organization (GNSO), “Final Report of the Temporary Specification for gTLD Registration Data Phase 2 Expedited Policy Development Process,” <https://gnso.icann.org/sites/default/files/file/field-file-attach/epdp-phase-2-temp-spec-gtld-registration-data-2-31jul20-en.pdf>

ADIEL AKPLOGAN is Vice President, Technical Engagement at ICANN. With more than 25 years of experience in the ICT industry (20 specifically in the Internet Technology Industry), Adiel previously served as CEO for AFRINIC (*The African Network Information Centre*), IT Director for Symbol Technology in France (2001–2003), and Director of New Technology at CAFÉ Informatique in Togo (1994–2000). He earned a graduate degree in Electrical Engineering and holds a Master’s degree in E-Business and New Technology Management from Paris Graduate School of Management. Recognized as one of the Internet technology pioneers in Africa, he contributed to technical capacity building and deployment of some of the first private Internet Service Providers in Africa from 1996 to 1999. He can be reached at: adiel.akplogan@icann.org

APNIC Releases Strategic Plan

The *Executive Council* (EC) of the *Asia Pacific Network Information Centre* (APNIC) is pleased to announce the availability of APNIC's new four-year strategy. The *APNIC Strategic Plan (2024-2027)*^[1] was created by the APNIC EC and Secretariat. It is informed by feedback from Members and the community. The plan sets out the future that APNIC wishes to see, the objectives and priorities that need to be achieved to help reach that future state, and the guiding principles underpinning APNIC's efforts.

The existing strategic pillars of activity (Membership, Registry, Development, Information, and Capability) have been re-cast into four new ones: Two Value Streams, Registry and Development; and two Enablers, Engagement and Capability.

The EC and Secretariat believe the new strategic pillars are the best way to group APNIC's priorities and activities over the coming four years, and the Secretariat is transitioning to a new operational staffing structure to mirror the plan's four pillars.

The strategy becomes the guide for APNIC's annual *Activity Plans*^[2], and the activities will align with the overall strategy. The first annual Activity Plan based on the strategy will be released in March 2024 at the APNIC 57 *Annual General Meeting* (AGM) in Bangkok, held in conjunction with APRICOT 2024: <https://2024.apricot.net/>

[1] https://www.apnic.net/wp-content/uploads/2023/12/APNIC_Strategic_Plan_2024-27.pdf

[2] <https://www.apnic.net/about-apnic/corporate-documents/plans-and-strategies/>

Randy Bush Receives Rob Blokzijl Award

The 2023 *Rob Blokzijl Award* was presented to Randy Bush at the RIPE 87 meeting in Rome in November for his many years of contributions to the Internet in the RIPE NCC service region and beyond, playing a vital role in establishing Internet networks in many developing countries in Africa, Latin America and the Caribbean. The award committee also recognised Randy's non-technical contributions as a dedicated mentor, for speaking the truth, and for passing on knowledge and values.

The award, bestowed by the Rob Blokzijl Foundation, honours the memory of Rob Blokzijl, the first Chair of RIPE. It recognises individuals who have made substantial technical and operational contributions to the development of the Internet in the RIPE NCC service region and supported or enabled others.

You can watch the presentation here: <https://ripe87.ripe.net/archives/video/1145/>

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Kjetil Aas	Lukasz Bromirski	Richard Dodsworth	Serge Van	Richard Johnson
Fabrizio Accatino	Václav Brožík	Ernesto Doelling	Ginderachter	Jim Johnston
Michael Achola	Christophe Brun	Michael Dolan	Greg Goddard	Jonatan Jonasson
Martin Adkins	Gareth Bryan	Eugene Doroniuk	Tiago Goncalves	Daniel Jones
Melchior Aelmans	Ron Buchalski	Michael Dragone	Ron Goodheart	Gary Jones
Christopher Affleck	Paul Buchanan	Joshua Dreier	Octavio Alfageme	Jerry Jones
Scott Aitken	Stefan Buckmann	Lutz Drink	Gorostiaga	Michael Jones
Jacobus Akkerhuis	Caner Budakoglu	Aaron Dudek	Barry Greene	Amar Joshi
Antonio Cuiñat Alario	Darrell Budic	Dmitriy Dudko	Jeffrey Greene	Javier Juan
William Allaire	BugWorks	Andrew Dul	Richard Gregor	David Jump
Nicola Altan	Scott Burleigh	Joan Marc Riera	Martijn Groenleer	Anders Marius Jørgensen
Shane Amante	Chad Burnham	Duocastella	Geert Jan de Groot	Merike Kaao
Marcelo do Amaral	Randy Bush	Pedro Duque	Ólafur Guðmundsson	Andrew Kaiser
Matteo D'Ambrosio	Colin Butcher	Holger Durer	Christopher Guemez	Naoki Kambe
Selva Anandavel	Jon Harald Bøvre	Karlheinz Dölger	Gulf Coast Shots	Christos Karayiannis
Jens Andersson	Olivier Cahagne	Mark Eanes	Sheryll de Guzman	Daniel Karrenberg
Danish Ansari	Antoine Camerlo	Andrew Edwards	Rex Hale	David Kekar
Finn Arildsen	Tracy Camp	Peter Robert Egli	Jason Hall	Stuart Kendrick
Tim Armstrong	Brian Candler	George Ehlers	James Hamilton	Robert Kent
Richard Artes	Fabio Caneparo	Peter Eisses	Darow Han	Thomas Kernen
Michael Aschwanden	Roberto Canonico	Torbjörn Eklöv	Handy Networks LLC	Jithin Kesavan
David Atkins	David Cardwell	Y Ertur	Stephen Hanna	Jubal Kessler
Jac Backus	Richard Carrara	ERNW GmbH	Martin Hannigan	Shan Ali Khan
Jaime Badua	John Cavanaugh	ESdatCo	John Hardin	Nabeel Khatri
Bent Bagger	Lj Cemerias	Steve Esquivel	David Harper	Dae Young Kim
Eric Baker	Dave Chapman	Jay Etchings	Edward Hauser	William W. H. Kimandu
Fred Baker	Stefanos Charchalakakis	Mikhail Evstiounin	David Hauweele	John King
Santosh Balagopalan	Molly Cheam	Bill Fenner	Marilyn Hay	Russell Kirk
William Baltas	Greg Chisholm	Paul Ferguson	Headcrafts SRLS	Gary Klesk
David Bandinelli	David Chosrova	Ricardo Ferreira	Hidde van der Heide	Anthony Klopp
A C Barber	Marcin Cieslak	Kent Fichtner	Johan Helsingius	Henry Kluge
Benjamin Barkin-Wilkins	Lauris Cikovskis	Ulrich N Fierz	Robert Hinden	Michael Kluk
Feras Batainah	Brad Clark	Armin Fisslthaler	Damien Holloway	Andrew Koch
Michael Bazarewsky	Narelle Clark	Michael Fiumano	Alain Van Hoof	Ia Kochiashvili
David Belson	Horst Clausen	The Flirble Organisation	Edward Hotard	Carsten Koempe
Richard Bennett	James Cliver	Jean-Pierre Forcioli	Bill Huber	Richard Koene
Matthew Best	Guido Coenders	Gary Ford	Hagen Hultzs	Alexader Kogan
Hidde Beumer	Robert Collet	Susan Forney	Kauto Huopio	Matthijs Koot
Pier Paolo Biagi	Joseph Connolly	Christopher Forsyth	Asbjørn Højmark	Antonin Kral
Arturo Bianchi	Steve Corbató	Andrew Fox	Kevin Iddles	Robert Krejčí
John Bigrow	Brian Courtney	Craig Fox	Mika Ilvesmaki	John Kristoff
Orvar Ari Bjarnason	Beth and Steve Crocker	Fausto Franceschini	Karsten Iwen	Terje Krogdahl
Tyson Blanchard	Dave Crocker	Erik Fredriksson	Joseph Jackson	Bobby Krupczak
Axel Boeger	Kevin Croes	Valerie Fronczak	David Jaffe	Murray Kuchera
Keith Bogart	John Curran	Tomislav Futivic	Ashford Jaggernaut	Warren Kumari
Mirko Bonadei	André Danthine	Laurence Gagliani	Thomas Jalkanen	George Kuo
Roberto Bonalumi	Morgan Davis	Edward Gallagher	Jozef Janitor	Dirk Kurfuerst
Lolke Boonstra	Jeff Day	Andrew Gallo	Martijn Jansen	Mathias Körber
Julie Bottorff	Fernando Saldana Del	Chris Gamboni	John Jarvis	Darrell Lack
Photography	Castillo	Xosé Bravo Garcia	Dennis Jennings	Andrew Lamb
Gerry Boudreaux	Rodolfo Delgado-Bueno	Oswaldo Gazzaniga	Edward Jennings	Richard Lamb
Leen de Braal	Julien Dhallenne	Kevin Gee	Aart Jochem	Yan Landriault
Kevin Breit	Freek Dijkstra	Rodney Gehrke	Nils Johansson	Edwin Lang
Thomas Bridge	Geert Van Dijk	Greg Giessow	Brian Johnson	Sig Lange
Ilia Bromberg	David Dillow	John Gilbert	Curtis Johnson	Markus Langenmair

Fred Langham	David Millsom	Harald Pilz	Philip Schneck	Kerry Thompson
Tracy LaQuey Parker	Desiree Miloshevic	Derrell Piper	James Schneider	Lorin J Thompson
Alex Latzko	Joost van der Minnen	Rob Pirnie	Peter Schoo	Fabrizio Tivano
Jose Antonio Lazaro	Thomas Mino	Jorge Ivan Pincay	Dan Schrenk	Peter Tomsu Fine Art
Lazaro	Rob Minshall	Ponce	Richard Schultz	Photography
Antonio Leding	Wijnand Modderman-	Marc Vives Piza	Timothy Schwab	Joseph Toste
Rick van Leeuwen	Lenstra	Victoria Poncini	Roger Schwartz	Rey Tucker
Simon Leinen	Mohammad Moghaddas	Blahoslav Popela	SeenThere	Sandro Tumini
Robert Lewis	Charles Monson	Andrew Potter	Scott Seifel	Angelo Turetta
Christian Liberale	Andrea Montefusco	Ian Potts	Paul Selkirk	Michael Turzanski
Martin Lillepuu	Fernando Montenegro	Eduard Llull Pou	Andre Serralheiro	Phil Tweedie
Roger Lindholm	Roberto Montoya	Tim Pozar	Yury Shefer	Steve Ulrich
Link Light Networks	Joel Moore	David Raistrick	Yaron Sheffer	Unitek Engineering AG
Art de Llanos	John More	Priyan R Rajeevan	Doron Shikmoni	John Urbanek
Mike Lochocki	Maurizio Moroni	Balaji Rajendran	Tj Shumway	Martin Urwaleck
Chris and Janet Lonvick	Brian Mort	Paul Rathbone	Jeffrey Sicuranza	Betsy Vanderpool
Sergio Loreti	Soenke Mumm	William Rawlings	Thorsten Sideboard	Surendran Vangadasalam
Eric Louie	Tariq Mustafa	Mujtiba Raza Rizvi	Greipur Sigurdsson	Ramnath Vasudha
Adam Loveless	Stuart Nadin	Bill Reid	Fillipe Cajaiba da Silva	Randy Veasley
Josh Lowe	Michel Nakhla	Petr Rejhon	Andrew Simmons	Philip Venables
Guillermo a Loyola	Mazdak Rajabi Nasab	Robert Remenyi	Pradeep Singh	Buddy Venne
Hannes Lubich	Krishna Natarajan	Rodrigo Ribeiro	Henry Sinnreich	Alejandro Vennera
Dan Lynch	Naveen Nathan	Glenn Ricart	Geoff Sisson	Luca Ventura
David MacDuffie	Darryl Newman	Justin Richards	John Sisson	Scott Vermillion
Sanya Madan	Mai Nguyen	Rafael Riera	Helge Skrivervik	Tom Vest
Miroslav Madić	Thomas Nikolajsen	Mark Risinger	Terry Slattery	Peter Villemoes
Alexis Madriz	Paul Nikolich	Fernando Robayo	Darren Sleeth	Vista Global Coaching
Carl Malamud	Travis Northrup	Michael Roberts	Richard Smit	& Consulting
Jonathan Maldonado	Marijana Novakovic	Gregory Robinson	Bob Smith	Dario Vitali
Michael Malik	David Oates	Ron Rockrohr	Courtney Smith	Rüdiger Volk
Tarmo Mamer	Ovidiu Obersterescu	Carlos Rodrigues	Eric Smith	Jeffrey Wagner
Yogesh Mangar	Jim Oplotnik	Magnus Romedahl	Mark Smith	Don Wahl
John Mann	Tim O'Brien	Lex Van Roon	Tim Sneddon	Michael L Wahrman
Bill Manning	Mike O'Connor	Marshall Rose	Craig Snell	Lakhinder Walia
Harold March	Mike O'Dell	Alessandra Rosi	Job Snijders	Laurence Walker
Vincent Marchand	John O'Neill	David Ross	Ronald Solano	Randy Watts
Normando Marcolongo	Carl Önn	William Ross	Asit Som	Andrew Webster
Gabriel Marroquin	Packet Consulting	Boudhayan	Ignacio Soto Campos	Jd Wegner
David Martin	Limited	Roychowdhury	Evandro Sousa	Tim Weil
Jim Martin	Carlos Astor Araujo	Carlos Rubio	Peter Spekrijse	Westmoreland
Ruben Tripiana Martin	Palmeira	Rainer Rudigier	Thayumanavan Sridhar	Engineering Inc.
Timothy Martin	Gordon Palmer	Timo Ruit	Paul Stancik	Rick Wesson
Carles Mateu	Alexis Panagopoulos	RustedMusic	Ralf Stempfer	Peter Whimp
Juan Jose Marin Martinez	Gaurav Panwar	Babak Saberi	Matthew Stenberg	Russ White
Ioan Maxim	Chris Parker	George Sadowsky	Martin Štěpánek	Jurrien Wijlhuizen
David Mazel	Alex Parkinson	Scott Sandefur	Adrian Stevens	Joseph Williams
Miles McCredie	Craig Partridge	Sachin Sapkal	Clinton Stevens	Derick Winkworth
Brian McCullough	Manuel Uruena Pascual	Arturas Satkovskis	John Streck	Pindar Wong
Joe McEachern	Ricardo Patara	PS Saunders	Martin Streule	Makarand Yerawadekar
Alexander McKenzie	Dipesh Patel	Richard Savoy	David Strom	Phillip Yialeloglou
Jay McMaster	Dan Paynter	John Sayer	Colin Strutt	Janko Zavernik
Mark Mc Nicholas	Leif Eric Pedersen	Phil Scarr	Viktor Sudakov	Bernd Zeimet
Olaf Mehlberg	Rui Sao Pedro	Gianpaolo Scassellati	Edward-W. Suor	Muhammad Ziad
Carsten Melberg	Juan Pena	Elizabeth Scheid	Vincent Surillo	Ziayuddin
Kevin Menezes	Luis Javier Perez	Jeroen Van Ingen	Terence Charles Sweetser	Tom Zingale
Bart Jan Menkveld	Chris Perkins	Schenau	T2Group	Jose Zumalave
Sean Mentzer	Michael Petry	Carsten Scherb	Roman Tarasov	Romeo Zwart
Eduard Metz	Alexander Peuchert	Ernest Schirmer	David Theese	廖明沂.
William Mills	David Phelan	Benson Schliesser	Douglas Thompson	

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at **ole@protocoljournal.org** or **olejacobsen@me.com**

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Diamond Sponsors

Your logo here!

Ruby Sponsors



Sapphire Sponsors

Your logo here!

Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

John Crain, Senior Vice President and Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Olaf Kolkman, Principal – Internet Technology, Policy, and Advocacy
The Internet Society

Dr. Jun Murai, Founder, WIDE Project
Distinguished Professor, Keio University
Co-Director, Keio University Cyber Civilization Research Center, Japan

Pindar Wong, Chairman and President
Verifi Limited, Hong Kong

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

March 2024

Volume 27, Number 1

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
Network Slicing.....	2
Ethernet History	12
Letter to the Editor.....	25
Fragments.....	27
Thank You.....	32
Call for Papers.....	34
Supporters and Sponsors	35

I have just returned from the annual *Asia Pacific Regional Internet Conference on Operational Technologies* (APRICOT), held this year in Bangkok, Thailand. Amongst the many interesting presentations given, there was one entitled “BGP in 2023,” by Geoff Huston. In his talk, he asked if we have reached “Peak IPv4,” noting that the overall IPv4 routing growth trends slowed down or even reversed through 2023. His presentation, as well as a YouTube video, are available on the APRICOT 2024 website.

In our two previous issues, we published a two-part set of articles under the heading “Introduction to 5G” by William Stallings. Part One introduced the standards, specifications, and usage scenarios for 5G. Part Two gave an overview of the structure and function of 5G networks. A third article, on *Network Slicing*, which is closely related to 5G, is included in this edition.

This journal, as well as its predecessor *ConneXions—The Interoperability Report*, has covered numerous networking technologies over the last 35 years. Some of these technologies have become important building blocks for all networks, for example, *Ethernet*, which for more than 50 years has seen further improvements and standardization. Our second article, by Mikael Holmberg, describes the history and future of Ethernet.

Pindar Wong has served on our Editorial Advisory Board since the inception of this journal. I have always appreciated his invaluable insight and advice, particularly on emerging technologies such as Blockchain. Pindar has indicated that he is moving on to pursue other interests and wishes to step down from his advisory role. I thank Pindar for all his contributions and wish him the best in his future endeavors.

I am also extremely honored to welcome Merike Kaeo as a new member of the Editorial Advisory Board. Merike has extensive experience in all aspects of network and information security, and I look forward to working with her on developing article topics for IPJ.

Publication of this journal is made possible by the generous support of our donors, supporters, and sponsors. We also depend on your feedback and suggestions. If you would like to comment on, donate to, or sponsor IPJ, please contact us at ipj@protocoljournal.org

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

Network Slicing

by William Stallings, Independent Consultant

One of the most important features of 5G is *Network Slicing*^[1,10,11]. Network slicing uses virtualization technologies, especially *Software Defined Networks* (SDN) and *Network Functions Virtualization* (NFV)^[0], which enable a 5G network operator to provide customized networks by creating multiple virtual and end-to-end networks, referred to as *network slices*. Each network slice can be defined according to different requirements on functionality, *Quality of Service* (QoS), and specific users.

The article “Network Slicing for 5G: Challenges and Opportunities,”^[2] lists the following advantages of slice-based networking compared with traditional networks:

- Network slicing can provide logical networks with better performance than one-size-fits-all networks.
- A network slice can scale up or down as service requirements and the number of users change.
- Network slices can isolate the network resources of one service from the others; the configurations among various slices don’t affect each other. Therefore, the reliability and security of each slice can be enhanced.
- A network slice is customized according to QoS requirements, which can optimize the allocation and use of physical network resources.

Network slicing is made possible by the “softwarization” techniques of NFV and SDN. NFV implements the *Network Functions* (NFs) in a network slice, enabling the isolation of each network slice from all other network slices. Isolation is achieved by (i) using a different physical resource; (ii) separating by virtualization, which may allow sharing of physical resources; or (iii) sharing a resource with the guidance of a respective policy that defines the access rights for each tenant. Isolation assures QoS and security requirements for that slice independent of other slices operating on the network from the same or different users. After a network slice is defined, SDN operates to monitor and enforce QoS requirements by controlling the behavior of the QoS flow for each slice.

Overview

Network slicing permits a physical network to be separated into multiple virtual networks (logical segments) that can support different *Radio Access Networks* (RANs) or several types of services for certain customer segments, greatly reducing network construction costs by using communication channels more efficiently. In essence, network slicing allows the creation of multiple virtual networks atop a shared physical infrastructure.

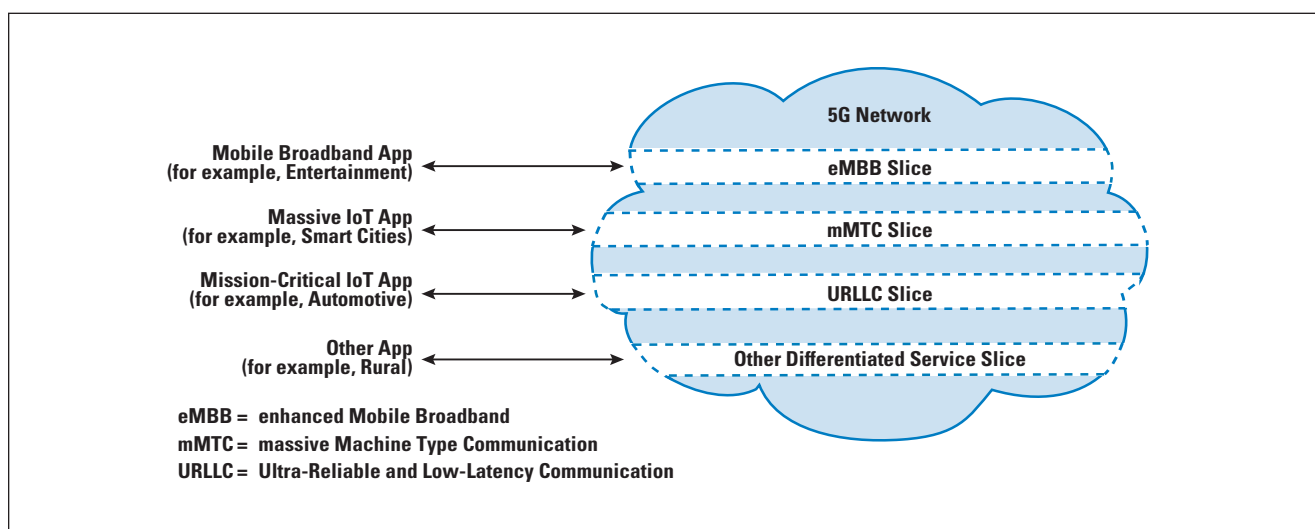
In this virtualized network scenario, physical components are secondary and logical (software-based) partitions are paramount, devoting capacity to certain purposes dynamically, according to need. As needs change, so can the devoted resources. Using common resources such as storage and processors, network slicing permits the creation of slices devoted to logical, self-contained, and partitioned network functions. Network slicing supports the creation of virtual networks to provide a given QoS, such as guaranteed delay, throughput, reliability, and/or priority.

The *International Telecommunication Union Telecommunication Standardization Sector* (ITU-T) is involved in the standardization of network slicing for 5G networks. ITU-T Recommendation Y.3112^[3] defines a network slice as a logical network that provides specific network capabilities and network characteristics. This recommendation lays out an overall framework for network slicing, defines high-level requirements, and describes core network functions relevant to network slicing.

Figure 1 illustrates the network slicing concept. The requirements of a particular application or user determine the physical and logical network resources needed to provide the desired QoS. The network slicing function dedicates the appropriate resources to support that QoS. Figure 1 illustrates the three major usage scenarios for 5G defined by *ITU Radiocommunication Sector* (ITU-R)^[4]. The scenarios include:

- *enhanced Mobile Broadband* (eMBB): Characterized by high data rates for mobile devices.
- *massive Machine-Type Communication* (mMTC): Characterized by the ability to support huge numbers of devices, such as in a large *Internet of Things* (IoT) deployment.
- *Ultra-Reliable and Low-Latency Communication* (URLLC): Characterized by the ability to support human-to-machine and machine-to-machine communications that require high reliability and/or low end-to-end delay.

Figure 1: Network Slicing Concept



Network Slicing Concepts

Network slicing permits you to separate a physical network into multiple virtual networks (logical segments) that can support different radio access networks or several types of services for certain customer segments, greatly reducing network construction costs by using communication channels more efficiently. In essence, network slicing allows you to create multiple virtual networks atop a shared physical infrastructure. This virtualized network scenario devotes capacity to certain purposes dynamically, according to need. As needs change, so can the devoted resources. Using common resources such as storage and processors, network slicing permits you to create slices devoted to logical, self-contained, and partitioned network functions. It supports the creation of virtual networks to provide a given QoS, such as guaranteed delay, throughput, reliability, and/or priority.

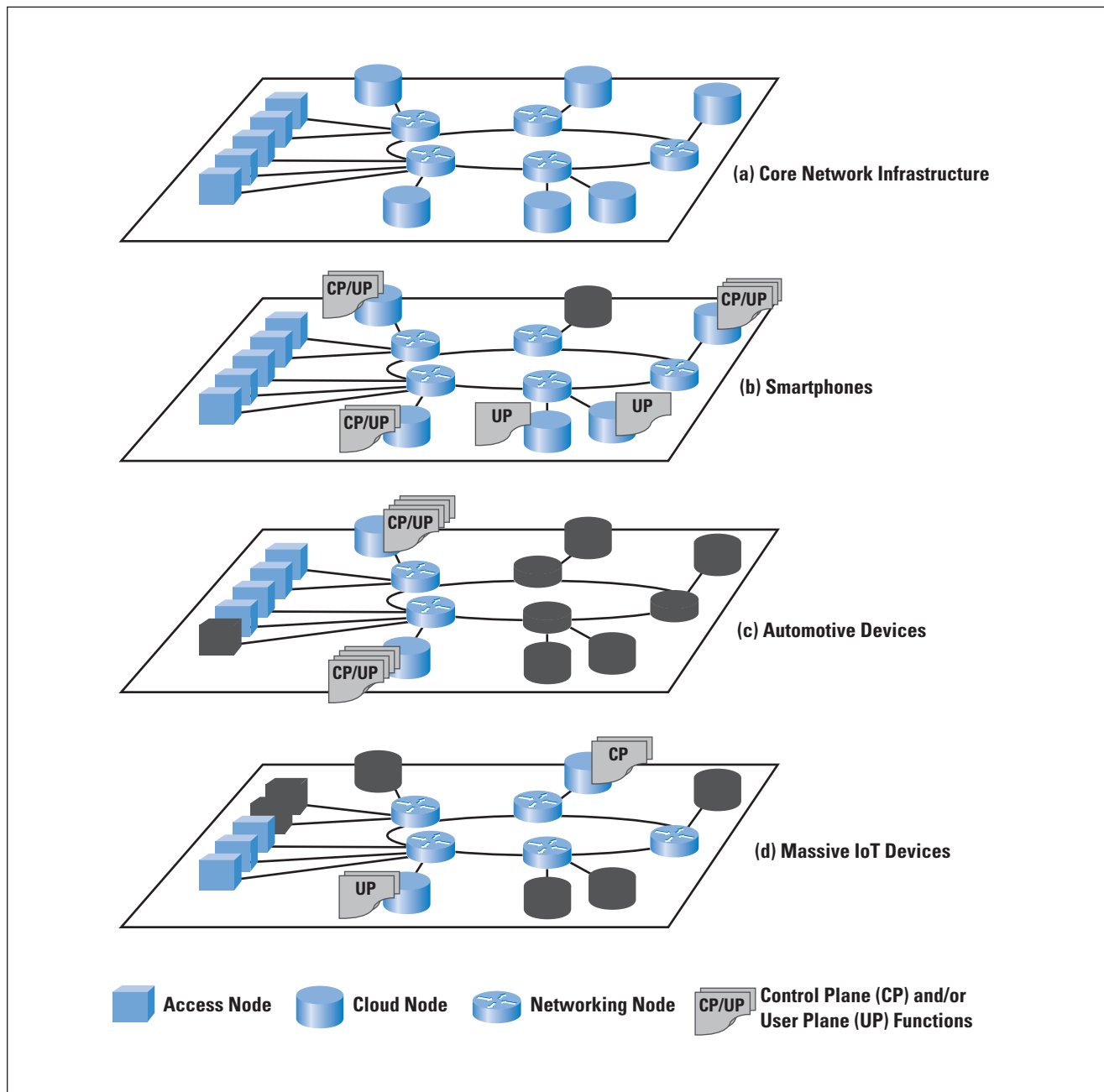
A network slice creates a partition of the core network consisting of virtualized network functions and resources running on some of the core network hardware resources. Figure 2, based on concepts in a *Next Generation Mobile Networks* (NGMN) document^[5], illustrates network slicing concepts. Figure 2a shows a simple core network configuration comprising three types of devices:

- *Cloud Nodes:* These nodes provide cloud services, software, and storage resources. There are likely to be one or more central cloud nodes that provide traditional cloud computing service. In addition, cloud-edge nodes provide low latency and higher security access to client devices at the edge of the network. All of these nodes include virtualization system software to support virtual machines and containers. NFV enables effective deployment of cloud resources to the appropriate edge node for a given application and given fixed or mobile user. The combination of SDN and NFV enables the movement of edge resources and services to dynamically accommodate mobile users.
- *Networking Nodes:* These nodes are IP routers and other types of switches for implementing a physical path through the network for a 5G connection. SDN provides for flexible and dynamic creation and management of these paths.
- *Access Nodes:* These nodes provide an interface to RANs, which in turn provide access to mobile *User Equipment* (UE). SDN creates paths that use an access node for one or both ends of a connection involving a wireless device.

The remainder of Figure 2 illustrates three use cases. The blacked-out core network resources represent resources not used to create the network slice. Cloud nodes that are part of the slice may include the following:

- Control-plane functions associated with one or more user-plane functions (for example, a reusable or common framework of control).
- Service- or service-category-specific control-plane and user-plane function pairs (for example, a user-specific multimedia application session).

Figure 2: 5G Network Slices Implemented on the Same Infrastructure



The first network slice depicted in Figure 2 is for a typical smartphone use case. Such a slice might have fully-fledged functions distributed across the network. The second network slice in Figure 2 indicates the type of support that may be allocated for automobiles in motion. This use case emphasizes the need for security, reliability, and low latency. A configuration to achieve these necessities would limit core network resources to nearby cloud-edge nodes, in addition to recruiting sufficient access nodes to support the use case.

The final use case illustrated in Figure 2 is for a massive IoT deployment, such as a huge number of sensors. The slice can contain just some specific *Control Plane* (CP) and *User Plane* (UP) functions with, for example, no mobility functions. The CP and UP functions might include filtering and preliminary data analysis at the edge and big data types of analysis at a more central node. This slice would need to engage only access nodes nearest to the IoT device deployment.

Requirements for Network Slicing

The *3rd Generation Partnership Project* (3GPP) is the organization responsible for developing specifications that are subsequently issued as ITU-T Recommendations. The 3GPP Technical Specification TS 22.261^[6] lists requirements for network slicing in two categories: *general requirements* and *management requirements*.

The general requirements for network slicing are the following:

- It must provide connectivity to home and roaming users in the same network slice.
- In a shared 5G network configuration, each operator must be able to apply all the requirements to their allocated network resources.
- It must support the *IP Multimedia Subsystem* (IMS) as part of a network slice.
- IMS support must be independent of network slices.

The IMS is a standards-based architectural framework for delivering multimedia communications services such as voice, video, and text messaging over IP networks^[7,12]. 3GPP originally developed the IMS specifications in the early 2000s to standardize access to multimedia services using cellular networks. The specifications define a complete framework and architecture that enables the convergence of video, voice, data, and mobile network technologies.

The management requirements of network slicing follow; it must:

- Allow the operator to create, modify, and delete a network slice.
- Allow the operator to define and update the set of services and capabilities supported in a network slice.
- Allow the operator to configure the information that associates a UE to a network slice.
- Allow the operator to configure the information which associates a service to a network slice.
- Allow the operator to assign a UE to a network slice, to move a UE from one network slice to another, and to remove a UE from a network slice based on subscription, UE capabilities, the access technology the UE uses, and the operator's policies and services the network slice provides.

- Support a mechanism for the *Visited Public Land Mobile Network* (VPLMN), as authorized by the *Home Public Land Mobile Network* (HPLMN), to assign a UE to a network slice with the needed services or to a default network slice.
- Enable a UE to be simultaneously assigned to and access services from more than one network slice of one operator.
- Ensure traffic and services in one network slice will have no impact on traffic and services in other network slices in the same network.
- Ensure that the creation, modification, and deletion of a network slice will have no or minimal impact on traffic and services in other network slices in the same network.
- Support scaling of a network slice, that is, adaptation of its capacity.
- Enable the network operator to define a minimum available capacity for a network slice. Ensure that scaling of other network slices on the same network will have no impact on the availability of the minimum capacity for that network slice.
- Enable the network operator to define a maximum capacity for a network slice.
- Enable the network operator to define a priority order between different network slices in case multiple network slices compete for resources on the same network.
- Support means by which the operator can differentiate policy control, functionality, and performance provided in different network slices.

Identification and Selection of a Network Slice

The *Single Network Slice Selection Assistance Information* (S-NSSAI) defines a single network slice. An S-NSSAI consists of two elements:

- *Slice/Service Type* (SST): An identifier that refers to the expected slice behavior in terms of features and services. Standardized SST values provide a way for establishing global interoperability for slicing so that 5G networks can support the roaming use case more efficiently for the most commonly used SSTs. Table 1 lists the standardized SSTs.
- *Slice Differentiator* (SD): Optional information that complements the SST to differentiate among multiple network slices of the same SST.

Table 1: Standardized Slice/Service Type Values.

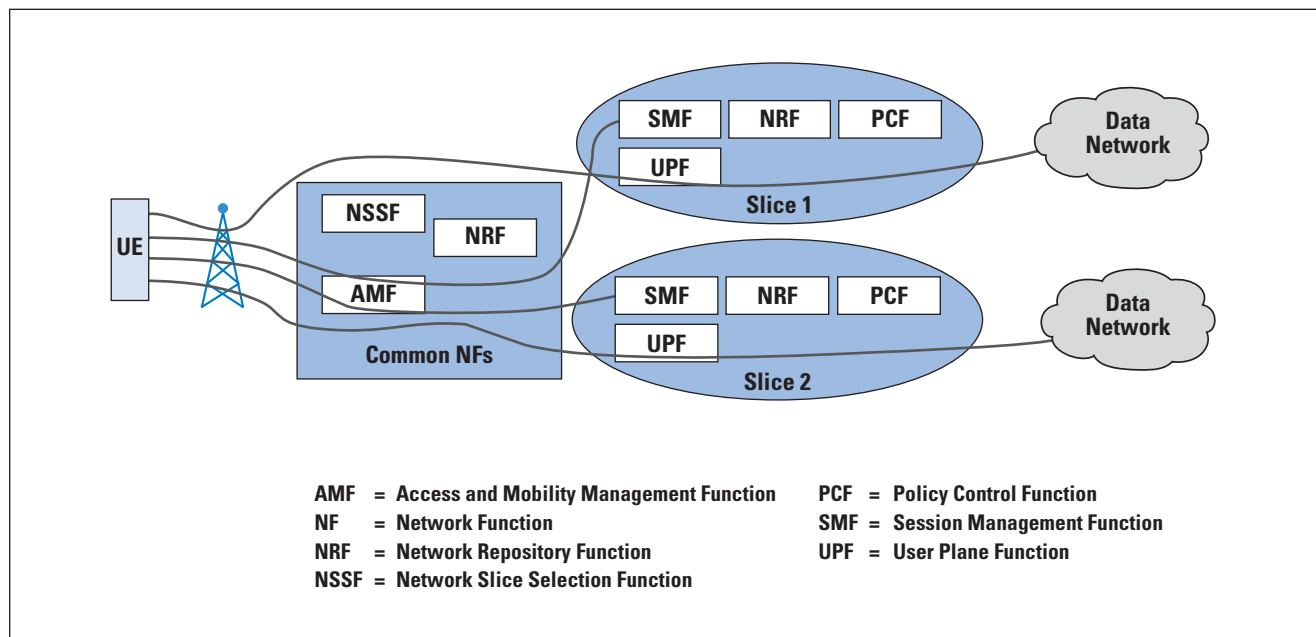
Slice/Service Type	SST Value	Characteristics
eMBB	1	Slice suitable to handle 5G-enhanced Mobile Broadband.
URLLC	2	Slice suitable to handle ultra-reliable low-latency communications.
Massive IoT (mIoT)	3	Slice suitable to handle massive IoT.
Vehicle-to-Everything (V2X)	4	Slice suitable to handle V2X services.

A UE may be served by up to eight network slices at a time, each identified by an S-NSSAI. The set of S-NSSAIs associated with a UE form a *Network Slice Selection Assistance Information* (NSSAI) data object.

Functional Aspects of Network Slicing

Figure 3 indicates the manner in which core NFs are used to implement network slices. A network function is a processing function in a network that has defined functional behavior and interfaces. You can implement a network function as a network element on dedicated hardware, as a software instance running on dedicated hardware, or as a virtualized function instantiated on an appropriate platform, for example, on a cloud infrastructure. Some NF instances support multiple network slices serving a UE, while others are specific to a given slice.

Figure 3: Network Functions that Support Network Slicing



The common NFs follow:

- *Access and Mobility Management Function* (AMF): Network slice instance selection is usually triggered as part of the registration procedure by the first AMF that receives the registration request from the UE. When a UE accesses the network, AMF provides functionalities to register and de-register the UE with the network, and it establishes the user context in the network. In the registration procedure, AMF performs, but is not limited to, network slice instance selection, UE authentication, authorization of network access and network services, and network access policy control. In addition, when AMF receives a session establishment request message from UE, it performs discovery and selection of the SMF that is the most appropriate to manage the session.

- *Network Slice Selection Function* (NSSF): The AMF retrieves the slices that the user subscription allows and interacts with the NSSF to select the appropriate network slice instance (for example, based on allowed S-NSSAIs, 5G network ID, and other parameters). The NSSF responds with a message including the list of appropriate network slice instances for the UE. As a result, the registration process may switch to another AMF if needed.
- *Network Repository Function* (NRF): During the AMF-NSSF interaction, the NSSF may return the identity of one or more NRFs to be used to select NFs and services within the selected network slice instance(s).

The slice-specific NFs follows:

- *Session Management Function* (SMF): The UE sends a message to the AMF requesting that a *Protocol Data Unit* (PDU) session be associated to one S-NSSAI and one *Data Network* (DN). The AMF selects the appropriate SMF, which manages the PDU session. The SMF sets up the PDU session for the UE and controls the user-plane operation. The SMF selects the UPF and invokes enforcement of QoS and charging policies.
- *User Plane Function* (UPF): Once a PDU session is established, QoS flows for this PDU session over this network slice pass through the UPF.
- *Policy Control Function* (PCF): The SMF gets policy information related to session establishment from the PCF.
- *Network Repository Function* (NRF): The SMF uses the NRF to discover the required NFs for the individual network slice.

Generic Slice Template

3GPP TS 28.531^[8] includes a description of the concept of the *Generic Slice Template* (GST). The *GSM Association* (GSMA) has specified the GST, which provides a standardized list of attributes that you can use to characterize different types of network slices^[9]. A *Network Slice Type* (NEST) is a GST filled with (ranges of) values. There may be two kinds of NESTs:

- *Standardized NEST* (S-NEST): Attributes are assigned (ranges of) values by *Standards-Developing Organizations* (SDOs), working groups, forums, and so forth, such as 3GPP, GSMA, *5G Automotive Association* (5GAA), and the *5G Alliance for Connected Industry and Automation* (5G-ACIA).
- *Private NEST* (P-NEST): Attributes are assigned (ranges of) values by the Network Slice Providers; these values are different from those assigned in S-NESTs.

Network Slice Providers can build their network slice product offering based on S-NESTs and/or their P-NESTs. GSMA has developed the GST to be a list of attributes sufficient for describing a wide range of NESTs that you can fully construct by allocating values (or ranges of values) to each relevant attribute in the GST. A network operator can use a NEST to identify the network resources and functions needed to instantiate network slices. The process to fill in the GST and to create a NEST comprises three steps:

1. Study use cases and derive service requirements based on discussions with the slice customers, such as vertical industries or specific enterprises.
2. Convert the service requirements identified in (1) into technical requirements.
3. Document the technical requirements produced in (2) using the NEST by filling in the values of each of the attributes of the GST.

The current version of the GST lists 35 attributes, shown in Figure 4.

Figure 4: Generic Network Slice Template Attributes

Availability	Network Functions Owned by Network Slice Customer	Supported Device Velocity
Area of Service	Maximum Number of PDU Sessions	Synchronicity
Delay Tolerance	Maximum Number of UEs	UE Density
Deterministic Communication	Performance Monitoring	Uplink Throughput per Network Slice
Downlink Throughput per Network Slice	Performance Prediction	Uplink Maximum Throughput per UE
Downlink Maximum Throughput per UE	Positioning Support	User Management Openness
Energy Efficiency	Radio Spectrum	User Data Access
Group Communication Support	Root Cause Investigation	V2X Communication Mode
Isolation Level	Session and Service Continuity Support	Latency from (last) User Plane Function (UPF) to Application Server
Maximum Supported Packet Size	Simultaneous Use of the Network Slice	Network Slice Specific Authentication and Authorization (NSSAA) Required
Mission-Critical Support	Slice Quality of Service Parameters	
Multimedia Telephony (MMTel) Support	Support for Non-IP Traffic	
NB-IoT Support		

Summary

Virtualization encompasses a variety of technologies for managing computing resources by providing a software translation layer, known as an abstraction layer, between the software and the physical hardware. Virtualization turns physical resources into logical, or virtual, resources. Virtualization enables users, applications, and management software operating above the abstraction layer to manage and use resources without needing to be aware of the physical details of the underlying resources. NFV is a key technology for implementing 5G wireless networks.

References

- [0] William Stallings, “Network Functions Virtualization,” *The Internet Protocol Journal*, Volume 24, No. 2, July 2021.
- [1] William Stallings, *5G Wireless: A Comprehensive Introduction*, ISBN-13: 9780136767299, Pearson, 2021.
- [2] Xin Li, Mohammed Samaka, H. Anthony Chan, Deval Bhamare, Lav Gupta, Chengcheng Guo, and Raj Jain, “Network Slicing for 5G: Challenges and Opportunities,” *IEEE Internet Computing*, September/October 2017.
- [3] ITU-T, “Framework for the support of network slicing in the IMT-2020 network,” ITU-T Y.3112, December 2018
- [4] ITU-R, “IMT Vision—Framework and overall objectives of the future development of IMT for 2020 and beyond,” ITU-R Recommendation M.2083, September 2015.
- [5] Next Generation Mobile Network Alliance, “5G End-to-End Architecture Framework,” Version 4.31, November 2020.
- [6] 3GPP TS 22.261, “Technical Specification Group Services and System Aspects; Service requirements for the 5G system; Stage 1 (Release 18),” January 2021.
- [7] Martin Koukal and Robert Bestak, “Architecture of IP Multimedia Subsystem,” Proceedings ELMAR Symposium, June 2006.
- [8] 3GPP TS 22.531, “Technical Specification Group Services and System Aspects; Management and Orchestration; Provisioning; (Release 16),” April 2020.
- [9] GSM Association, “Generic Network Slice Template Version 3.0,” May 22, 2020.
- [10] William Stallings, “Introduction to 5G Part One: Standards, Specifications, and Usage Scenarios,” *The Internet Protocol Journal*, Volume 26, No. 2, September 2023.
- [11] William Stallings, “Introduction to 5G Part Two: Core Network, Radio Access Network, and Air Interface,” *The Internet Protocol Journal*, Volume 26, No. 3, December 2023.
- [12] Mark Grayson, “Lessons Learned from 20 Years of Cellular and Wi-Fi Integration,” *The Internet Protocol Journal*, Volume 26, No. 3, December 2023.

WILLIAM STALLINGS is an independent consultant and author of numerous books on computer networking, security, and computer architecture. His latest book is *Wireless 5G: A Comprehensive Introduction* (Pearson, 2021). He maintains a computer science resource site for computer science students and professionals at **ComputerScienceStudent.com** and is on the editorial board of *Cryptologia*. He has a Ph.D. in computer science from M.I.T. and can be reached at **wllmst@icloud.com**

The History and Future of Ethernet

by Mikael Holmberg, *Extreme Networks*

Initially developed by *Xerox Palo Alto Research Center* (PARC) in the 1970s and ratified by the *Institute of Electrical and Electronics Engineers* (IEEE) as a standard in 1983, the evolution of Ethernet has taken this technology through many specifications and standardizations during its 50-year history.

Ethernet technology has become the backbone of modern communication and connectivity, linking billions of devices to each other and the Internet. Today, Ethernet connects *Local Area Networks* (LANs), *Wide Area Networks* (WANs), Internet, Cloud, *Internet of Things* (IoT) devices, Wi-Fi, and many other systems into one seamless global communications network.

The name *Ethernet* is based on the word “ether” as a way of describing an essential feature of the system: the physical medium (that is, a cable) carries bits to all stations, much the same way that the old “luminiferous ether” was once thought to propagate electromagnetic waves through space.

In its early days, Ethernet competed with other technologies like *Token Ring*. It was eventually chosen as the ubiquitous technology used in computer networks because of the simplicity by which the communication protocol can be deployed and its ability to incorporate modern advancements without losing backward compatibility. Ethernet continues to reign as the de facto standard for computer networking and many newly evolved applications and use cases. Just to choose one of interest, the topic that everybody talks about today is *Artificial Intelligence* (AI). As AI workloads increase, network industry giants are teaming to ensure Ethernet networks can keep pace and satisfy the AI high-performance networking requirements, among many other new use cases and applications. I will cover AI as well as a few other interesting use cases around the evolved Ethernet in this article.

In 1975, Xerox filed a patent application listing Bob Metcalfe, David Boggs, Chuck Tucker, and Butler Lampson as inventors. Then, in 1976, after the system was deployed at PARC, Metcalfe and Boggs published a seminal paper.^[1] Four gentlemen, Yogen Dalal, Ron Crane, Bob Garner, and Roy Ogus, facilitated the upgrade from the original 2.94-Mbps to the 10-Mbps protocol, which was released to the market in 1980 and ratified by the IEEE as a standard in 1983^[2]. Ethernet has become the dominant LAN technology, and five decades after its initial specification its evolution continues.^[31]

Taking a step back in time, let’s look at the progress of Ethernet technology over the past five decades and explore where experts think it could be heading in the years to come.

The Early Days of Ethernet

The evolution of Ethernet officially began in 1973 when engineer Robert Metcalfe introduced the concept in a memo he wrote while working at Xerox PARC. Metcalfe initially described Ethernet as interconnecting computing workstations, and it enabled them to communicate with each other as well as with devices like laser printers. These interconnected endpoints became the environment we now recognize as the world's first LAN.

Metcalfe was inspired by ALOHAnet^[1], an earlier networking project that began at the University of Hawaii in 1968 and aimed to connect remote workstations across the Hawaiian Islands to a central computer at the main Oahu campus.

ALOHAnet was realized by using a quite rudimentary *Additive Links On-line Hawaii Area* (ALOHA) protocol, where an end station would transmit a frame over a common data channel and then wait for confirmation that it had reached its destination successfully. If the end station didn't receive confirmation within a given period, it assumed a *collision* had occurred with another frame sent by a different end station simultaneously. In that case, that station would continue to resend the data until it achieved successful transmission. But as amounts of end stations and transmissions increased, more collisions would occur, and the network would become less efficient. An ALOHA variation named *Slotted ALOHA* aimed to minimize network contention problems by precisely coordinating individual transmissions for the end stations and assigning them designated timeslots via a beacon signal schema.

Metcalfe's Ethernet experiment, at that time referred to as the *Alto ALOHA Network*, included many revolutionary features that enabled significantly more efficient use of a computer network. This set of rules, which became known as the *Carrier Sense Multiple Access/Collision Detect* (CSMA/CD) protocol, allowed end stations to monitor the availability of a shared communication path and detect possible collisions when two end stations sent data at the same time. When frames collided, the system would discard them, leading each end station to wait for a randomly assigned length of time before trying to resend. The end station would continue this schema to pause and try to resend as many times as necessary. This process is known as *exponential back-off*.^[33]

So, the original Ethernet technology was based on a shared medium that was collision-prone, where all computers trying to communicate shared the same cable and, as such, competed with each other. The modern Ethernet implementation has a collision-free switched connection, where each computer communicates with only its own switch port, without competing for the cable with others.

By 1973, Metcalfe thought that the technology had outgrown its original name and renamed it *Ethernet*. Four years later, Metcalfe and Boggs, together with Charles Thacker and Butler Lampson, who also worked at Xerox, successfully patented Ethernet technology.^[3]

How Does Ethernet Work?

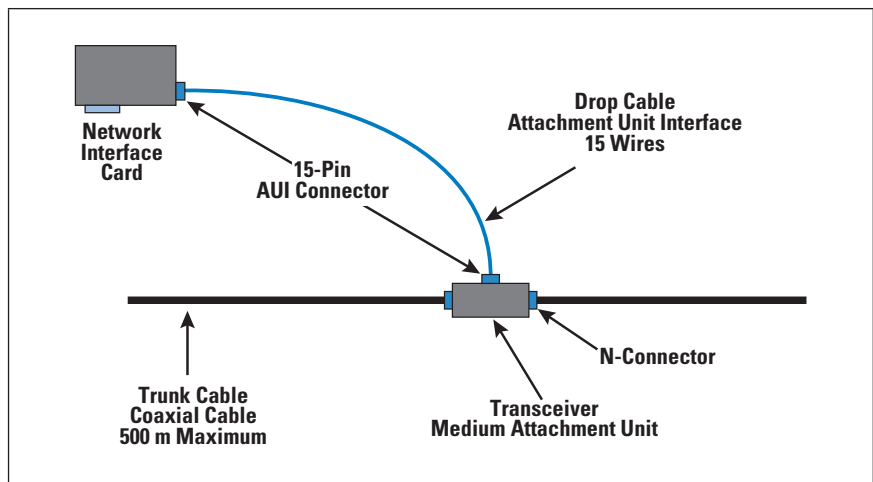
Ethernet works by breaking up data being sent to or from devices, like a personal computer, into short pieces of different-sized bits of information called *frames*. Those frames contain information such as the source and destination address that helps the frame route its way through a network.

In the past, computers on a LAN typically shared a single connection. Ethernet was built around the principle of CSMA/CD, as was briefly explained earlier in this article, where the protocol ensures that the cable is not in use before sending any frames out. Now that feature is far less important than it was in the early days of networking, as devices generally have their own private connection to a network through a *switch*. Ethernet now operates using *Full Duplex* (FDX), where the sending and receiving channels are separate, so it is impossible for collisions to occur over the same connection. As there is no error correction in Ethernet, the communication relies on upper-layer advanced protocols to ensure that everything is transmitted perfectly. Ethernet provides the basis for most digital communications, and it integrates quite easily with most higher-level protocols.

Ethernet IEEE Standardization

Xerox worked with two other vendors, Digital Equipment Corporation and Intel, to publish the first 10-Mbps Ethernet specification in 1980. Meanwhile, the *Local and Metropolitan Area Networks* (LAN/MAN) Standards Committee at the IEEE set out to develop a similar open standard. The IEEE LAN/MAN committee, which applies the number 802 to all its standards, formed an Ethernet subcommittee and named it the *IEEE 802.3 Working Group*. Through the first half of the 1980s, the Ethernet 10BASE-5 implementation used a coaxial cable 0.375 inches (9.5 mm) in diameter, also referred to as *Thick Ethernet* or *Thicknet*, and it was standardized in 1982 as 10BASE-5.

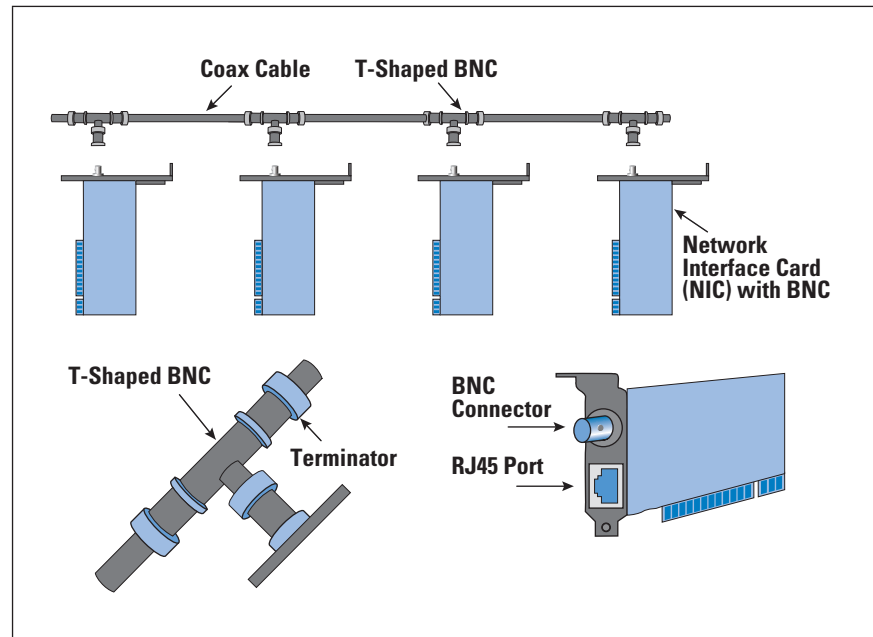
Figure 1: 10BASE-5 Ethernet



Some of us who have worked in this industry for many years might recall that on 10BASE-5, you drilled the transceiver pin (the so-called “vampire tap”) into the core of the thick coaxial cable, and if you were not careful, you might end up shortcutting the wire.

In the late 1980s, 10BASE-5 was replaced by 10BASE-2, called *Thin Ethernet* or *Thinnet*, and it used a BNC connector to connect the Ethernet *Network Interface Cards* (NICs) to a BNC-T splitter to ensure that the Ethernet segment stayed intact. It used the RG-58 coaxial cable, which is 0.2 inches (5 mm) in diameter, as media. The emphasis was on making the installation of the cable easier and less costly. This thin Ethernet was followed by twisted pair (10BASE-T) and fiber-optic (10BASE-FL).

Figure 2: 10BASE-2 Ethernet



In 1995, with the Fast Ethernet standard, the speed was upgraded to 100 Mbps, and no such speed improvement was ever made for Thinnet. By 2001, prices for Fast Ethernet cards had fallen to under \$50, and by 2003 Wi-Fi (802.11) networking equipment was widely available and affordable. Because of the immense demand for high-speed networking, the low cost of *Category 5* (Cat5) cable, and the popularity of 802.11 wireless networks, both 10BASE-2 and 10BASE-5 have become obsolete, though devices using those standards might still exist in some locations.

Also, in 1995, 100-Mbps Ethernet introduced *auto-negotiation*, which allowed for two network devices to signal each other and establish the best-shared mode of operation, including speed and duplex mode.

Three years later, a new milestone was reached when the 802.3 working group introduced *Gigabit Ethernet* [GE] (100BASE-T)^[4, 5], which was first realized over fiber-optic cable and, subsequently, over twisted-pair copper cable.

The evolution of Ethernet continued with 10-Gbps speeds in 2002, initially over fiber, then over twisted-pair copper cable, and finally, over unshielded twisted-pair cable. Then, in 2010, IEEE approved 40 GE and 100 GE, which was realized by aggregating multiple 10-Gbps lanes.

In 2016, driven by the rising demand from hyperscalers (web companies), the IEEE ratified 25 GE, which was 2.5 times faster than 10 GE but more cost-efficient than 40 GE. This standard improved throughput by increasing the capacity of a single lane, rather than aggregating multiple lower-capacity lanes, and meant that 25 GE required less cable and power and had higher port density than 40 GE. In some cases, an upgrade to 25 GE lets data-center operators extend the life of top-of-rack switches and avoid full “rip-and-replace” upgrades of cabling infrastructure. Hence, hyperscalers upgraded to 25 GE speeds in their data centers.

The following year, in late 2017, the networking industry saw the ratification of 200 and 400 GE. These standards were both based on 50-Gbps single lanes, as the cloud providers and hyperscale data centers, *Internet Service Providers* (ISPs), and specialized organizations like *Network Operations Centers* (NOCs) needed and wanted more bandwidth. Some of the challenges with 400-Gbps speeds include new cabling requirements because the current Category 5 and 5e cables don’t support such speeds.

In 2019, *Communication Service Providers* (CSPs) began deploying (or, more likely, testing) 5G^[26, 27] networks, the fifth-generation technology standard for broadband cellular networks. It is defined by the *3rd Generation Partnership Project* (3GPP), and it is the planned successor to 4G networks. Like its predecessors, 5G networks are cellular networks. All 5G wireless devices in a cell connect to the Internet and telephone network via radio waves through a local antenna in the cell. These new networks boost higher download speeds, eventually up to 10 Gbps. In addition to being faster than existing networks, 5G offers higher bandwidth, enabling it to connect a greater number of devices and improve the quality of Internet services in crowded areas. Naturally, Ethernet acts as the packet-based solution within 5G, accommodating all the essential containerized microservices required for 5G functions running on computers in all sizes of data centers with Ethernet fabric technologies.

Ethernet-based 5G cloud data-center fabrics come in various sizes, from small edge data-center fabrics implemented as Layer 2 network infrastructures to truly scalable three- and five-stage large data-center fabrics. These larger fabrics deployed as Layer 3 infrastructure with dozens or even hundreds of Ethernet switches connected in a *spine and leaf* architecture, also known as CLOS. The CLOS architecture has its origins in Charles Clos’ crossbar switches for telephone-call switching, and it is composed of leaf and spine layers where switches are used.

The most prevalent design for these cloud data-center fabrics consists of Ethernet switches that use *Virtual Extensible LAN* (VXLAN) with *Multiprotocol Extensions for BGP* (MP-BGP) and an *Ethernet Virtual Private Network* (EVPN) control plane.

All Ethernet switches are deployed in pairs to provide dual-homed redundant connectivity to computers and other switches. The leaf switch pairs interconnect to form a cluster, providing redundancy for the attached computers. *Border Leaf* (BL) switches are also deployed in pairs, ensuring dual-homed redundant connectivity to external *Provider Edge* (PE) routers and the Internet. This connectivity presents yet another interesting use case in which Ethernet serves as the foundation for cloud-native 5G mobile network applications and workloads.

Technically, the specification for 800-Gbps Ethernet also exists but is not really used outside of test environments. The interesting thing about Ethernet is that because it is such an open protocol, there is no reason to think that even the 800-Gbps speeds are anywhere near the theoretical maximum. Research is being done to set the groundwork for a 1.6-Tbps standard. Speeds like that will probably be useful only in highly specific applications.

Ethernet Cables

You can't talk about Ethernet without also talking about various cables used for Ethernet. As I previously described, the early days of Ethernet relied on coaxial cables, basically the same as were used for cable television. Coaxial cable is robust in design, having a thick internal copper wire, but it does have trade-offs, as it is heavy and difficult to work with and not very flexible. Ethernet changed to use twisted-pair cables that are still used when deploying Ethernet networks today, as are fiber-optic cables.

Many companies that manufacture Ethernet cables have moved away from the dull gray color scheme and instead offer them in a wide variety of colors that allow for improving data-center racks with different-colored cables. It also enables color-coding, so technicians can group their different network connections visually into groups based on different colors for quick troubleshooting.

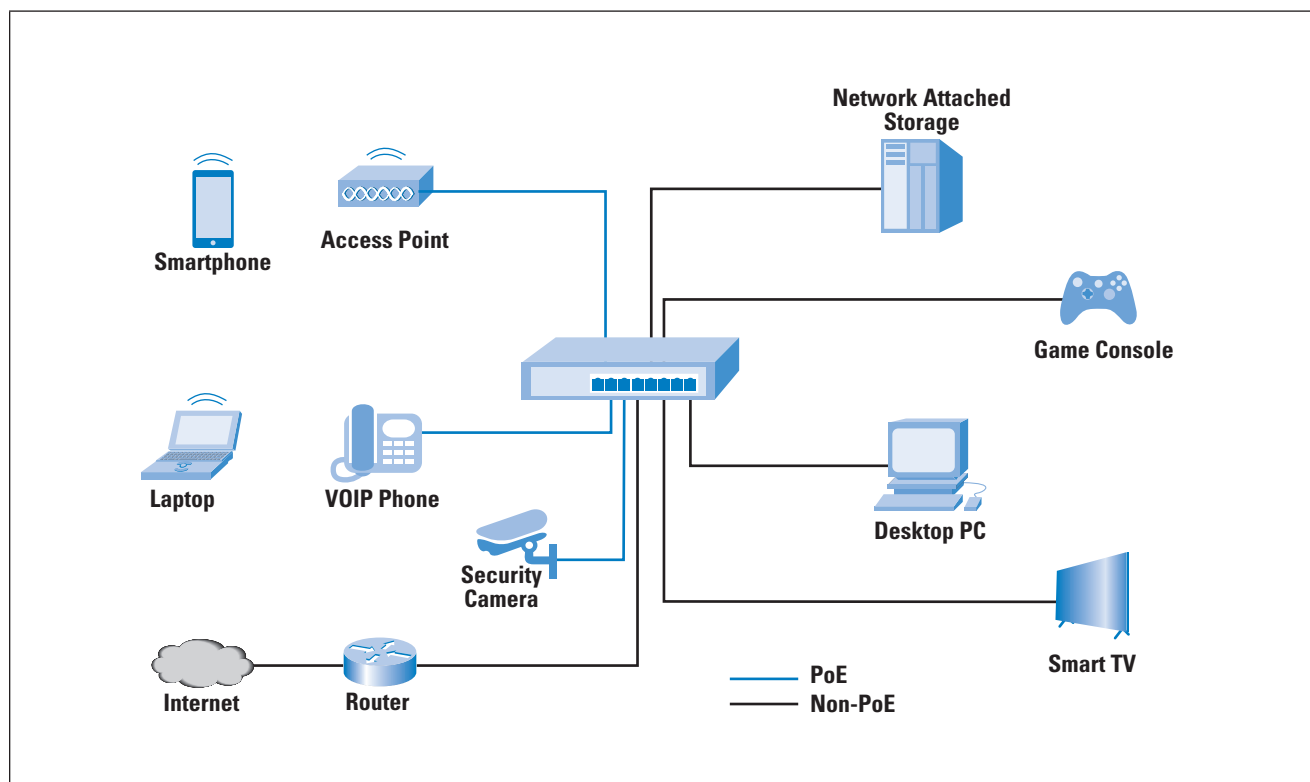
The standard plug on both ends of twisted-pair cables (RJ-45), which is very similar to the same kind of connector that wired telephone systems use (RJ-11), made it easy to just click the cables into any device that supports Ethernet connectivity. So, simply plugging in a device and attaching it to a network using one of those colored Ethernet cables is the only step required to gain connectivity. The long-time standard for Ethernet cables is *Category 5* (Cat5). The Cat5 standard has been used since 2001, and a slightly more advanced cable called *Category 5e* (Cat5e) is also used today for faster Ethernet applications. Category 5e cables are targeted at 100-Mbps Ethernet, but the design also supports higher speeds, such as Gigabit Ethernet. New *Category 6* (Cat6) cables have been introduced to support higher speeds than the Cat5 and Cat5e cables.

Several standards have been defined for *Power over Ethernet* (PoE), which allows you to connect devices with a single Ethernet cable without the need for additional power sources. Table 1 lists Ethernet cabling standards.

Table 1: Ethernet Cabling Standards

ETHERNET TYPE	BANDWIDTH	CABLE TYPE	MAXIMUM DISTANCE
10BASE-T	10Mbps	Cat 3/Cat 5 UTP	100m
100BASE-TX	100Mbps	Cat 5 UTP	100m
100BASE-TX	200Mbps	Cat 5 UTP	100m
100BASE-FX	100Mbps	Multi-mode fiber	400m
100BASE-FX	200Mbps	Multi-mode fiber	2Km
1000BASE-T	1Gbps	Cat 5e UTP	100m
1000BASE-TX	1Gbps	Cat 6 UTP	100m
1000BASE-SX	1Gbps	Multi-mode fiber	550m
1000BASE-LX	1Gbps	Single-mode fiber	2Km
10GBASE-T	10Gbps	Cat 6a/Cat 7 UTP	100m
10GBASE-LX	10Gbps	Multi-mode fiber	100m
10GBASE-LX	10Gbps	Single-mode fiber	10Km

Figure 3: Ethernet Switching, Wi-Fi, and PoE



Ethernet and Time Synchronization – Are We in Sync?

As applications continue to advance, latency has become a significant concern that we need to address. The solution to this problem lies in using *Precision Time Protocol* (PTP)^[28, 29], as we are addressing timing accuracy in the range of hundreds of nanoseconds. The use of Ethernet in mission-critical networks and the telecom industry showcases that Ethernet has now emerged as the de facto transport technology, using protocols like PTP to synchronize clocks throughout the network.

To appreciate the significance of PTP, it's important to understand that we're addressing timing accuracy in a range of hundreds of nanoseconds. This problem represents extremely tight timing requirements for certain applications and use cases. Maintaining precise timing is crucial for operating distributed systems at scale while ensuring that various operations remain synchronous. Additionally, precise timing is especially crucial when handling critical processes that govern infrastructure operations.

An excellent example of such systems can be seen in the telecom industry's 5G networks. Coordinating time between multiple servers or base stations in 5G could be compared to synchronized swimming in the Olympics, where all swimmers must perform their part of the routine at the same pace. If they perform at different paces, the routine will not look as it should. Ensuring that all servers and base stations operate in sync is vital for efficient network performance, much like how synchronized swimming relies on all swimmers to perform their parts of the routine at the same pace.

In addition to managing latency, new use cases for Ethernet also require it to become a deterministic networking technology, and as I briefly discussed 5G and PTP previously, *Time-Sensitive Networking* (TSN) comes into play, where PTP is one of the requirements for Ethernet to become a deterministic networking technology.

Ethernet and TSN for Deterministic Networking

Ethernet has evolved to incorporate various applications and technologies. TSN enables the synchronization of network elements and endpoints, such as switches and routers, to prioritize traffic classes and provide accountable delay and guaranteed bandwidth reservation. TSN is based on numerous international standards that are integrated with the Ethernet standard IEEE 802.3, where punctuality is ensured, allowing for transmission within a given period while simultaneously accommodating a mix of other communication protocols.

When discussing use cases and Ethernet, I need to mention *Industry 4.0* needs as the fourth industrial revolution that is transforming the manufacturing industry towards more efficient, connected, and flexible factories of the future. With Industry 4.0, factories will be able to rely on cloud-native technologies and connectivity based on Ethernet and TSN.

The goal of Industry 4.0 is to create full transparency across all processes and assets at all times, including both the *Information Technology* (IT) and *Operational Technology* (OT) domains based on architectures that require communication between production systems, logistics chains, people, and processes to unite these two domains into a single domain.

Ethernet has been used as the wired solution in both computer and automation networks. Ethernet open standard allows you to connect end devices quickly and simply as well as easily scale them to exchange data with other devices and functions. Ethernet was not originally designed to meet the requirements set by automation technology regarding guaranteed and real-time communication. Various bus systems in automation have evolved using Ethernet on a physical level while implementing proprietary real-time protocols such as PROFIBUS, PROFINET, and EtherCAT, to name a few. These systems often lead to the exclusive use of the network infrastructure as well as vendor dependencies. Today such networks handling time-critical data traffic are separated from networks directing less-critical data traffic. In the future, Industry 4.0 applications will require increasingly more consistent Ethernet networks, and TSN will provide a solution for that.

Traditional Ethernet networks involving sectors like manufacturing are based on a hierarchical automation model that separates information technology from operational technology. The IT domain includes enterprise-like communication with typical end devices such as computers, while the OT domain includes systems, machines, and software used for process control and automation. These two areas are fundamentally different in how they communicate, where the IT domain requires bandwidth while the OT domain requires high availability. Data traffic in the IT domain can be classified more as non-critical, while data traffic in the OT domain is critical and time-sensitive, and as such each domain uses a particular communication standard. Ethernet, as we know it in the enterprise or IT domain, relies on TCP/IP, while the OT domain relies on various bus systems, also known as *fieldbus systems*.

In the IT domain today, wireless technologies like Wi-Fi are used and could be used in some parts of the OT domain, but because of the nature of the technology, which is based on unlicensed spectra, it cannot guarantee bounded low latency with high reliability as the load increases. In certain use cases, Wi-Fi does not perform that well during uncontrolled interference because it uses an unlicensed spectrum. That may not be relevant for less-critical applications, because there will be a variety of applications with different traffic profile demands in both the Ethernet-enabled IT and OT domains within Industry 4.0.

Future of Ethernet

Our new world of AI workloads is expected to put unprecedented performance and capacity demands on networks that are based on Ethernet. Hence, we are possibly looking at a new enhancement to the well-known Ethernet technology to handle the scale and speed required by AI.

A group of vendors and operators have teamed to form the *Ultra Ethernet Consortium* (UEC)^[30], as there are concerns that today's traditional network interconnects cannot provide the required performance, scale, and bandwidth to keep up with AI demands. The consortium aims to address these concerns by adding new capabilities to the known and proven Ethernet technology specification, adding numerous new features and capabilities including:

- Multipathing and packet spraying to ensure AI workflows have access to a destination simultaneously.
- Flexible delivery order to make sure Ethernet links are optimally load-balanced while ordering is enforced only when the AI workload requires it in bandwidth-intensive operations.
- Congestion control mechanisms to ensure AI workloads avoid hotspots and spread the load evenly across multipaths within the network. These mechanisms can be designed to work in conjunction with multipath packet spraying, thus enabling a reliable transport of AI traffic.
- End-to-end telemetry to manage congestion, where information originating from the network can advise the participants of the location and cause of the congestion. In addition, shortening the congestion signalling path and providing more information to the endpoints allow more responsive congestion control.

After this journey covering some highlights after Ethernet has enjoyed five decades of existence, one might contend that Ethernet is one of the most crucial technologies today, even though it often goes unnoticed. Ethernet, as the ubiquitous network technology, powers infrastructure across the cosmos as it is used in space as well as in the deepest ocean trenches.

I named just a few examples, including the new era of cloud-native 5G data centers that provide the infrastructure for 5G applications and workloads, the industry revolution (Industry 4.0), as well as the AI challenges. But as we all know, the number of applications being developed that have substantial requirements not only around bandwidth but also latency, etc. is increasing.

This demand requires that the underlying Ethernet transport technologies can cater to such requirements; consequently, 400 GE is a reality today, and 800 GE is expected to become commonplace in the near future. Given this trend, it wouldn't be surprising to see 1-Terabyte Ethernet in use by 2030.

Based on our unending appetite for bandwidth, Ethernet, a 50-year-old technology, will soon reinvent itself once more.

References and Further Reading

- [0] Robert Metcalfe and David Boggs, “Ethernet: Distributed Packet Switching for Local Computer Networks,” *Communications of the ACM*, Volume 19, Issue 7, July 1976.
- [1] Norman Abramson, “The ALOHA System – Another Alternative for Computer Communications,” *Proceedings of the 1970 Fall Joint Computer Conference, American Federation of Information Processing Societies (AFIPS)*, Houston, Texas, November 17–19, 1970.
- [2] John Shoch, Yogen K. Dalal, David D. Redell, and Ronald C. Crane, “Evolution of the Ethernet Local Computer Network,” *Computer*, August 1982.
- [3] Robert M. Metcalfe, David R. Boggs, Charles P. Thacker, and Butler W. Lampson, “Multipoint Data Communication System with Collision Detection,” United States Patent US4063220A, December 13, 1977.
- [4] William Stallings, “Gigabit Ethernet,” *The Internet Protocol Journal*, Volume 2, No. 3, September 1999.
- [5] William Stallings, “Gigabit Ethernet: From 1 to 100 Gbps and Beyond,” *The Internet Protocol Journal*, Volume 18, No. 1, March 2015.
- [6] Howard Frazier and Howard Johnson, “Gigabit Ethernet: From 100 to 1,000 Mbps,” *IEEE Internet Computing*, January/February 1999.
- [7] Serag Gadelrab, “10-Gigabit Ethernet Connectivity for Computer Servers,” *IEEE Micro*, May-June 2007.
- [8] Shamus McGillicuddy, “40 Gigabit Ethernet: The Migration Begins,” *Network Evolution E-Zine*, December 2012.
- [9] Gautam Chanda and Yinglin (Frank) Yang, “40 GbE: What, Why & Its Market Potential,” *Ethernet Alliance White Paper*, November 2010.
- [10] Mark Nowell, Vijay Vusirikala, and Robert Hays, “Overview of Requirements and Applications for 40 Gigabit and 100 Gigabit Ethernet,” *Ethernet Alliance White Paper*, August 2007.
- [11] John D’Ambrosia, David Law, and Mark Nowell, “40 Gigabit Ethernet and 100 Gigabit Ethernet Technology Overview,” *Ethernet Alliance White Paper*, November 2008.
- [12] Hidehiro Toyoda, Goichi Ono, and Shinji Nishimura, “100GbE PHY and MAC Layer Implementation,” *IEEE Communications Magazine*, Volume 48, Issue 3, March 2010.
- [13] Rick Rabinovich, “40 Gb/s and 100 Gb/s Ethernet Short Reach Optical and Copper Host Board Channel Design,” *IEEE Communications Magazine*, Volume 50, Issue 4, April 2012.

- [14] Timothy Prickett Morgan, “IEEE Gets Behind 25G Ethernet Effort,” *Enterprise Tech*, July 27, 2014.
- [15] Rick Merritt, “50G Ethernet Debate Brewing,” *EE Times*, September 3, 2014.
- [16] Tom Nolle, “Will We Ever Need 400 Gigabit Ethernet Enterprise Networks?” *Network Evolution E-Zine*, December 2012.
- [17] John D’Ambrosia, Paul Mooney, and Mark Nowell, “400 Gb/s Ethernet: Why Now?” *Ethernet Alliance White Paper*, April 2013.
- [18] Stephen Hardy, “400 Gigabit Ethernet Task Force Ready to Get to Work,” *Lightwave*, March 28, 2014.
- [19] John D’Ambrosia, “400GbE and High Performance Computing,” *Scientific Computing Blog*, April 18, 2014.
- [20] Jim Duffy, “100-Gigabit Ethernet: Bridge to Terabit Ethernet,” *Network World*, April 20, 2009.
- [21] John D’Ambrosia, “TEF 2014: The Rate Debate,” *Ethernet Alliance Blog*, June 23, 2014.
- [22] Scott Kipp, “5 New Speeds – 2.5, 5, 25, 50 and 400 GbE,” *Ethernet Alliance Blog*, August 8, 2014.
- [23] William Stallings, “Gigabit Wi-Fi,” *The Internet Protocol Journal*, Volume 17, No. 1, September 2014.
- [24] David Chalupsky and Adam Healey, “Datacenter Ethernet: Know Your Options,” *Network Computing*, March 28, 2014.
- [25] Rich Seifert, *Gigabit Ethernet: Technology and Applications for High-Speed LANs*, ISBN 0-201-18553-9, Addison-Wesley, 1998. (Reviewed in *The Internet Protocol Journal*, Volume 1, Number 2, September 1998.)
- [26] William Stallings, “Introduction to 5G Part One: Standards, Specifications, and Usage Scenarios,” *The Internet Protocol Journal*, Volume 26, No. 2, September 2023.
- [27] William Stallings, “Introduction to 5G Part Two: Core Network, Radio Access Network, and Air Interface,” *The Internet Protocol Journal*, Volume 26, No. 3, December 2023.
- [28] IEEE Standards Association, “IEEE Standard for a Precision Clock Synchronization Protocol for Networked Measurement and Control Systems,” IEEE 1588-2008, July 24, 2008.
- [29] IEEE Standards Association, “IEEE Standard for a Precision Clock Synchronization Protocol for Networked Measurement and Control Systems,” IEEE 1588-2019, June 16, 2020.
- [30] Ultra Ethernet Consortium: <https://ultraethernet.org/>

- [31] Wikipedia article on the IEEE 802.3 working group. Lists all current standards, as well as standards under development:
https://en.wikipedia.org/wiki/IEEE_802.3
- [32] Charles Spurgeon and Joann Zimmerman, *Ethernet: The Definitive Guide: Designing and Managing Local Area Networks 2nd Edition*, ISBN-13 978-1449361846, O'Reilly Media, 2014.
- [33] Wikipedia article on Exponential Backoff:
https://en.wikipedia.org/wiki/Exponential_backoff

MIKAEL HOLMBERG is a Distinguished Engineer and member of the office of the CTO at Extreme Networks. He is an experienced professional in networking architectures and technologies including cloud who has worked in the industry for over 30 years and is active in a variety of industry committees.
E-mail: mikael@extremenetworks.com

Our Privacy Policy

The *General Data Protection Regulation* (GDPR) is a regulation for data protection and privacy for all individual citizens of the *European Union* (EU) and the *European Economic Area* (EEA). Its implementation in May 2018 led many organizations worldwide to post or update privacy statements regarding how they handle information collected in the course of business. Such statements tend to be long and include carefully crafted legal language. We realize that we may need to provide similar language on our website and in the printed edition, but until such a statement has been developed here is an explanation of how we use any information you have supplied relating to your subscription:

- The mailing list for *The Internet Protocol Journal* (IPJ) is entirely “opt in.” We never have and never will use mailing lists from other organizations for any purpose.
- You may unsubscribe at any time using our online subscription system or by contacting us via e-mail. We will honor any request to remove your name and contact information from our database.
- We will use your contact information only to communicate with you about your subscription; for example, to inform you that a new issue is available, that your subscription needs to be renewed, or that your printed copy has been returned to us as undeliverable by the postal authorities.
- We will never use your contact information for any other purpose or provide the subscription list to any third party other than for the purpose of distributing IPJ by post or by electronic means.
- If you make a donation in support of the journal, your name will be listed on our website and in print unless you tell us otherwise.

Letter to the Editor

Craig Partridge’s “Why ATM Failed” article in *The Internet Protocol Journal*, Volume 26, No 3, December 2023, is excellent. It prompted me to write down a few of my own recollections at the time:

I had thought that it was the work of Sandy Fraser in Bell Labs in 1974 that transformed the *Time Division Multiplexor* (TDM) theory with a label attached to the data segment that identified the now virtual timeslot of each data stream that was the precursor of ATM. My reading on this topic had noted that when this technology was presented to the Bell telephone operating companies as a scalable switching architecture that had far greater flexibility and efficiency than TDM-based architectures in the mid-1980s, the response was positive, and many expected the migration to ATM in the telephone switching fabric to be completed by 2020!

I also recall some material from Dave Sincowski and Bob Lyon from the late 1980s about the early days of ATM. At that time the experience of changing the local network architecture from common-bus 10-Mbps Ethernet to 100-Mbps *Fiber Distributed Data Interface* (FDDI) rings was challenging: It involved large-scale replacement of both the physical network media of the *Local-Area Network* (LAN) and the network interfaces in the attached workstations. The underlying concern was that the next incremental step in LAN speeds would require a similar comprehensive replacement. They were searching for a scalable LAN architecture that could offer increasing capacity without enforced replacement of large parts of the network infrastructure. Bell Lab’s Dave Sincowski proposed ATM to Sun Microsystem’s Bob Lyons as an answer to that problem, as Sun Microsystem’s workstation products had just gone through the Ethernet-to-FDDI NIC transformation. That proposal appears to be why it was the computing sector, not the telephone sector, that was behind the initial adoption of ATM in digital networking.

It wasn’t just the small buffer size in ATM switches that was an issue here for high-speed computer networks, it was the choice of the ATM cell size of 53 bytes that proved to be a problem. The cell size decision was already a compromise between a small cell size that more closely matched the TDM slot size that was ideally suited to low-jitter, low-volume voice data streams and a much larger cell size that reduced the per-cell processing overheads of a higher-speed data stream. Switching equipment initially struggled to achieve switching performance of a 10-Mbps Ethernet with a theoretical maximum packet rate of some 1,440 packets per second, so a smaller cell size would require faster processing capability in the switches.

I also recall at that time extensive debate about the “correct” internal buffer dimensioning for ATM switches. Low-jitter objectives call for very shallow internal buffers, while the congestion-loss algorithms of TCP called for delay-bandwidth-sized internal buffers.

I dimly recall a report on Doug Comer's experiments on achievable throughputs using a 155-Mbps ATM switch where the shallow buffers in ATM switches resulted in an achieved 3-Mbps data throughput from a 155-Mbps switch. The LAN market had already gained extensive experience with Ethernet switching, and ATM simply was not an effective alternative in price and performance for local networks.

Despite these issues, for a while in the early to mid 1990s ATM had some success. I vaguely recall the first 155-Mbps backbone circuits in the US used Fore ATM switches with ATM as the only available technology with a clock that was faster than 45 Mbps.

The telephone companies clung to the dream of a single multipurpose digital switched foundation for a suite of data products. The company I worked for at the time in Australia, Telstra, was using Nortel Magellan Passport ATM switches as the basis for their suite of data products as well as telephone trunks. But it was useful for only a brief period of time. When IP services wanted to use circuit transmission rates in excess of 622 Mbps, the Passport switches could not deliver, and at that point the IP product moved to sit beside the Passports and connect to the *Synchronous Digital Hierarchy* (SDH) network at first, and then directly to optical transponders shortly thereafter.

—Geoff Huston, gih@apnic.net

Check your Subscription Details!

Make sure that both your postal and e-mail addresses are up-to-date since these are the only methods by which we can contact you. If you see the words "Invalid E-mail" on your printed copy this means that we have been unable to contact you through the e-mail address on file. If this is the case, please contact us at ipj@protocoljournal.org with your new information. The subscription portal is located here: <https://www.ipjsubscription.org/>

IAB Statement on Encryption and Mandatory Client-side Scanning of Content

The *Internet Architecture Board* (IAB) recently issued the following statement: “A secure, resilient, and interoperable Internet benefits the public interest and supports human rights^[1] to privacy and freedom of opinion and expression. This is endangered by technologies, such as recent proposals for client-side scanning, that mandate unrestricted access to private content and therefore undermine end-to-end encryption and bear the risk to become a widespread facilitator of surveillance and censorship.

This statement is a reaction to recent policy proposals in the United Kingdom^[2], European Union^[3], United States^[4], and other countries that are mandating client-side scans that require access to otherwise end-to-end encrypted content. These proposals envision client-side scanning technologies that search content on devices before it is encrypted or after decryption on receipt. This would potentially be accomplished by comparison against a database maintained by an authority or by leveraging machine learning to identify previously unseen but potentially prohibited content. These envisioned mechanisms fail to consider their broader implications for Internet security.

The *Internet Engineering Task Force* (IETF) is the leading standards development organization for the global Internet. The IAB provides long-range technical direction for Internet development, ensuring the Internet continues to grow and evolve as a platform for global communication and innovation. To create and maintain the Internet as the bedrock of current secure communication, the IETF and the IAB serve as stewards of the Internet’s communication protocols and its core values of trust, openness, and fairness that underpin secure online communication. This is accomplished through a transparent process backed by consensus that is open for anybody to participate in. We encourage the continued deployment and strengthening of mechanisms that enhance privacy and security for all users of the Internet.

The IETF and the IAB have published concerns about standardizing wiretaps^[5], backdoors^[6,7], and surveillance^[8], because these technologies reduce the security of the Internet as a whole, fail to curtail malicious actors, and reduce security for Internet users. To ensure all communication can remain properly protected, the IETF continues to develop and enhance encrypted protocols like *Internet Protocol Security* (IPsec)^[9] at the IP layer, *Transport Layer Security* (TLS)^[10] at the transport layer which is further incorporated into the *Hypertext Transfer Protocol Version 2* (HTTP/2)^[11] and QUIC^[12] protocols, and inside many application protocols such as email *Secure/Multipurpose Internet Mail Extensions* (S/MIME)^[13], *Open Specification for Pretty Good Privacy* (OpenPGP)^[14] or instant messaging *Messaging Layer Security* (MLS)^[15] and *End-to-End Signing and Object Encryption for the Extensible Messaging and Presence Protocol* (XMPP)^[16]. Recognizing that management of increasingly encrypted networks can pose operational challenges, the IAB has recently held a workshop on techniques for managing encrypted networks in ways that intend not to sacrifice security for the Internet’s end-users^[17].

The IAB has recognized surveillance of any form as a threat to Internet user privacy, where “surveillance is the observation or monitoring of an individual’s communications or activities”^[18]. As the IAB and *Internet Engineering Steering Group* (IESG) documented in 1996^[6], instituting governmental control into communication “provide[s] only a marginal or illusory benefit to law enforcement agencies” as any seemingly beneficial purpose can be equally used by malevolent actors or future authoritarian shifts in government administrations. The IETF community still holds true to these principles today.

For technologies where the intended purpose is scanning of user communication, there is by design no technical way to limit the scope and intent of scanning, nor curtail subsequent changes in scope or intent. Further, specifically when scanning for illegal content, transparency cannot be provided. Mandating such technologies impacts all users of the global Internet and creates a tool that is straightforward to abuse as a widespread facilitator of surveillance and censorship, presenting real-world dangers to the free flow of information and the security and privacy of people. Without privacy, users cannot benefit from the Internet’s virtue to connect people and support freedom of expression.

Additionally, one of the founding principles of the Internet has been its openness; the ability for any standards-compliant software to access the network of networks has been the catalyst for world-changing innovations over many decades. Mandatory use of client-side scanning, and the regulatory burden it would impose, would negatively impact this, restrict use of open-source software, and lead to a stagnant landscape where users lose choice.

The IAB shares concerns about societal harms through the distribution of illegal content and criminal action on the Internet and recognizes the need to protect Internet users from such threats. However, the IAB believes that mandating client-side scanning is in direct opposition to the safe, secure and open communication platform that the Internet provides today and undermines the core principles applied by the IAB and the IETF^[5, 6, 18] in order to secure the Internet through encryption. The IAB opposes technologies that foster surveillance as they weaken the user’s expectations of private communication which decreases the trust in the Internet as the core communication platform of today’s society. Mandatory client-side scanning creates a tool that is straightforward to abuse as a widespread facilitator of surveillance and censorship. Mandating on-device scanning of content will compromise privacy, weaken security, and imperil human rights to communication, freedom of expression and freedom of opinion.”

References

- [1] United Nations Human Rights Office of The High Commissioner, “A/HRC/29/32: Report on encryption, anonymity, and the human rights framework,” May 2015.
- [2] UK Parliament, “Online Safety Act 2023.”

- [3] European Commission, “Proposal for a Regulation of The European Parliament and of The Council laying down rules to prevent and combat child sexual abuse,” May, 2022.
- [4] United States Congress, S.1207, “Eliminating Abusive and Rampant Neglect of Interactive Technologies Act of 2023.”
- [5] IAB and IETF, “IETF Policy on Wiretapping,” RFC 2804, May 2000. *This document articulates why the IETF stated that it was not appropriate to accommodate wiretapping.*
- [6] IAB and IESG, “IAB and IESG Statement on Cryptographic Technology and the Internet,” RFC 1984, August 1996. *This document stated the IESG and IAB’s position regarding legal constraints on encryption in 1996, with a focus on the effects on the Internet. The publication of the document was prompted in part by the controversy surrounding the US government’s promotion of the Clipper Chip.*
- [7] Jeffrey I. Schiller, “Strong Security Requirements for Internet Engineering Task Force Standard Protocols,” RFC 3365, August 2002. *This document set a requirement for IETF standard protocols to use ‘appropriate strong security mechanisms,’ including encryption. It was published as Best Current Practice in 2002.*
- [8] Stephen Farrell and Hannes Tschofenig, “Pervasive Monitoring Is an Attack,” RFC 7258, May 2014. *This RFC documents the IETF consensus that pervasive monitoring is an attack, and thus should be mitigated in IETF protocols (often, using encryption). It was a response to the Snowden revelations and an output of the workshop on Strengthening the Internet Against Pervasive Monitoring (STRINT), held jointly by the W3C and IAB.*
- [9] IETF IPsec Working Group
- [10] Eric Rescorla, “The Transport Layer Security (TLS) Protocol Version 1.3,” RFC 8446, August 2018.
- [11] Mike Belshe, Roberto Peon, and Martin Thomson, “Hypertext Transfer Protocol Version 2 (HTTP/2),” RFC 7540, May 2015.
- [12] IETF QUIC Working Group
- [13] IETF S/MIME Mail Security Working Group
- [14] IETF Open Specification for Pretty Good Privacy Working Group
- [15] IETF Messaging Layer Security Working Group
- [16] Peter Saint-Andre, “End-to-End Signing and Object Encryption for the Extensible Messaging and Presence Protocol (XMPP),” RFC 3923, October 2004.
- [17] Mallory Knodel, Wes Hardaker, and Tommy Pauly, “Report from the IAB Workshop on Management Techniques in Encrypted Networks (M-TEN),” RFC 9490, January 2024.
- [18] Alissa Cooper, Hannes Tschofenig, Bernard Aboba, Jon Peterson, John B. Morris, Jr., Marit Hansen, and Rhys Smith, “Privacy Considerations for Internet Protocols,” RFC 6973, July 2013.

The full IAB statement, including additional references, can be found on the IETF datatracker website.

Achieving Greater Heights for MANRS

The *Global Cyber Alliance* (GCA) recently announced a new phase for *Mutually Agreed Norms for Routing Security* (MANRS)^[0,1]. The *Internet Society* has partnered with the GCA an international non-profit specializing in addressing cybersecurity challenges at scale by mobilizing stakeholders toward collective action. As part of this partnership, the GCA will take on the functions of the MANRS secretariat and operations, while the Internet Society will maintain significant funding, advocacy, and training functions over the next five years.

In 2014, the Internet Society recognized the industry's willingness for collaborative agreement on best practices for routing security and helped initial participants to capture and share those practices in what became MANRS. Since then, the Internet Society has advocated globally for MANRS uptake, encouraged industry collaboration, supported the evolution of the norms, and evolved to become the secretariat of MANRS.

Fast forward a decade, and MANRS has grown from nine original operators to a community of more than 1,000 participants ranging from small enterprise networks to Tier-1 transit providers, from *Internet Exchange Points* (IXPs) of various sizes to *Content Delivery Network* (CDN) and cloud providers publicly professing their commitment to the MANRS requirements. MANRS is now globally recognized as a beacon for securing global routing.

As MANRS matured, so did the community-led governance model with the establishment of the community-elected Steering Committee. The Internet Society has proudly served as the secretariat, in addition to supporting the initiative with both financial and staff resources as well as operations support to ensure MANRS' growth. In 2019 the *MANRS Observatory*, a conformance measurement tool for routing security, was launched. Since then, many new features have been added to the MANRS Observatory, such as alerts and monthly MANRS readiness reports. Growth also happened through capacity building, and over the years, thousands of network engineers have gone through online courses, virtual labs, and on-site workshops. In 2020, the Internet Society, together with the MANRS community, launched the *Mentors and Ambassadors* program promoting routing security in the areas of research, policy, and training.

MANRS has more than one thousand participating operators across three programs, as well as six network equipment vendors. The initiative has been a tremendous success, but the task of supporting MANRS has grown well beyond the scope of what was a startup initiative 10 years ago. This partnership is an important evolution of a successful initiative that the Internet Society launched, incubated, and nurtured. GCA is honored and excited to step into this role and provide the basis for the long-term sustainability and evolution of MANRS.

Routing security is one of the focus areas of GCA, and the Internet Society and GCA have been collaborating around MANRS since 2021 with excellent results. GCA conducted a survey of network operators to learn more about the state of routing security implementation, the level of concern within network operations and business decision-making, and plans for next steps.

The Internet Society is dedicated to improving routing security and ensuring the best future for MANRS. Over the next five years, the Internet Society will focus on funding and support through training and global advocacy activities, while GCA will provide the secretariat function and operate the MANRS Observatory. GCA is uniquely placed to lead the next evolution of MANRS as its focus on building communities to collectively drive action towards addressing cybersecurity challenges at scale allows it to step into this role and provide the best future home for the operational growth MANRS is experiencing.

GCA is committed to maintaining the vision of MANRS and continuing to expand its global impact. With this partnership, MANRS will continue to achieve greater heights and be further established as the globally recognized benchmark for global routing security.

The partnership builds on the strengths of both organizations—GCA’s global footprint of mobilizing communities towards collective action to deliver a secure, trustworthy Internet that enables social and economic progress for all, and the Internet Society’s training and advocacy expertise. Together, the Internet Society and GCA are committed to maintaining and expanding the vision of MANRS to continue to increase the awareness and uptake of MANRS principles and improve the Internet’s functional integrity.

Everyone who runs a network has a responsibility to ensure a globally robust and secure routing infrastructure. Your network’s safety depends on a routing infrastructure that stops bad actors and mitigates accidental misconfigurations that wreak havoc on the Internet. The more network operators work together, the fewer incidents there will be, and the less damage they can do.

For more information about this partnership visit:

<https://www.globalcyberalliance.org/achieving-greater-heights-manrs/>

References

[0] MANRS: <https://manrs.org/>

[1] Andrei Robachevsky, “Improving Routing Security,” *The Internet Protocol Journal*, Volume 22, No. 2, July 2019.

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Kjetil Aas	Lukasz Bromirski	Richard Dodsworth	John Gilbert	Brian Johnson
Fabrizio Accatino	Václav Brožík	Ernesto Doelling	Serge Van	Curtis Johnson
Michael Achola	Christophe Brun	Michael Dolan	Ginderachter	Richard Johnson
Martin Adkins	Gareth Bryan	Eugene Doroniuk	Greg Goddard	Jim Johnston
Melchior Aelmans	Ron Buchalski	Michael Dragone	Tiago Goncalves	Jonatan Jonasson
Christopher Affleck	Paul Buchanan	Joshua Dreier	Ron Goodheart	Daniel Jones
Scott Aitken	Stefan Buckmann	Lutz Drink	Octavio Alfageme	Gary Jones
Jacobus Akkerhuis	Caner Budakoglu	Aaron Dudek	Gorostiaga	Jerry Jones
Antonio Cuñat Alario	Darrell Budic	Dmitriy Dudko	Barry Greene	Michael Jones
William Allaire	BugWorks	Andrew Dul	Jeffrey Greene	Amar Joshi
Nicola Altan	Scott Burleigh	Joan Marc Riera	Richard Gregor	Javier Juan
Shane Amante	Chad Burnham	Duocastella	Martijn Groenleer	David Jump
Marcelo do Amaral	Randy Bush	Pedro Duque	Geert Jan de Groot	Anders Marius Jørgensen
Matteo D'Ambrosio	Colin Butcher	Holger Durer	Ólafur Guðmundsson	Merike Kao
Selva Anandavel	Jon Harald Bøvre	Karlheinz Dölger	Christopher Guemez	Andrew Kaiser
Jens Andersson	Olivier Cahagne	Mark Eanes	Gulf Coast Shots	Vladislav Kalinovskiy
Danish Ansari	Antoine Camerlo	Andrew Edwards	Sheryll de Guzman	Naoki Kambe
Finn Arildsen	Tracy Camp	Peter Robert Egli	Rex Hale	Akbar Kara
Tim Armstrong	Brian Candler	George Ehlers	Jason Hall	Christos Karayiannis
Richard Artes	Fabio Caneparo	Peter Eisses	James Hamilton	Daniel Karrenberg
Michael Aschwanden	Roberto Canonico	Torbjörn Eklöv	Darow Han	David Kekar
David Atkins	David Cardwell	Y Ertur	Handy Networks LLC	Stuart Kendrick
Jac Backus	Richard Carrara	ERNW GmbH	Stephen Hanna	Robert Kent
Jaime Badua	John Cavanaugh	ESdatCo	Martin Hannigan	Thomas Kernen
Bent Bagger	Lj Cemerax	Steve Esquivel	John Hardin	Jithin Kesavan
Eric Baker	Dave Chapman	Jay Etchings	David Harper	Jubal Kessler
Fred Baker	Stefanos Charchalakakis	Mikhail Evstiounin	Edward Hauser	Shan Ali Khan
Santosh Balagopalan	Molly Cheam	Bill Fenner	David Hauweele	Nabeel Khatri
William Baltas	Greg Chisholm	Paul Ferguson	Marilyn Hay	Dae Young Kim
David Bandinelli	David Chosrova	Ricardo Ferreira	Headcrafts SRLS	William W. H. Kimandu
A C Barber	Marcin Cieslak	Kent Fichtner	Hidde van der Heide	John King
Benjamin Barkin-Wilkins	Lauris Cikovskis	Ulrich N Fierz	Johan Helsingius	Russell Kirk
Feras Batainah	Brad Clark	Armin Fisslthaler	Robert Hinden	Gary Klesk
Michael Bazarewsky	Narelle Clark	Michael Fiumano	Michael Hippert	Anthony Klopp
David Belson	Horst Clausen	The Flirble Organisation	Damien Holloway	Henry Kluge
Richard Bennett	James Cliver	Jean-Pierre Forcioli	Alain Van Hoof	Michael Kluk
Matthew Best	Guido Coenders	Gary Ford	Edward Hotard	Andrew Koch
Hidde Beumer	Robert Collet	Susan Forney	Bill Huber	Ia Kochiashvili
Pier Paolo Biagi	Joseph Connolly	Christopher Forsyth	Hagen Hultsch	Carsten Koempe
Arturo Bianchi	Steve Corbató	Andrew Fox	Kauto Huopio	Richard Koene
John Bigrow	Brian Courtney	Craig Fox	Asbjørn Højmark	Alexander Kogan
Orvar Ari Bjarnason	Beth and Steve Crocker	Fausto Franceschini	Kevin Iddles	Matthijs Koot
Tyson Blanchard	Dave Crocker	Erik Fredriksson	Mika Ilvesmaki	Antonin Kral
Axel Boeger	Kevin Croes	Valerie Fronczak	Karsten Iwen	Robert Krejčí
Keith Bogart	John Curran	Tomislav Futivic	Joseph Jackson	John Kristoff
Mirko Bonadei	André Danthine	Laurence Gagliani	David Jaffe	Terje Krogdahl
Roberto Bonalumi	Morgan Davis	Edward Gallagher	Ashford Jaggernaut	Bobby Krupczak
Lolke Boonstra	Jeff Day	Andrew Gallo	Thomas Jalkanen	Murray Kucherawy
Julie Bottorff	Fernando Saldana Del	Chris Gamboni	Jozef Janitor	Warren Kumari
Photography	Castillo	Xosé Bravo Garcia	Martijn Jansen	George Kuo
Gerry Boudreaux	Rodolfo Delgado-Bueno	Oswaldo Gazzaniga	John Jarvis	Dirk Kurfuerst
Leen de Braal	Julien Dhallenne	Kevin Gee	Dennis Jennings	Mathias Körber
Kevin Breit	Freek Dijkstra	Rodney Gehrke	Edward Jennings	Darrell Lack
Thomas Bridge	Geert Van Dijk	Radu Cristian Gheorghiu	Aart Jochem	Andrew Lamb
Ilia Bromberg	David Dillow	Greg Giessow	Nils Johansson	Richard Lamb

Yan Landriault	Kevin Menezes	Chris Perkins	Carsten Scherb	Douglas Thompson
Edwin Lang	Bart Jan Menkveld	Michael Petry	Ernest Schirmer	Kerry Thompson
Sig Lange	Sean Mentzer	Alexander Peuchert	Benson Schliesser	Lorin J Thompson
Markus Langenmair	Eduard Metz	David Phelan	Philip Schneck	Fabrizio Tivano
Fred Langham	William Mills	Harald Pilz	James Schneider	Peter Tomsu Fine Art
Tracy LaQuey Parker	David Millsom	Derrell Piper	Peter Schoo	Photography
Christian de Larrinaga	Desiree Miloshevic	Rob Pirnie	Dan Schrenk	Joseph Toste
Alex Latzko	Joost van der Minnen	Jorge Ivan Pincay	Richard Schultz	Rey Tucker
Jose Antonio Lazaro	Thomas Mino	Ponce	Timothy Schwab	Sandro Tumini
Lazaro	Rob Minshall	Marc Vives Piza	Roger Schwartz	Angelo Turetta
Antonio Leding	Wijnand Modderman-	Victoria Poncini	SeenThere	Michael Turzanski
Rick van Leeuwen	Lenstra	Blahoslav Popela	Scott Seifel	Phil Tweedie
Simon Leinen	Mohammad Moghaddas	Andrew Potter	Paul Selkirk	Steve Ulrich
Robert Lewis	Charles Monson	Ian Potts	Andre Serralheiro	Unitek Engineering AG
Christian Liberale	Andrea Montefusco	Eduard Llull Pou	Yury Shefer	John Urbanek
Martin Lillepuu	Fernando Montenegro	Tim Pozar	Yaron Sheffer	Martin Urwaleck
Roger Lindholm	Roberto Montoya	David Preston	Doron Shikmoni	Betsy Vanderpool
Link Light Networks	Joel Moore	David Raistrick	Tj Shumway	Surendran Vangadasalam
Art de Llanos	Joseph Moran	Priyan R Rajeevan	Jeffrey Sicuranza	Ramnath Vasudha
Mike Lochocki	John More	Balaji Rajendran	Thorsten Sideboard	Randy Veasley
Chris and Janet Lonvick	Maurizio Moroni	Paul Rathbone	Greipur Sigurdsson	Philip Venables
Mario Lopez	Brian Mort	William Rawlings	Fillipe Cajaiba da Silva	Buddy Venne
Sergio Loreti	Soenke Mumm	Mujtiba Raza Rizvi	Andrew Simmons	Alejandro Vennera
Eric Louie	Tariq Mustafa	Bill Reid	Pradeep Singh	Luca Ventura
Adam Loveless	Stuart Nadin	Petr Rejhon	Henry Sinnreich	Scott Vermillion
Josh Lowe	Michel Nakhla	Robert Remenyi	Geoff Sisson	Tom Vest
Guillermo a Loyola	Mazdak Rajabi Nasab	Rodrigo Ribeiro	John Sisson	Peter Villemoes
Hannes Lubich	Krishna Natarajan	Glenn Ricart	Helge Skrivervik	Vista Global Coaching
Dan Lynch	Naveen Nathan	Justin Richards	Terry Slattery	& Consulting
David MacDuffie	Darryl Newman	Rafael Riera	Darren Sleeth	Dario Vitali
Sanya Madan	Mai Nguyen	Mark Risinger	Richard Smit	Rüdiger Volk
Miroslav Madić	Thomas Nikolajsen	Fernando Robayo	Bob Smith	Jeffrey Wagner
Alexis Madriz	Paul Nikolich	Michael Roberts	Courtney Smith	Don Wahl
Carl Malamud	Travis Northrup	Gregory Robinson	Eric Smith	Michael L Wahrman
Jonathan Maldonado	Marijana Novakovic	Ron Rockrohr	Mark Smith	Lakhinder Walia
Michael Malik	David Oates	Carlos Rodrigues	Tim Sneddon	Laurence Walker
Tarmo Mammers	Ovidiu Obersterescu	Magnus Romedahl	Craig Snell	Randy Watts
Yogesh Mangar	Jim Oplotnik	Lex Van Roon	Job Snijders	Andrew Webster
John Mann	Tim O'Brien	Marshall Rose	Ronald Solano	Jd Wegner
Bill Manning	Mike O'Connor	Alessandra Rosi	Asit Som	Tim Weil
Harold March	Mike O'Dell	David Ross	Ignacio Soto Campos	Westmoreland
Vincent Marchand	John O'Neill	William Ross	Evandro Sousa	Engineering Inc.
Normando Marcolongo	Carl Önné	Boudhayan	Peter Spekrijse	Rick Wesson
Gabriel Marroquin	Packet Consulting	Roychowdhury	Thayumanavan Sridhar	Peter Whimp
David Martin	Limited	Carlos Rubio	Paul Stancik	Russ White
Jim Martin	Carlos Astor Araujo	Rainer Rudigier	Ralf Stempffer	Jurrien Wijlhuizen
Ruben Tripiana Martin	Palmeira	Timo Ruitert	Matthew Stenberg	Joseph Williams
Timothy Martin	Gordon Palmer	RustedMusic	Martin Štěpánek	Derick Winkworth
Carles Mateu	Alexis Panagopoulos	Babak Saberi	Adrian Stevens	Pindar Wong
Juan Jose Marin Martinez	Gaurav Panwar	George Sadowsky	Clinton Stevens	Makarand Yerawadekar
Ioan Maxim	Chris Parker	Scott Sandefur	John Streck	Phillip Yialeloglou
David Mazel	Alex Parkinson	Sachin Sapkal	Martin Streule	Janko Zavernik
Miles McCredie	Craig Partridge	Arturas Satkovskis	David Strom	Bernd Zeimet
Gavin McCullagh	Manuel Uruena Pascual	PS Saunders	Colin Strutt	Muhammad Ziad
Brian McCullough	Ricardo Patara	Richard Savoy	Viktor Sudakov	Ziauddin
Joe McEachern	Dipesh Patel	John Sayer	Edward-W. Suor	Tom Zingale
Alexander McKenzie	Dan Paynter	Phil Scarr	Vincent Surillo	Jose Zumalave
Jay McMaster	Leif Eric Pedersen	Gianpaolo Scassellati	Terence Charles Sweetser	Romeo Zwart
Mark Mc Nicholas	Rui Sao Pedro	Elizabeth Scheid	T2Group	廖明沂.
Olaf Mehlberg	Juan Pena	Jeroen Van Ingen	Roman Tarasov	
Carsten Melberg	Luis Javier Perez	Schenau	David Theese	

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at **ole@protocoljournal.org** or **olejacobsen@me.com**

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Diamond Sponsors

Your logo here!

Ruby Sponsors



Sapphire Sponsors

Your logo here!

Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

John Crain, Senior Vice President and Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Merike Kaeo, Founder and vCISO
Double Shot Security

Olaf Kolkman, Principal – Internet Technology, Policy, and Advocacy
The Internet Society

Dr. Jun Murai, Founder, WIDE Project
Distinguished Professor, Keio University
Co-Director, Keio University Cyber Civilization Research Center, Japan

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

July 2024

Volume 27, Number 2

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
Starlink and TCP	2
DNS Evolution	14
Fragments	23
Thank You.....	28
Call for Papers.....	30
Supporters and Sponsors	31

Internet access by means of *Low Earth Orbit* (LEO) satellites has become very popular in recent years, particularly in rural areas where alternative solutions are limited. We covered this technology in an article in our September 2023 issue (Volume 26, No. 2). The benefits of LEO systems include a much lower cost to launch and place the satellites into a low orbit, and a shorter *Round Trip Time* (RTT) as compared to solutions involving geosynchronous satellites. However, since LEO satellites move across the sky, a complex system of tracking and handoffs is deployed in order to provide continuous connectivity to the end user. In our first article, Geoff Huston examines the performance of Starlink from the point of view of the *Transmission Control Protocol* (TCP).

When I joined the *Network Information Center* (NIC) at SRI International in 1984, I was handed two *Request For Comments* (RFCs) describing the *Domain Name System* (DNS), and I was told that the DNS would soon be deployed across the Internet (mainly known as ARPANET and MILNET at the time). The NIC was still maintaining and publishing a host table in 1984, and it would take a couple of years before the DNS became fully operational. Our second article, also by Geoff Huston, looks at how the DNS has evolved in the last 40 years with various enhancements and extensions. The DNS is still one of the most active areas of work within the *Internet Engineering Task Force* (IETF).

As the Internet has evolved, interest by governments and intergovernmental organizations has grown to legislate and regulate various aspects of the system. These efforts, often collectively referred to as *Internet Governance*, are sometimes developed in ways that do not fully include input from the Internet technical community. One example is the *Global Digital Compact* (GDC) currently being drafted by the United Nations. In our Fragments section you will find a letter from individuals concerned about the latest GDC draft.

Publication of this journal is made possible by the generous support of our donors, supporters, and sponsors. We also depend on your feedback and suggestions. If you would like to comment on, donate to, or sponsor IPJ, please contact us at ipj@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

A View of Starlink from a Transport Protocol

by Geoff Huston, APNIC

Digital communications systems always represent a collection of design trade-offs. Maximising one characteristic of a system may impair other characteristics, and various communications services may offer different performance characteristics based on the intersection of these design decisions with the physical characteristics of the communications medium. In this article I'll look at the Starlink service^[0,1], and how the *Transmission Control Protocol* (TCP)—the transport-protocol workhorse of the Internet—interacts with the Starlink service.

To start, it's useful to recall a small piece of Newtonian physics from some 340 years ago^[2]. On the surface of the earth, assuming that you are high enough to clear various mountains that may be in the way—and also assuming that the earth has no friction-inducing atmosphere—if you fire a projectile horizontally fast enough it will not return to the earth, but head into space. There is, however, a critical velocity where the projectile will be captured by the earth's gravity and neither fall to ground nor head out into space. That orbital velocity at the surface of the earth is some 40,320 km/sec. The orbital velocity decreases with altitude, and at an altitude of 35,786 km above the surface of the earth the orbital velocity of the projective relative to a point on the surface of the spinning earth is 0 km/sec. This altitude is of a geosynchronous equatorial orbit, where the object appears to sit at a fixed location in the sky.

Geosynchronous Services

Geosynchronous satellites were the favoured approach for the first wave of satellite-based communications services. Each satellite could “cover” an entire hemisphere. If the satellite was on the equatorial plane, then it was at a fixed location in the sky with respect to the earth, allowing the use of large antennas. These antennas could operate at a low signal-to-noise ratio, allowing the signal modulation to use a high density of discrete phase amplitude points, which lifted the capacity of the service. All these advantages have to be offset against the less-favourable aspects of this service.

Consideration of crosstalk interference between adjacent satellites in geosynchronous orbits resulted in international agreements that require a 2° spacing for geosynchronous satellites that use the same frequency, so this orbital slot is a limited resource: it is limited to just 180 spacecraft if they all use K band (18–27 GHz) radio. At any point on the earth there is an upper bound to the signal capacity that can be received (and sent) using geosynchronous services.

It is relatively expensive to place satellites into this orbit because it generally requires three-stage rockets to propel them into this high orbit.

Depending on whether the observer is on the equator directly beneath the satellite or further away from this point, a geosynchronous orbit satellite is between 35,760 and 42,664 km away, so a signal *Round-Trip Time* (RTT) to the satellite and back will be between 238 and 284 ms in terms of signal propagation time. In IP terms, a RTT will be between 477 and 569 ms, and signal encoding and decoding times must be added to that. In addition, the delay for the signal to be passed between the endpoints and the satellite earth station must also be added. In practice, a RTT of around two-thirds of a second (660 ms) for Internet paths that use geosynchronous satellite services is common.

This extended latency means that the endpoints need to use large buffers to hold a copy of all the unacknowledged data, as is required by the TCP protocol. TCP is a feedback-governed protocol that uses ACK pacing. The longer the RTT the greater the lag in feedback, and the slower the response from endpoints to congestion or to available capacity. The congestion considerations lead to the common use of large buffers in the systems that drive the satellite circuits, which can further exacerbate congestion-induced instability. In geosynchronous service contexts, the individual TCP sessions are more prone to instability and they experience longer recovery times following low events^[3].

Low Earth Orbit Systems

A response to this situation is to bring the satellite closer to earth. This approach has several benefits. The spinning iron core of the earth generates a magnetic field, which traps energetic charged solar particles and redirects them through what is called the *Van Allen Belt*, thus deflecting solar radiation. Not only does this deflection allow the earth to retain its atmosphere, but it also protects the electronics of orbiting satellites that use an orbital altitude below 2,000 km or so from the worst effects of solar radiation. It's far cheaper to launch satellites into a *Low Earth Orbit* (LEO), and these days SpaceX can do so using reusable rocket boosters. The reduced distance between the earth and the orbiting satellite reduces the latency in sending a signal to the satellite and back, which can improve the efficiency of the end-to-end packet-transport protocols using such satellite circuits.

This group of orbital altitudes, from some 160 to 2,000 km, are collectively termed LEOs^[4]. The objective is to keep the orbit of the satellite high enough to prevent its slowing down by grazing the denser parts of the earth's ionosphere, but not so high that it loses the radiation protection afforded by the Inner Van Allen belt. At a height of 550 km, the minimum signal propagation delay to reach the satellite and return to the surface of the earth is just 3.7 ms.

But all of these facts come with some different issues. At a height of 550 km, an orbiting satellite can be seen from only a small part of the earth. If the minimum effective elevation to establish communication is 25 degrees of elevation above the horizon, then the footprint of the satellite is a circle with a radius of 940 km, or a circle of area 2M km².

To provide continuous service to any point on the surface of the earth (510.1M km^2), the number of orbiting satellites must be a minimum of 500. This reality implies that a satellite-based service is not a simple case of sending a signal to a fixed point in the sky and having that single satellite mirror that signal down to some outer earth location. A continuous LEO satellite service must use a continual sequence of satellites as they pass overhead and switch the circuit path across to successive satellites as they come into view.

At this altitude, the satellite orbits with a relative speed of 27,000 km/hour and it passes across the sky from horizon to horizon in less than 5 minutes. Some implications for the design of the radio component of the service are evident. The satellites are close enough that there is no need to use larger dish antennas that require some mechanised steering arrangement, but this situation itself is not without its downsides. An individual signal carrier might be initially received as a weak signal (in relative terms), increase in strength as the satellite transponder and the earth antenna move into alignment, and weaken again as the satellite moves on. Starlink's services use a phased-array arrangement with a grid of smaller antennas on a planar surface, which allows the antennas to be electronically steered by altering the phase difference between each of the antennas in the grid. Even so, this arrangement is relatively coarse, so the signal quality is not consistent, implying a constantly variable signal-to-noise ratio as the phased-array antenna tracks each satellite.

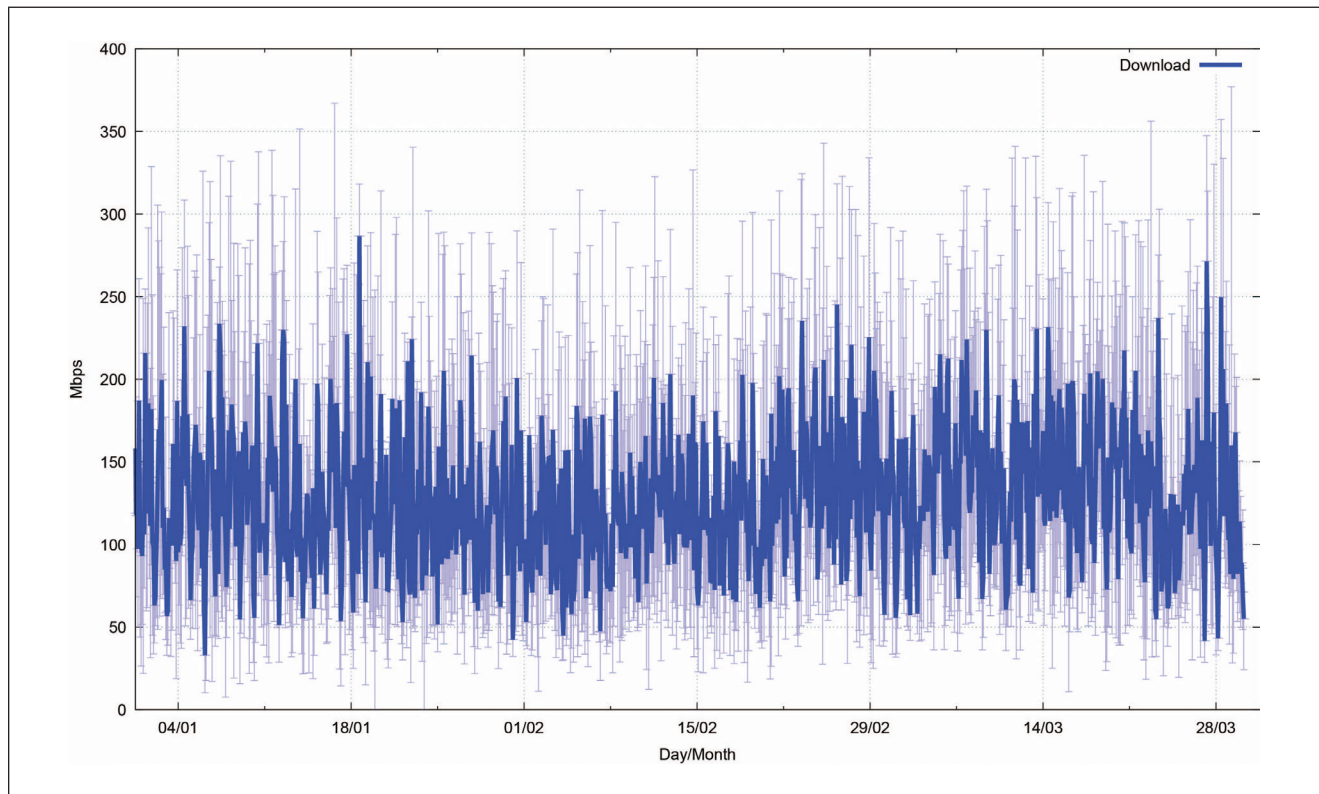
The modulation of this signal uses adaptive phase amplitude modulation, and as the signal-to-noise level improves, the modulator can use a larger number of discrete code points in this phase amplitude space, thus increasing the effective capacity of the service even while using a constant-frequency carrier signal. The implication is that if the satellite service attempts to always operate at peak efficiency, then it must constantly adapt its signal modulation to take advantage of the instantaneous signal-to-noise ratio, which results in a constantly varying service capacity.

Now we have four major contributory factors for variability of the capacity of the Starlink service, namely the variance in signal modulation capability, which is a direct outcome of the variable signal-to-noise ratio of the signal, the variance in the satellite path latency due to the relative motion of the satellite and the earth antennas, and the need to perform satellite switching constantly, and the variability induced by sharing the common medium with other users.

One way to see how this variability affects the service characteristics is to use a capacity measurement tool to measure the service capacity regularly. The results of such regularity of testing are shown in Figure 1. Here the test is a Speedtest measurement test^[5], performed on a 4-hourly basis for the period January 2024 through March 2024.

The service appears to have a median value of around 120 Mbps of download capacity, with individual measurements reading as high as 370 Mbps and as low as 10 Mbps, and 15 Mbps of upload capacity, with variance of between 5 and 50 Mbps.

Figure 1: Starlink Performance



In Internet terms, *ping*^[6] is a very old tool. However, at the same time it is very useful which probably explains its longevity. Figure 2 shows a plot of a continuous (flood) *ping* across a Starlink connection from the customer-side terminal to the first IP endpoint behind the Starlink earth station.

The first major characteristic of this data is that the minimum latency changes every 15 seconds. It appears that this change correlates to the user's being assigned to a different satellite, which implies that the user equipment "tracks" each spacecraft for 15-second intervals. This period corresponds to a tracking angle of 11 degrees of arc.

The second characteristic is that loss events are seen to occur at times of switchover between satellites (as shown in Figure 3), as well as occurring less frequently as a result of obstruction, signal quality, or congestion.

Figure 2: Starlink Ping Profile

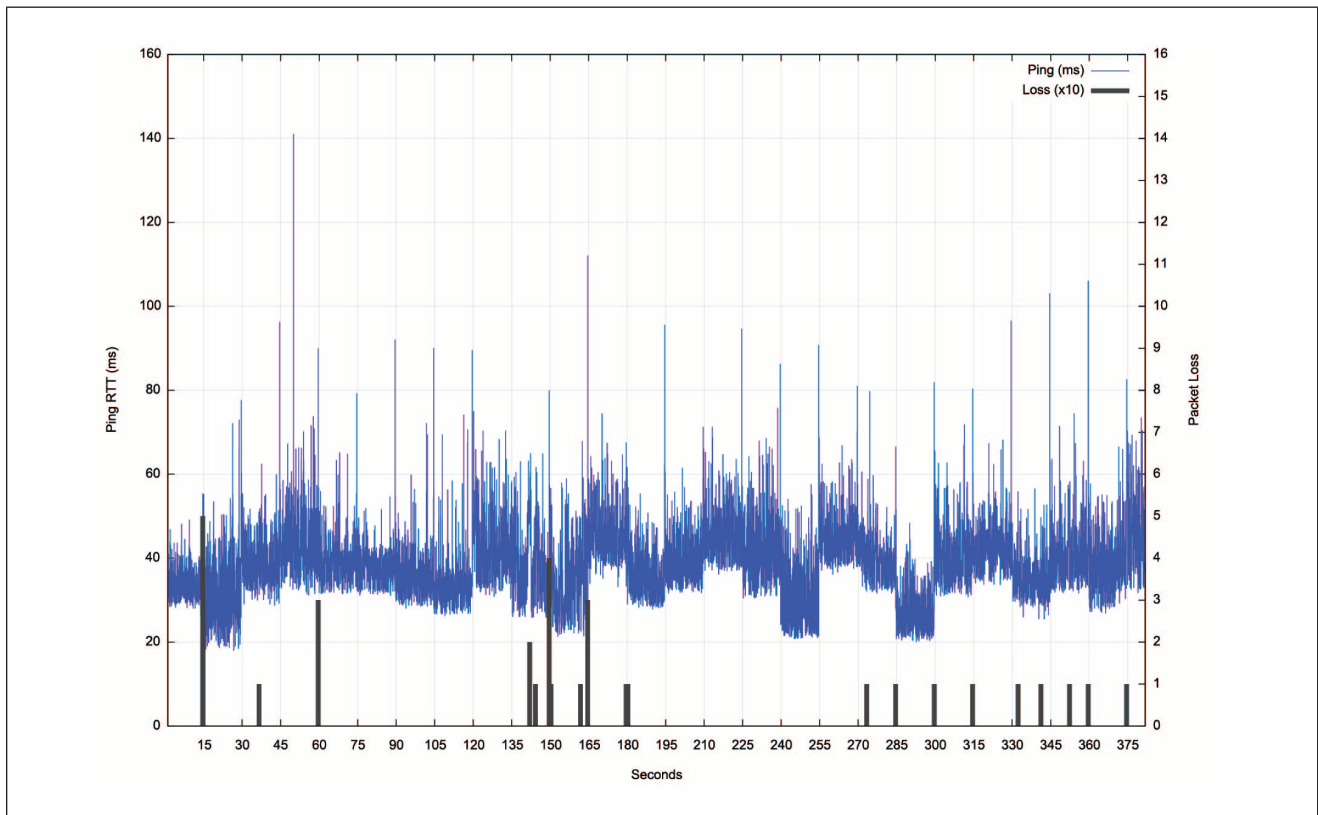
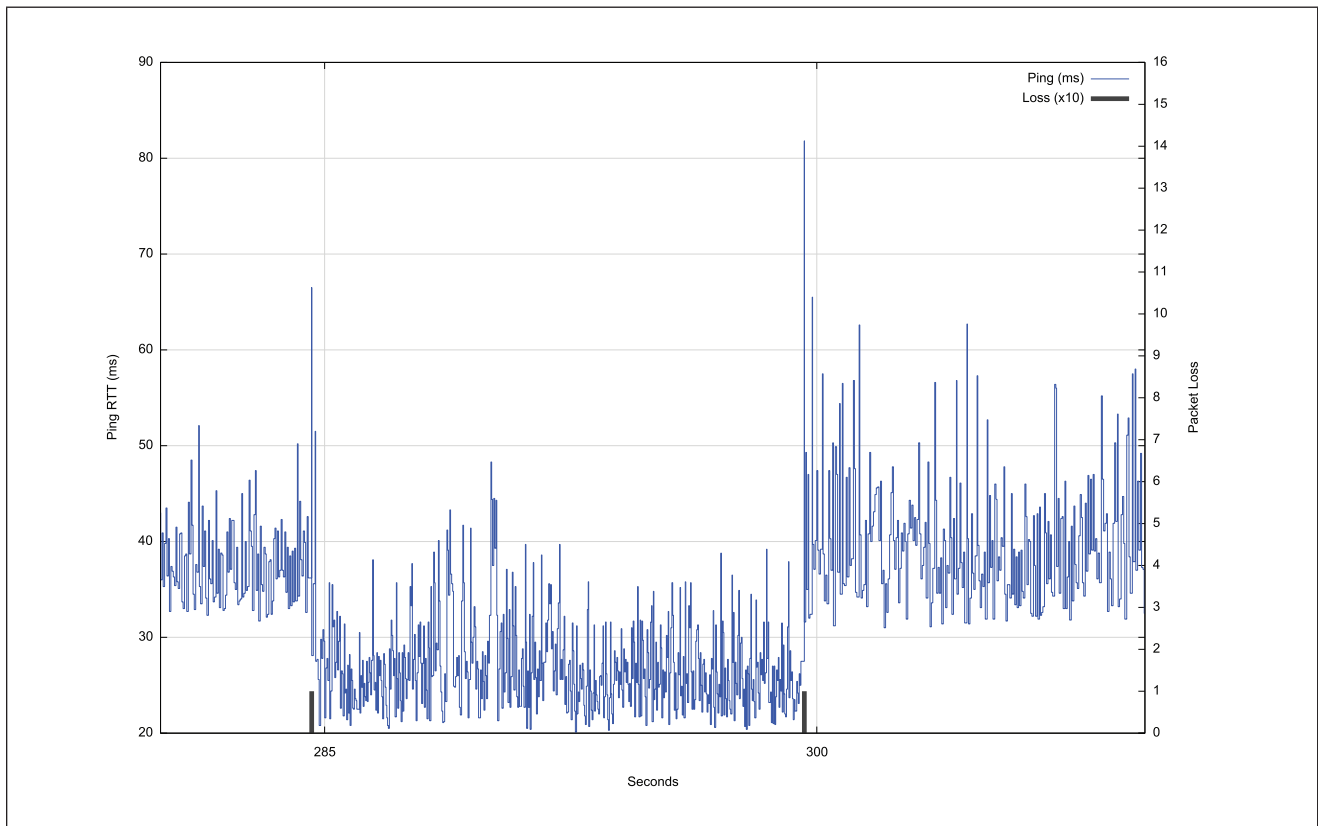


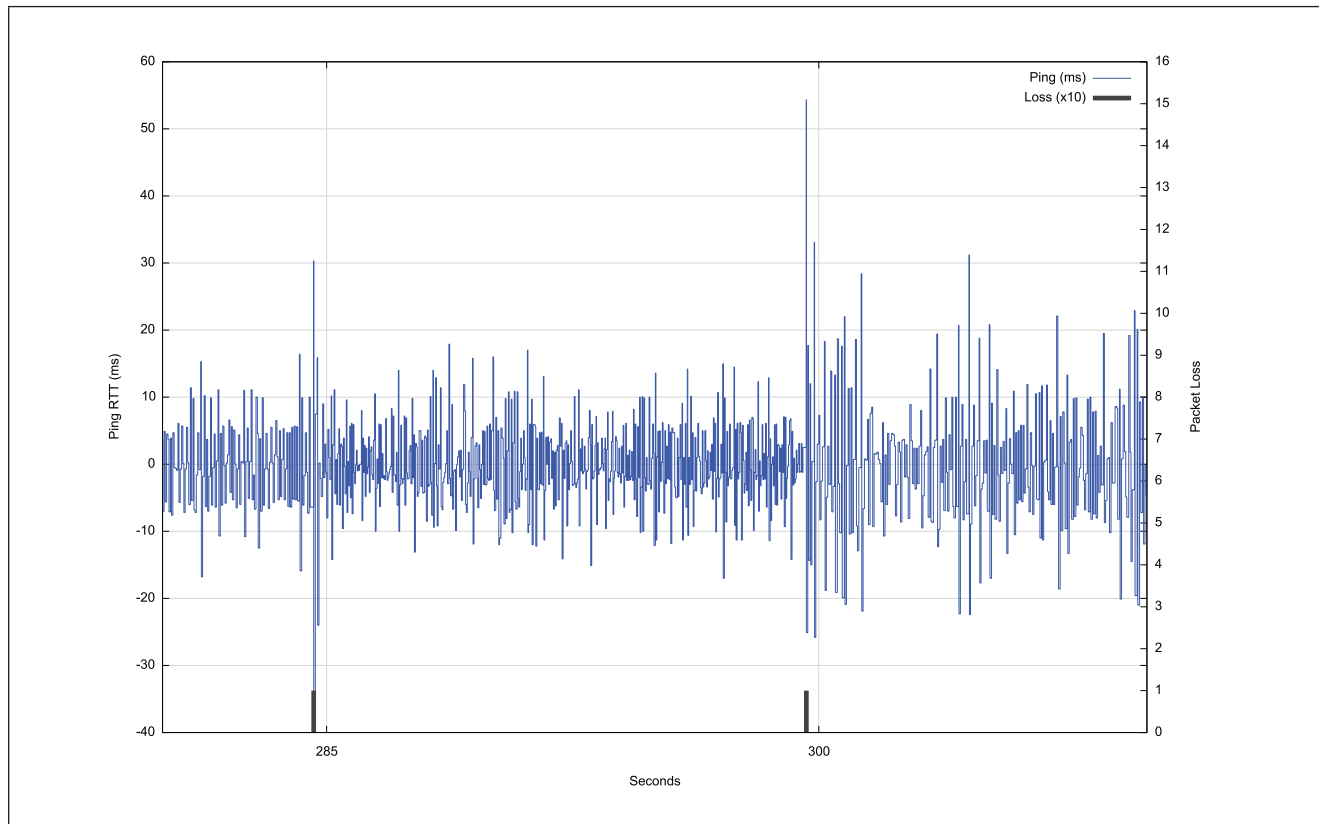
Figure 3: Starlink Ping Profile Showing Satellite Handover



The third characteristic is a major shift in latency when the user is assigned to a different spacecraft. The worst case in this data set is a shift from a minimum of 20 ms to a minimum of 40 ms.

Finally, within each satellite tracking interval the latency variation is relatively high. The average variation of jitter over successive RTT intervals is 6.7 ms. The latency spikes at handover impose an additional 30 to 50 ms, indicating the presence of deep buffers in the system to accommodate the transient issues associated with satellite handover (Figure 4).

Figure 4: Starlink Ping Profile Showing Latency Variance



The overall packet-loss rate when measured using 1-second paced *pings* over an extended period is a little over 1% as a long-term average loss rate.

TCP Protocol Performance

TCP^[7] is an instance of a sliding window positive acknowledgement protocol. The sender maintains a local copy of all data that has been passed into the communications systems and discards that data only when it has received a positive acknowledgement from the receiver.

Variants to TCP are based on the variations in the sender's control of the rate of passing data into the network and variations in the response to data loss. The classic version of TCP is one that uses a linear inflation of the sending window size while there is no loss, and halves the window in response to loss.

The algorithm is called the RENO TCP control algorithm. Its use in today's Internet has been largely supplanted by the CUBIC TCP control algorithm^[8], which uses a varying window inflation rate that attempts to stabilise the sending rate at a level just below a level that causes the buildup of network queues, which ultimately leads to packet loss.

In general terms, there is a small set of common assumptions about the characteristics of the network path for such control algorithms:

- There is a *stable* maximal capacity of the path, where the term stability describes a situation where the available path capacity is relatively constant across a number of RTT intervals.
- The amount of *jitter* (variation in end-to-end delay) is low in proportion to the RTT.
- The average packet-loss rate is low. In the case of congestion-based loss, the TCP control algorithm generally interprets packet loss as a sign that the network buffers have filled and the loss is an indication of buffer overflow.

Obviously, as we've noted, the first two conditions do not hold for end-to-end paths that include a Starlink component. The loss profile is also different. There is the potential for congestion-induced packet loss, as is the case in any non-synchronous packet-switched medium, but an additional loss component can occur during satellite handover, and other impairments can further affect the radio signal.

TCP tends to react to such environments by using conservative choices.

The RTT estimate is a smoothed average value of RTT measurements to which is added the mean deviation of individual measurements from this average. For Starlink, with its relatively high level of individual variance in RTT measurements, this estimate means that the TCP sender may operate with a RTT estimate that is too high, which in turn will result in a sending rate that is lower than the available end-to-end capacity of the system.

The occurrence of non-congestion-based loss can also detract from TCP performance. Conventionally, loss will cause the sender to quickly reduce its sending window on the basis that if this loss is caused by network buffer overflow, then the sender needs to allow these buffers to drain. The sender will then resume sending at a lower rate, which should restore coherency of the feedback control loop.

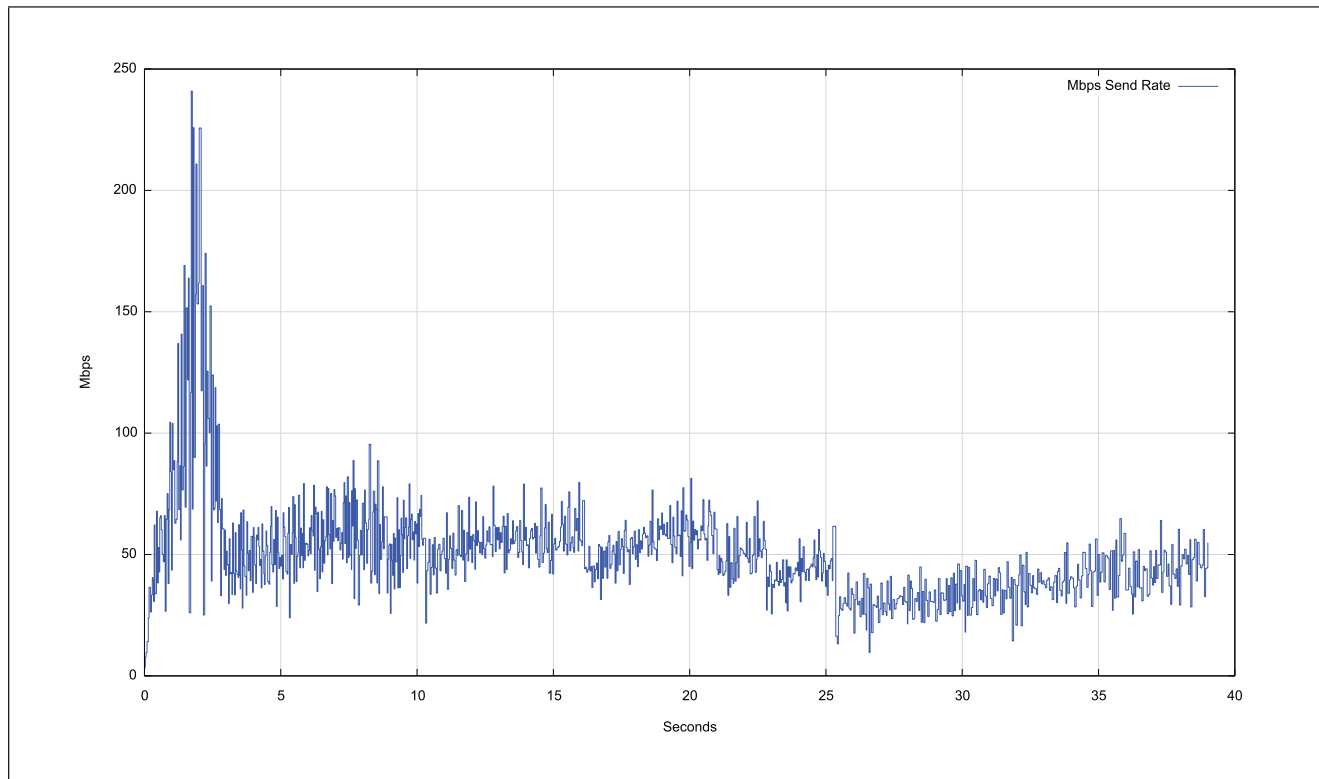
How does this mechanism work in practice?

Figure 5 shows a detailed view of a TCP CUBIC session over a Starlink circuit. The initial 2 seconds show the *slow start* TCP sending rate inflation, where the sending window doubles in size for each RTT interval, reaching a peak of 250 Mbps in 2 seconds. The sender then switches to a rapid reduction of the sending window in the next second, dropping to 50 Mbps within 1 second.

At this point the CUBIC congestion-avoidance phase appears to kick in, and the sending rate increases to 70 Mbps over the ensuing 5 seconds. A single loss event occurs that causes the sending rate to drop back to 40 Mbps in second 8. The remainder of the trace shows this same behaviour of slow sending rate inflation and intermittent rate reductions that are typical of CUBIC.

This CUBIC session managed an average transfer rate of some 45 Mbps, which is well below the peak circuit capacity of 250 Mps.

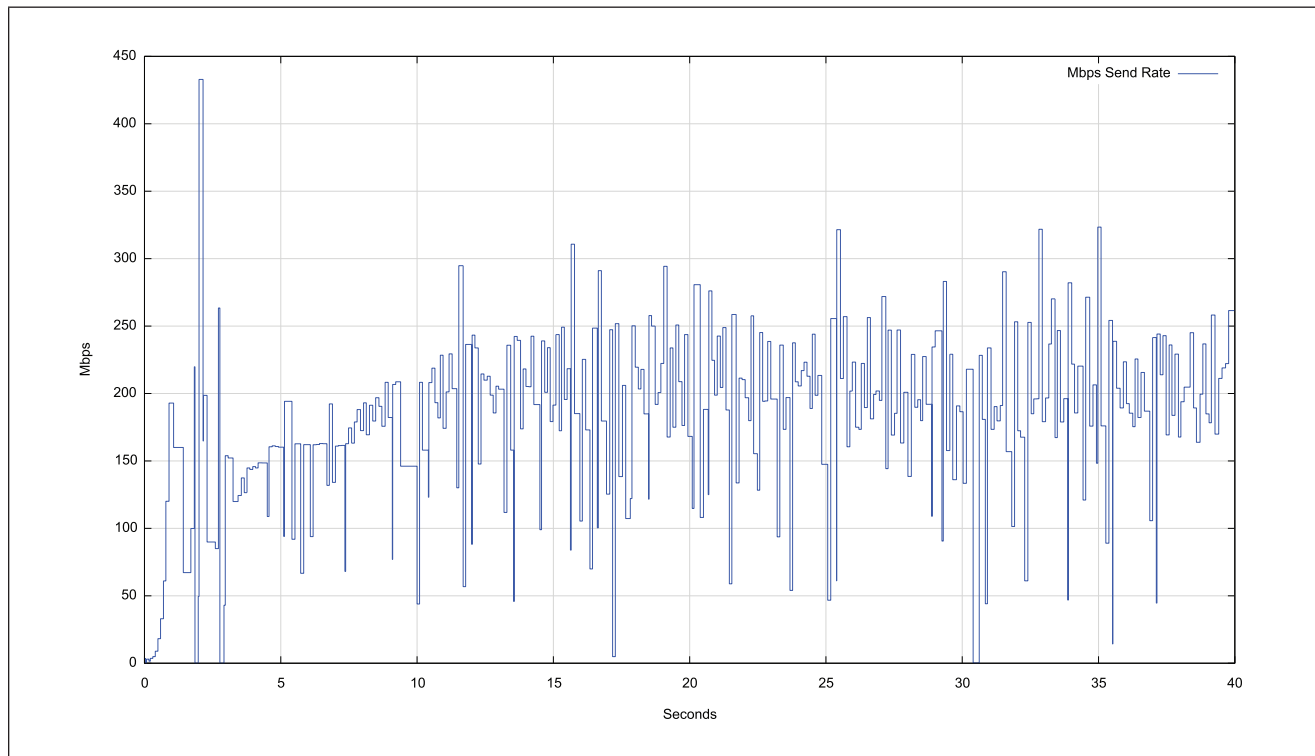
Figure 5: TCP CUBIC Over Starlink



Starlink is a shared medium, and the performance of the system in local times of light use (off peak) is significantly different from that of performance in peak times. Figure 6 shows the CUBIC performance profile during an off-peak time.

The difference between the achievable throughput between peak and off-peak times is quite significant, with the off-peak performance reaching a throughput level some 3 to 4 times greater than the peak-load performance. The slow-start phase increased the throughput to some 200 Mbps within the first second. The flow then oscillated for a second, then started a more stable congestion-avoidance behaviour by second 4. The CUBIC window inflation behaviour is visible up to second 12 and then the flow oscillates around some 200 Mbps of throughput.

Figure 6: TCP CUBIC Over Starlink – Off-Peak



Is the difference between these two profiles in Figures 5 and 6 a result of active flow management by Starlink equipment, or the result of the way in which CUBIC reaches a flow equilibrium with other concurrent flows?

We can attempt to answer this question by using a different TCP control protocol that has a completely different response to contention with other concurrent flows.

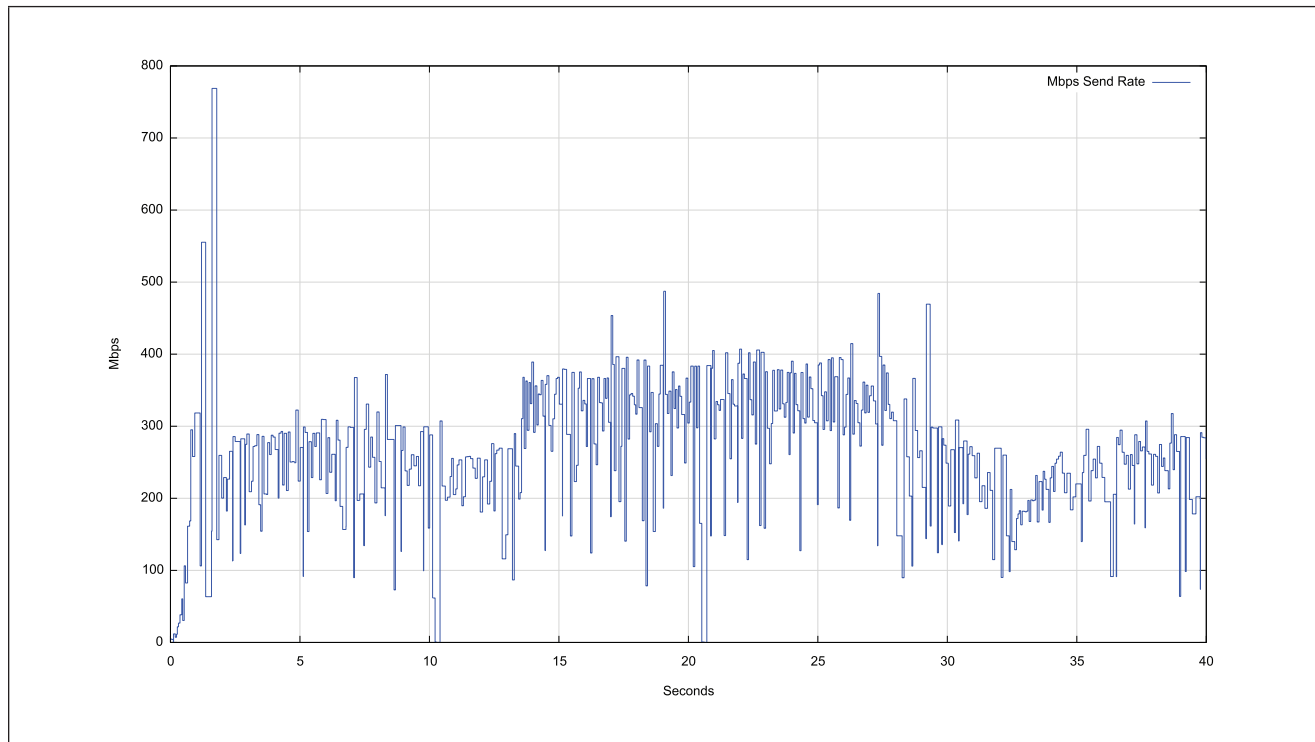
The *Bottleneck Bandwidth and Round-trip propagation time* (BBR)^[9] is a TCP congestion-control algorithm developed at Google a decade ago. BBR attempts to position the TCP flow at the onset of network queue formation rather than oscillating between full and empty queue states (as is the case in loss-based congestion-control algorithms).

Briefly, BBR makes an initial estimate of the delay-bandwidth product of the network path, and then drives the sender to send at this rate for 6 successive RTT intervals. It performs repair for dropped packets without adjusting its sending rate. The 7th RTT interval sees the sending rate increase by 25%, and the end-to-end delay is carefully measured in this interval. The final RTT interval in the cycle sees the sending rate drop by 25% from the original rate, intended to drain any network queues that may have formed in the previous RTT interval. If the end-to-end delay increases in the inflate interval, the original sending rate is maintained.

If the increased sending window does not impact the end-to-end delay, it indicates that the network path has further capacity and the delay-bandwidth estimate is increased for the next 8-RTT cycle. (There have been a couple of subsequent revisions to the BBR protocol, but in this case, I'm using the original (v1) version of BBR.)

Figure 7 shows the results of a Starlink performance test using BBR.

Figure 7: TCP BBR Over Starlink



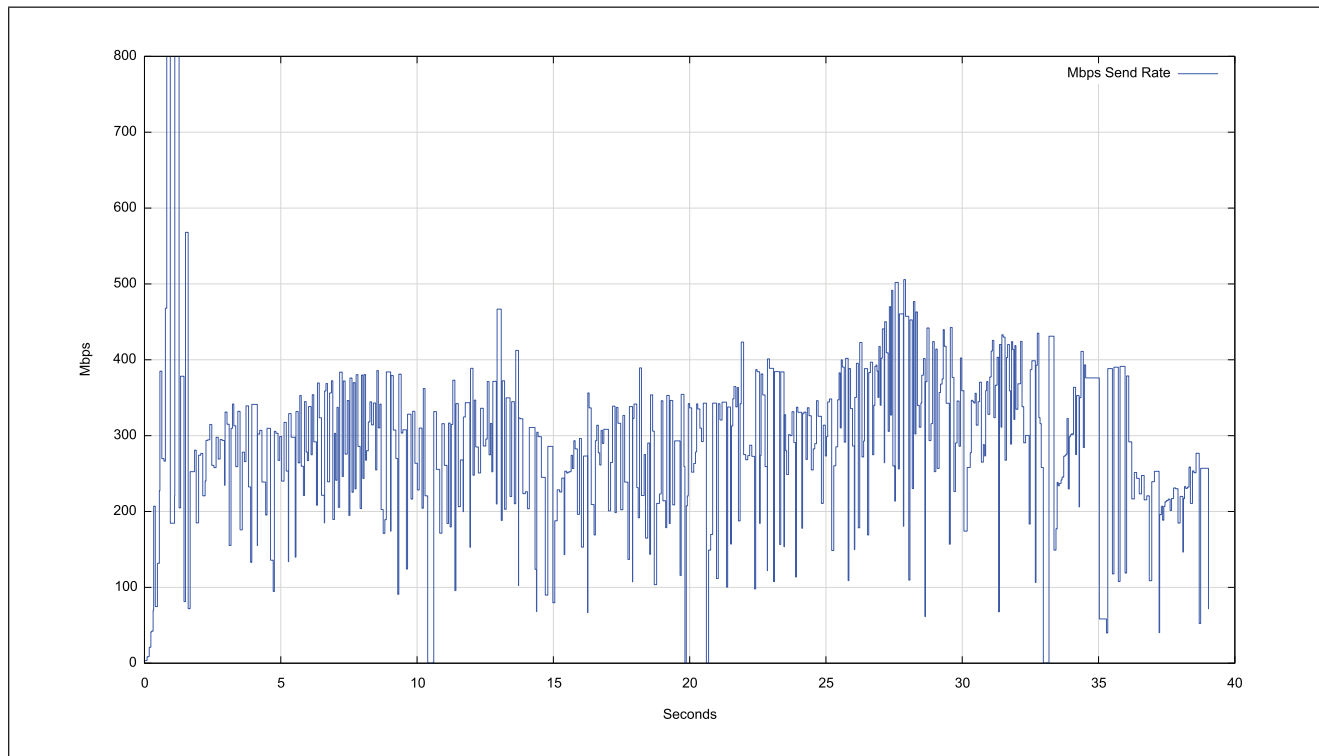
In this case, BBR has made an initial estimate of some 250 Mbps for the path bandwidth. This estimate appears to have been revised at second 14 to 350 Mbps, and then dropped to 200 Mbps 15 seconds later for the final 10 seconds of this test. It is likely that these changes are the result of BBR responding to satellite handover in Starlink.

The same BBR test was performed in an off-peak time and had a very similar outcome (Figure 8 on the following page).

If BBR is sensitive to changes in latency, and latency is so variable in Starlink, then why does BBR perform so well?

I suspect that here BBR is not taking a single latency measurement, but measuring the RTT for all packets that are sent in this 7th RTT interval, and using the minimum RTT value as the “loaded” RTT value to determine whether to perform a send-rate adjustment. As long as the minimum RTT levels are consistent, and they—as shown in Figure 3—are consistent across each 15-second scheduling interval, then BBR will set its sending rate close to the maximum sending rate that Starlink supports.

Figure 8: TCP BBR Over Starlink – Off-Peak



Protocol Tuning for Starlink

Could you tune a variant of TCP to optimise its performance over a path that includes a Starlink component?

A promising approach would appear to be a variant of BBR. The reason for the choice of BBR is its ability to maintain its sending rate in the face of individual packet-loss events. Starlink performs a satellite handover at regular 15-second intervals. If the regular sending-rate inflation in BBR occurs at the same time as scheduled satellite handover, the BBR sender could defer its rate inflation, maintaining its current sending rate across the scheduled handover.

The issue with BBR is that, for version 1 of this protocol, it is quite aggressive in claiming network resources, and this aggression can starve other concurrent sessions of capacity. One possible response is to use the same 15-second satellite handover timer with version 3 of the BBR protocol, which is intended to be less aggressive when working with concurrent data flows.

In theory, it would be possible to adjust CUBIC in a similar manner, performing a lost packet repair using *Selective Acknowledgement* (SACK)^[10] if the packet loss occurred at the time of a scheduled satellite handover. While CUBIC is a fairer protocol with respect to sharing the path capacity with other concurrent sessions, it tends to react conservatively when faced with high jitter paths. Even with some sensitivity to scheduled satellite handovers, CUBIC is still prone to reduced efficiency in the use of network resources.

References and Resources

- [0] Dan York and Geoff Huston, “Low Earth Orbit Satellite Systems for Internet Access,” *The Internet Protocol Journal*, Volume 26, No. 2, September 2023.
- [1] Starlink: <https://www.starlink.com>
- [2] Isaac Newton, *Philosophiæ Naturalis Principia Mathematica*, July 1687.
- [3] Mark Allman, Daniel R. Glover, and Luis A. Sanchez, “Enhancing TCP Over Satellite Channels Using Standard Mechanisms,” RFC 2488, January 1999.
- [4] Wikipedia, Low Earth Orbit:
https://en.wikipedia.org/wiki/Low_Earth_orbit
- [5] Speedtest: <https://www.speedtest.net>
- [6] Ping network utility: [https://en.wikipedia.org/wiki/Ping_\(networking_utility\)](https://en.wikipedia.org/wiki/Ping_(networking_utility))
- [7] Jon Postel, “Transmission Control Protocol,” RFC 793, September 1981.
- [8] Sangtae Ha, Injong Rhee, and Lisong Xu, “CUBIC: a new TCP-friendly high-speed TCP variant,” *ACM SIGOPS Operating Systems Review*, Volume 42, No. 5, July 2008.
- [9] Neal Cardwell, Yuchung Cheng, C. Stephen Gunn, Soheil Hassas Yeganeh, and Van Jacobson, “BBR: Congestion-Based Congestion Control: Measuring bottleneck bandwidth and round-trip propagation time,” *ACM Queue*, Volume 14, No. 5, October 2016.
- [10] Sally Floyd, Jamshid Mahdavi, Matt Mathis, and Matthew Podolsky, “An Extension to the Selective Acknowledgement (SACK) Option for TCP,” RFC 2883, July 2000.
- [11] Josh Fomon, “New Speedtest Data Shows Starlink Users Love Their Provider,” *Ookla Insights Articles*, May 8, 2023.
- [12] Kevin Hurler, “Starlink Is Now Connecting Remote Antarctic Research Camps to the Internet,” *Gizmodo*, January 23, 2023.
- [13] “Perspectives on LEO Satellites: Using Low Earth Orbit Satellites for Internet Access,” Internet Society, 2022.
- [14] Ulrich Speidel, “Satellite still a necessity for many Pacific Islands,” *APNIC Blog*, September 18, 2018.

GEOFF HUSTON AM, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

DNS Evolution

by Geoff Huston, APNIC

The *Domain Name System* (DNS) is a crucial part of today's Internet. With the fracturing of network address space as a byproduct of IPv4 address rundown and the protracted IPv6 transition, the namespace of the Internet is now the defining attribute that makes it one network. However, the DNS is not a rigid and unchanging technology. It has changed considerably over the lifetime of the Internet, and here I'd like to look at what has changed and what has remained the same.

The Early DNS

The early Internet architecture used names as a convenient alias for an IP address. Each host used a local list of name and address pairs, and an application would look up the name in this file (**hosts.txt**) and use the associated address in the subsequent packet exchange. In many ways, this file was a direct analogy to the telephone directory in a telephone network.

This simple framework has one major drawback: *scalability*. As the number of connected hosts on the network increased, the burden of distributing updated copies of the name file increased and the task of maintaining loose coherence across all these local copies of this file became more challenging. The document IEN 61^[1], describing an *Internet Name Server*, was released in 1978, and it appears to be a basic predecessor of today's DNS.

Some five years later, in 1983, RFC 882^[2] defined a hierarchical namespace using a tree-structure name hierarchy. It also defined a name server as a service that holds information about a part of the name hierarchy, and also refers to other name servers that hold information about lower parts of the name hierarchy. The document also defined a resolver that can resolve names into their stored attributes by following referrals to find the appropriate name server to query, and then obtaining this information from the server. RFC 883^[3] defined the DNS query and response protocol, a simple stateless protocol.

And that's about it.

Most of what is in today's DNS was defined in these early specifications, and what we've been doing over the intervening forty years has been filling in the details. The DNS has not really changed in any substantive manner over the intervening period.

Evolutionary Pressures

However, I think that such a perspective ignores a large body of refinement in the DNS world that has occurred. The DNS is by no means perfect; it can be extremely slow to resolve a name, and even slower to incorporate changes into the distributed data framework.

The resolution of DNS queries pays scant regard to concerns about user privacy, and any party who can observe a user's DNS query stream can readily piece together an accurate picture of the user's activities. The distributed stateless method used to resolve names is prone to various efforts to eavesdrop DNS transactions and manipulate the information being provided in DNS responses. The DNS cannot easily protect itself from disruptive attack and has been regularly used in highly effective denial-of-service attacks. It's also insecure, in that a client cannot verify the authenticity and currency of a response.

The operation of the DNS in resolving a name can be extremely opaque. The use of parallel servers and resolvers to improve the resilience of the DNS creates combinatorial explosion in the number of paths that can be used to navigate through the distributed data structure. It is not possible to tell in advance which servers may be used in the resolution of a query, or the number of additional queries a single original query may trigger. Given that resolvers can respond directly to a query with a locally cached response, it is not possible to tell in advance where the response will come from, or if the response is authentic.

For a common and fundamental service that every user not only uses, but implicitly relies upon, the DNS in practice is far from a paragon of sound operational engineering.

The evolutionary efforts have been intended to remedy some of these shortcomings, with goals to improve the speed of DNS resolution, improve aspects of privacy of DNS transactions, increase the level of trust in DNS responses, and resist efforts to subvert the integrity of DNS name-resolution transactions.

DNS Privacy

The DNS is not what you might call a discrete protocol. By default, queries are made in the clear. The IP addresses of the querier, the server being queried, and the name being queried are visible to any party that is in a position to inspect DNS traffic. These parties include not only potential eavesdroppers in the network, but also the operating system platform that hosts the application making the DNS query, the recursive resolver that receives the query, and any forwarding agent that the recursive resolver uses. Depending on the state of the local cache in the recursive resolver, the recursive resolver may need to perform some level of top-down navigation through the nameserver hierarchy, asking an authoritative server at each level the full original query name. The recursive resolver normally lists itself as the source of these queries, so the identity of the original user is occluded, but the query name is still visible.

RFC 7858 provides a specification for DNS over a *Transport Layer Security* (TLS) session (DoT)^[4]. This specification allows the client and server to securely set up a shared session key that is then used to encrypt all subsequent transactions between the two parties. TLS can also be used to authenticate the server name in order to assure the client that it is connecting to an instance of the named server.

There is some overhead to setting up a TLS session, and the most efficient use of this approach is in the stub-to-recursive DNS environment where a single TLS session can be kept open and reused for subsequent queries, amortizing the initial setup overheads across these queries. The standard specification of DoT defines the use of TCP port 853, which allows an onlooker to identify that DoT is being used and identify the two end parties by their IP addresses, but not the DNS queries or responses.

Subsequent standards work has defined *DNS over QUIC* (DoQ), RFC 9250^[5]. The encryption that QUIC provides has properties similar to those that TLS provides, while QUIC transport eliminates the head-of-line blocking issues inherent with TCP and provides more efficient packet-loss recovery than *User Datagram Protocol* (UDP).

In addition, it is possible to add a *Hypertext Transfer Protocol* (HTTP) wrapper to the DNS data object, defining *DNS over HTTPS* (DoH), RFC 8484^[6]. DoH uses port 443, using either TCP in the case of HTTP/2 or UDP with the QUIC-based HTTP/3, so the DNS transactions would be largely indistinguishable from Web traffic. HTTP adds its own ability to perform object caching, redirection, proxying, authentication, and compression beyond that provided in the conventional DNS model, although the use of such HTTP capabilities in the DNS context is not well understood. HTTP also allows a server to push content to a client. In the DoH scenario this possibility could permit the use of queryless DNS, where the server pushes DNS responses to a client without any initial triggering DNS query.

In these approaches to encrypted transport for the DNS, the remote server is aware of the client's IP address and the queries that the client is making. In the stub-to-recursive scenario, this awareness allows the recursive resolver to be privy to the user's DNS actions, even when the network path between the two parties is secure. A stronger level of privacy is obtained by using *Oblivious DNS over HTTPS*, RFC 9230^[7], where no single DNS server is simultaneously aware of the client's IP address and the content of the DNS queries. Here a double level of encryption is used in conjunction with two independent agents within the network. The client sends an encrypted DNS query to the first proxy using DoH. This proxy is aware of the client's IP identity, but is not able to decrypt the DNS query. The proxy makes its own query using the encrypted query to a separate target, again using DoH, but this time there is no record of the original client. The target can decrypt the query and function as a conventional recursive resolver.

These four specifications show that it is possible to cloak DNS transactions within a secure veil of secrecy, but it remains a topic of speculation as to the extent of uptake of these technologies. Encrypted transport sessions impose higher costs on the operation of DNS infrastructure (recursive resolvers and authoritative servers), and it is unclear how the current DNS economic models, where individual DNS queries are essentially unfunded by the client, can absorb these higher costs.

An entirely different approach to improving DNS privacy is described in DNS *Query Name Minimisation*, RFC 7816^[8]. The observation is that as a recursive resolver navigates its path through the DNS hierarchy, it uses the original query name to query authoritative name servers, essentially sharing the knowledge of the name being queried with a set of servers. The rationale for this approach is that the client does not necessarily know where a zone cut may exist in advance. Query Name Minimisation proposes to minimise the amount of information being disclosed to authoritative name servers by sending a request to the nameserver authoritative for the closest known ancestor of the original query name, and asking for a *Name Server* (NS) delegation record rather than the original query type. This approach does not impose additional overheads on DNS server infrastructure. It does not offer channel security, but it does limit the amount of information “leakage” that is a feature of the DNS name-resolution process.

On a more general level, none of these DNS privacy measures can assure users of the authenticity of the DNS response that they receive. These measures limit the ability of other parties to eavesdrop on DNS queries and responses, but detecting (and presumably rejecting) DNS responses that are inauthentic is a separate issue for the DNS.

DNS Authenticity – DNSSEC

Domain Name System Security Extensions (DNSSEC) is an extension to the DNS that associates a cryptographically-generated digital signature with each record in a DNSSEC-signed zone, specified in RFC 4033^[9]. DNSSEC does not change the DNS namespace, nor the DNS name-resolution protocol. Clients who are aware of DNSSEC can request that a DNS response should include a DNSSEC signature, if one is available for the zone, and may then validate the response using that signature.

You might think that a tool that allows the client to verify a DNS response would be immediately popular. If the relationship between the names that applications use and services and IP addresses that are used at the protocol level is disrupted, then users can be readily deceived. Yet, after close to three decades from its initial specification, DNSSEC is still struggling to achieve mainstream adoption. Part of the issue is that the strong binding of the DNS protocol to a UDP transport causes a set of problems when responses bloat in size because of attached signatures and keys. Another part of the issue lies in the care and attention required to manage cryptographic keys and the unforgiving nature of cryptographic validation. And a large part of the problem is that when the Web began using TLS as a means of verifying the identity of a remote server, many didn’t consider any marginal incremental benefit of DNSSEC in the DNS part of session creation to be worth the incremental effort and cost of using DNSSEC.

For these reasons DNSSEC continues in the DNS environment as a “work in progress.”

Evolution of Query Mechanisms

The base DNS specification uses a limited repertoire, where queries contain a query name and a query type, and, if carried over the UDP, DNS responses are limited to 512 bytes in length. The restrictions in the size of several flag fields, return codes, and label types available in the basic DNS protocol were hindering the development of DNSSEC. The chosen path to resolve this dilemma was to use a so-called *Pseudo Resource Record*, the OPT (for “options”) record that is included in the additional data section of a DNS message. To ensure backward compatibility, a responder does not use the OPT record unless it was present in the query. This is the general *Extension Mechanism for DNS*, or EDNS^[10].

EDNS options have been used so far to support DNSSEC functions, padding, TCP keepalive settings, and Client Subnet fields. It has also been used to extend the maximum size of UDP messages in the DNS by using a EDNS Buffer Size.

It is often desirable to separate the name of a service and the location of the service platform that delivers the service, and service record type that was intended to achieve that outcome. *Service Records*, or SRV records, can provide that form of flexibility, where the service is defined by a host name, a port identifier, and a protocol identifier, and the associated resource record provides the TCP or UDP port number and the canonical service name of the target service platform. Multiple service targets can be specified with an associated preference for use. The functional shift in the use of the SRV record was loading the DNS query with a service profile rather than a plain domain name, and in return receiving enough information to enable the user to then connect to the desired service without making further DNS queries.

This functional shift was further extended in the *Service Binding and Parameter Specification via the DNS* (SVCB and HTTPS Resource Records) specification, RFC 9460^[11]. By providing more information to the client before it attempts to establish a connection, these records offer potential benefits to both performance and privacy. These enhancements represent a shift in the design approach of the DNS, where the prior use of DNS resource record types was to segment the information associated with a DNS name, so that a complete collection of information about a service name was obtained by making a set of queries. The SVCB record effectively provides an “omnibus” response to a service query, so that the client can gather sufficient information to connect to a service with a single DNS transaction.

Delegation Records

One of the fundamental parts of the DNS data structure is the *delegation record*, which passes the control of an entire subtree in the DNS hierarchy from one node to another.

While this NS record has served the DNS since its inception, it has a few limitations. The target of the delegation record is one or more DNS server names, not their IP addresses.

Conventionally the IP addresses are provided as “glue records” contained in the *Additional Section* of a DNS referral response. However, the veracity of such glue records cannot be established, and this weakness has been the focal point of numerous DNS attacks over the years. The target of a NS record cannot be a CNAME alias. The NS record is shared across both the parent and child zones, and the child zone is deemed to be authoritative for this record. The implication is that while the parent-zone name servers can (and must) respond with referral responses with this NS record, it cannot provide a DNSSEC-signed response. It is not possible to provide a DNS service profile in a referral response. If the zone authoritative servers can be accessed using an encrypted transport protocol, this capability cannot be signalled by the NS record.

Work is underway in the *Internet Engineering Task Force* (IETF) in the *DNS Delegation* (deleg) Working Group to take the existing specification of service binding mapping for DNS servers, RFC 9461^[12], and see how it could be used as a more flexible delegation record that addresses some or all of these identified shortcomings in the existing NS form of delegation.

Alternate Name Systems

The Internet protocol suite can be regarded as a collection of elements, including addressing, routing, forwarding, and naming, and it's possible to substitute a different technology for one element without necessarily impacting the others. For example, the transition from IP version 4 to IP version 6 in the addressing realm does not necessitate any fundamental changes to routing, forwarding, or naming. The same can be said of the DNS name system. Alternate name systems can be used and to some extent they can coexist with the DNS.

In the traditional model of DNS resolution, users have little control over their DNS settings. Some technically literate users may choose settings that differ from the defaults, but there has been little incentive to do so, and the vast majority of users have their DNS settings configured for them by administrators via a protocol such as the *Dynamic Host Configuration Protocol* (DHCP).

Many alternative naming systems in use today come bundled with the specific applications that use them: a particular alternative naming system is often tied to a corresponding application, and this application often bypasses administrator-controlled settings and any preconfigured DNS settings. For example, the *Tor Project* uses its own naming system that bypasses traditional DNS resolution. Users can install the *Tor Browser*, and it will use the Tor naming system for names ending in **.ONION**, while forwarding any other names to the local DNS library. The application developer makes the choice of which naming system to use without users even knowing that they are using an alternative naming system, nor do they understand potential implications.

Various forms of experimentation have used decentralised models that eschew a single name hierarchy and allow individual names to exist in an unstructured flat namespace. The underlying registry framework that associates a name with an “owner” has often relied on some blockchain-like approach, where the association of a name and a public-key value is placed into the blockchain. Numerous such alternate name systems exist today, including the *Ethereum Name Service* (ENS), which uses so-called “smart contracts” in its blockchain, and *Unstoppable Domains*, which uses a blockchain platform but operates the namespace as a centrally operated space. The *GNU Name System* (GNS) is also a decentralised platform that offers name persistence, but it has no concept of a root zone. Instead GNS uses the concept of a “start zone” that is configured locally and determines where to begin resolution. Since local users have complete control over their own start zone, every GNS user can potentially use a different namespace. Thus, there is no guarantee that names will be globally unique, or that a given name will resolve the same for different users. The only guarantee is that users with the same start zone will have the same view of the namespace. Every unique start zone defines its own namespace. This scenario is similar in practice to DNS resolution using different root zones. The key innovation in GNS is to replace a search hierarchy with a distributed hash table that can include links to other hash tables.

Such alternate name systems interact with the existing DNS-defined namespace in a variety of ways. Some attempt to coexist with the DNS with the alternate names being some form of extension to the DNS namespace, potentially using a different name-resolution protocol. Other systems are completely self-contained and make no effort to coexist with the DNS. This situation is more commonly seen in an application-specific context where the application environment is exclusively associated with an alternate namespace.

Conclusions

Only a completely moribund technology is impervious to change! As digital technologies and services evolve, the demands placed on the associated namespaces also evolve in novel and unpredictable ways.

The DNS is an interesting case in that so far it has been able to respond to the evolving Internet without requiring fundamental changes to the structure of its namespace, the distributed information model, or the name-resolution protocol. Most of the evolutionary changes that have been folded into the DNS to date have been undertaken in a way that preserves backward compatibility, and the cohesion of the underlying namespace has been largely preserved.

However, maintaining this cohesion across the Internet is not an assured outcome for the future. The pressures to impose barriers to the access to content at national and regional levels are often expressed by imposing selective barriers to the resolution of content service names, and the DNS is left carrying the burden of supporting such selective fragmentation in the Internet.

The camel has undeniably poked its nose into the tent of name coherence in the form of EDNS *Client Subnet*^[13], where the response given to a query may be dependent on who is querying, as much as the name that is being used in the query, and it's likely that this more qualified and fragmented model of a namespace will persist and support an increasingly fragmented Internet.

References and Further Reading

- [1] Jon Postel, "Internet Name Server," IEN 61, October 1978.
- [2] Paul Mockapetris, "Domain names: Concepts and facilities," RFC 882, November 1983.
- [3] Paul Mockapetris, "Domain names: Implementation specification," RFC 883, November 1983.
- [4] Zi Hu, Liang Zhu, John Heidemann, Allison Mankin, Duane Wessels, and Paul Hoffman, "Specification for DNS over Transport Layer Security (TLS)," RFC 7858, May 2016.
- [5] Christian Huitema, Sara Dickinson, and Allison Mankin, "DNS over Dedicated QUIC Connections," RFC 9250, May 2022
- [6] Paul Hoffman, and Patrick McManus, "DNS Queries over HTTPS (DoH)," RFC 8484, October 2018.
- [7] Eric Kinnear, Patrick McManus, Tommy Pauly, Tanya Verma, and Christopher A. Wood, "Oblivious DNS over HTTPS," RFC 9230, June 2022.
- [8] Stephane Bortzmeyer, "DNS Query Name Minimisation to Improve Privacy," RFC 7816, March 2016.
- [9] Roy Arends, Rob Austein, Matt Larson, Dan Massey, and Scott Rose, "DNS Security Introduction and Requirements," RFC 4033, March 2005.
- [10] Joao Damas, Michael Graff, and Paul Vixie, "Extension Mechanisms for DNS (EDNS(0))," RFC 6891, April 2013.
- [11] Ben Schwartz, Mike Bishop, and Erik Nygren, "Service Binding and Parameter Specification via the DNS (SVCB and HTTPS Resource Records)," RFC 9460, November 2023.
- [12] Ben Schwartz, "Service Binding Mapping for DNS Servers," RFC 9461, November 2023.
- [13] Carlo Contavalli, Wilmer van der Gaast, David C. Lawrence, and Warren Kumari, "Client Subnet in DNS Queries," RFC 7871, May 2016.
- [14] Miek Gieben, "DNSSEC: The Protocol, Deployment, and a Bit of Development," *The Internet Protocol Journal*, Volume 7, No. 2, June 2004.
- [15] Richard Barnes, "Let the Names Speak for Themselves: Improving Domain Name Authentication with DNSSEC and DANE," *The Internet Protocol Journal*, Volume 15, No.1, March 2012.

- [16] Geoff Huston, “A Question of DNS Protocols,” *The Internet Protocol Journal*, Volume 17, No. 1, September 2014.
- [17] Geoff Huston, “What’s in a DNS Name?” *The Internet Protocol Journal*, Volume 19, No. 1, March 2016.
- [18] Geoff Huston and Joao Luis Silva Dama, “DNS Privacy,” *The Internet Protocol Journal*, Volume 20, No. 1, March 2017.
- [19] Geoff Huston, “The Root of the DNS,” *The Internet Protocol Journal*, Volume 20, No. 2, June 2017.
- [20] Geoff Huston, “DNS Privacy and the IETF,” *The Internet Protocol Journal*, Volume 22, No. 2, July 2019.
- [21] Geoff Huston, “DNS Trends,” *The Internet Protocol Journal*, Volume 24, No. 1, March 2021.
- [22] Burton Kaliski Jr., “Minimized DNS Resolution: Into the Penumbra,” *The Internet Protocol Journal*, Volume 25, No. 3, December 2022.
- [23] Wikipedia article on the DNS:
https://en.wikipedia.org/wiki/Domain_Name_System

GEOFF HUSTON AM, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator.
E-mail: **gih@apnic.net**

Check your Subscription Details!

Make sure that both your postal and e-mail addresses are up-to-date since these are the only methods by which we can contact you. If you see the words “Invalid E-mail” on your printed copy this means that we have been unable to contact you through the e-mail address on file. If this is the case, please contact us at **ipj@protocoljournal.org** with your new information. The subscription portal is located here:
<https://www.ipjsubscription.org/>

An Open Letter to the United Nations

1 July 2024

Secretary-General António Guterres
and Envoy on Technology Amandeep Singh Gill,

Since its inception more than fifty years ago, the Internet's technical architecture has evolved and been collaboratively maintained through multistakeholder processes. While it was born in government laboratories, the Internet became a network of networks that kept expanding and required continuous work. Much of that was coordinated in the *Internet Engineering Task Force* (IETF)^[1], an open, consensus-based, bottom-up, voluntary and global standards body.

More than thirty-five years ago, the World Wide Web was born in the laboratories of CERN. It, too, quickly evolved into a global public tool, maintained and developed by a collaboration of like-minded engineers and other stakeholders at the *World Wide Web Consortium* (W3C)^[2]. It, too, is an open, bottom-up, consensus-driven, voluntary and global standards body.

The success of both IETF's and W3C's work can be measured by where the Internet is today and what it has achieved: global communication has flourished, bringing education, entertainment, information, connectivity and commerce to most of the world's population. The Internet has been a catalyst for advancing development. These communities and the way they have structured themselves have paid off.

We recognize that governments take seriously their responsibility to protect their citizens. So, as harms associated with the Internet and the Web become more apparent, there is a desire on the part of governments to act through regulation and legislation. Technical architecture can enable and influence how the Internet is used, but on its own it cannot address abuse, misinformation, inequality, or many other issues. There is nevertheless a potential danger in regulation and legislation, if it undermines the fundamentally empowering nature of the Internet.

The Internet is an unusual technology because it is fundamentally distributed. It is built up from all of the participating networks. Each network participates for its own reasons according to its own needs and priorities. And this means, necessarily, that there is no center of control on the Internet. This feature is an essential property of the Internet, and not an accident. Yet over the past few years we have noticed a willingness to address issues on the Internet and Web by attempting to insert a hierarchical model of governance over technical matters. Such proposals concern us because they represent an erosion of the basic architecture.

In particular, some proposals for the *Global Digital Compact* (GDC)^[3] can be read to mandate more centralized governance. If the final document contains such language, we believe it will be detrimental to not only the Internet and the Web, but also to the world's economies and societies.

Furthermore, we note that the GDC is being developed in a multi-lateral process between states, with very limited application of the open, inclusive and consensus-driven methods by which the Internet and Web have been developed to date. Beyond some high-level consultations, non-government stakeholders (including Internet technical standards bodies and the broader technical community) have had only weak ways to participate in the GDC process. We are concerned that the document will be largely a creation only of governments, disconnected from the Internet and the Web as people all over the world currently experience them.

Therefore, we ask that member states, the Secretary-General and the Tech Envoy seek to ensure that proposals for digital governance remain consistent with the enormously successful multistakeholder Internet governance practice that has brought us the Internet of today. Government engagement in digital and Internet governance is needed to deal with many abuses of this global system but it is our common responsibility to uphold the bottom-up, collaborative and inclusive model of Internet governance that has served the world for the past half century.

Signed,

All signatures are in a personal capacity; affiliations are informational only.

Daniel Appelquist, W3C TAG co-chair

David Baron, former W3C TAG

Hadley Beeman, W3C TAG

Robin Berjon, former W3C TAG; former
W3C HTML Activity Lead

Andrew Betts, former W3C TAG

Sir Tim Berners-Lee, inventor of the World
Wide Web; founder & emeritus director,
W3C

Tim Bray, former W3C TAG; Editor of XML
(W3C), JSON (IETF)

Randy Bush, former IESG, former ISO/WG13

Dr. Brian E. Carpenter, former Group
Leader, Communication Systems, CERN;
former IAB chair; former ISOC BoT chair;
former IETF chair

Vint Cerf, Internet Pioneer

David Conrad, former IANA general
manager; former ICANN CTO

Martin Duke, former IESG

Dr. Lars Eggert, former IETF chair;
former IRTF chair

David Jack Farber, former IAB; former ISOC
BoT; former Chief Technologist USA FCC

Dr. Stephen Farrell, Trinity College Dublin;
former IESG; former IAB

Demi Getschko, .br

Christian Huitema, former IAB chair

Geoff Huston, former ISOC BoT chair;
former IAB

Erik Kline, IESG

Mallory Knodel, former IAB

Olaf Kolkman, former IAB chair

Konstantinos Komaitis, senior resident
fellow, Internet Governance lead,
Democracy and Tech Initiative,
Atlantic Council

Chris Lilley, W3C Technical Director;
former W3C TAG

Peter Linss, W3C TAG co-chair

Sangwhan Moon, former W3C TAG

Jun Murai, former IAB; WIDE Project
founder; former W3C steering
committee; former ISOC BoT

Mark Nottingham, former IAB;
former W3C TAG

Lukasz Olejnik, former W3C TAG

Colin Perkins, IRTF chair

Pete Resnick, former IAB; former IESG

Alex Russell, former W3C TAG

Peter Saint-Andre, former IESG

David Schinazi, IAB

Melinda Shore, IRSG; former IAB

Robert Sparks, former IAB; former IESG

Lynn St. Amour, former Internet Society
President and CEO; former UN IGF

Multistakeholder Advisory Group chair

Andrew Sullivan, former IAB chair

Martin Thomson, W3C TAG; former IAB

Brian Trammell, IRSG; former IAB

Léonie Watson, W3C Web Applications
Working Group Chair

Paul Wouters, IESG

References and Acronym Expansions

- [1] Internet Engineering Task Force (IETF): <https://ietf.org/>
- [2] World Wide Web Consortium (W3C): <https://w3.org/>
- [3] Global Digital Compact (GDC):
<https://www.un.org/techenvoy/global-digital-compact>
- [4] Internet Architecture Board (IAB): <https://iab.org/>
- [5] Internet Engineering Steering Group (IESG): <https://iesg.org/>
- [6] Internet Research Steering Group (IRSG):
<https://www.irtf.org/irsg.html>
- [7] Internet Research Task Force (IRTF): <https://www.irtf.org/>
- [8] ISOC BoT: Internet Society Board of Trustees:
<https://www.internetsociety.org/board-of-trustees/>
- [9] World Wide Web Consortium Technical Architecture Group (W3C TAG): <https://w3ctag.org/>
- [10] Originally posted here:
<https://open-internet-governance.org/letter>

Call for Papers: IAB Workshop on AI-Control

The *Internet Architecture Board* (IAB) is planning a workshop to explore practical opt-out mechanisms for *Artificial Intelligence* (AI), and build an understanding of use cases, requirements, and other considerations in this space. The workshop will be held in September 2024 in the Washington, DC area. Exact dates and location to be confirmed soon. The deadline for submissions is August 2nd, 2024 and invitations will be issued by August 15th, 2024.

Large Language Models (LLM) and other machine learning techniques require voluminous input data, and one common source of such data is the Internet—usually, “crawling” Web sites for publicly available content, much in the same way that search engines crawl the Web. This similarity has led to an emerging practice of allowing the *Robots Exclusion Protocol*, defined in RFC 9309, to control the behavior of AI-oriented crawlers.

This emerging practice raises many design and operational questions. It is not yet clear whether **robots.txt** (the mechanism specified by RFC 9309) is well-suited to controlling AI crawlers. A content creator or host may not be able to distinguish a crawler used for search indexing from a crawler used for LLM ingest—and indeed some crawlers may be used for both purposes. Potential use cases may extend across many different units of content, policies to be signaled, and types of content creators. Before **robots.txt** becomes a de facto solution to AI crawling opt-out, it is necessary to examine whether it is an appropriate mechanism: in particular, whether the creator of a particular unit of content can realistically and fully exercise their right to opt-out, and the scope of data ingest to which that opt-out applies.

This workshop aims to explore practical opt-out mechanisms for AI, and build an understanding of use cases, requirements, and other considerations in this space. It will focus on mechanisms to communicate the opt-out choice and their associated data models. Technical enforcement of opt-out signals is not in scope. The IAB is looking for short position papers on the following topics; however, this list is non-exhaustive and should be interpreted broadly:

- User stories, use cases, and requirements for opting content out of inclusion in large language models, from a variety of sources including but not limited to the Web
- Interactions between opt-out mechanisms and different use cases for AI
- Advantages and/or deficiencies of reusing robots.txt for controlling AI crawlers on the Web
- Comparisons of use cases for crawling opt-out
- Desired properties of an AI opt-out mechanism
- Potential developments in AI that may require adjustments in opt-out mechanisms
- Implications of legal/policy frameworks (for example, copyright, privacy, research ethics) and requirements on the design of opt-out mechanisms
- Evolution of opt-out signals

Because **robots.txt** is emerging as a solution in this space, the discussion will be anchored on it as a starting point, but not limited to that mechanism. Proposals for alternative solutions may be made, but time will not be available for a detailed presentation or discussion.

Interested participants are invited to submit position papers on the workshop topics. Participants can choose their preferred format, including Internet-Drafts, text- or Word-based documents, or papers formatted similar as used by academic publication venues. Submission as PDF is preferred. Paper size is not limited, but brevity is encouraged. By default, submissions that are considered relevant will be published on the workshop website. If you wish for your submission to be anonymous or withheld from such publication, please indicate that clearly in the submission. The organizers will issue invitations based on the submissions received. Sessions will be organized according to the submissions received, and not every accepted submission or invited attendee will have an opportunity to present; the intent is to foster an active discussion and not simply to have a sequence of presentations.

Discussion at the workshop will be held under *Chatham House Rule*, and therefore will not be recorded or minuted. However, a workshop report will be published afterwards.

It is anticipated that the workshop report will include:

- A list of participants (unless they request to be withheld)
- Documentation of use cases and requirements discussed
- Recommendations for IETF standards work to be considered (if any)
- Recommendations for non-IETF standards work to be considered (if any)

The workshop will be by invitation only. Those wishing to attend should submit a position paper to ai-control-workshop-pc@iab.org. Position papers from those not planning to attend the workshop themselves are also encouraged. Feel free to contact the Program Committee with any further questions: ai-control-workshop-pc@iab.org.

For more information, visit:

<https://datatracker.ietf.org/group/aicontrolws/about/>

Our Privacy Policy

The *General Data Protection Regulation* (GDPR) is a regulation for data protection and privacy for all individual citizens of the *European Union* (EU) and the *European Economic Area* (EEA). Its implementation in May 2018 led many organizations worldwide to post or update privacy statements regarding how they handle information collected in the course of business. Such statements tend to be long and include carefully crafted legal language. We realize that we may need to provide similar language on our website and in the printed edition, but until such a statement has been developed here is an explanation of how we use any information you have supplied relating to your subscription:

- The mailing list for *The Internet Protocol Journal* (IPJ) is entirely “opt in.” We never have and never will use mailing lists from other organizations for any purpose.
- You may unsubscribe at any time using our online subscription system or by contacting us via e-mail. We will honor any request to remove your name and contact information from our database.
- We will use your contact information only to communicate with you about your subscription; for example, to inform you that a new issue is available, that your subscription needs to be renewed, or that your printed copy has been returned to us as undeliverable by the postal authorities.
- We will never use your contact information for any other purpose or provide the subscription list to any third party other than for the purpose of distributing IPJ by post or by electronic means.
- If you make a donation in support of the journal, your name will be listed on our website and in print unless you tell us otherwise.

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Kjetil Aas	Václav Brožík	Ernesto Doelling	Greg Goddard	Jose Enrique Diaz Jolly
Fabrizio Accatino	Christophe Brun	Michael Dolan	Tiago Goncalves	Jonatan Jonasson
Michael Achola	Gareth Bryan	Eugene Doroniuk	Ron Goodheart	Daniel Jones
Martin Adkins	Ron Buchalski	Michael Dragone	Octavio Alfageme	Gary Jones
Melchior Aelmans	Paul Buchanan	Joshua Dreier	Gorostiaga	Jerry Jones
Christopher Affleck	Stefan Buckmann	Lutz Drink	Barry Greene	Michael Jones
Scott Aitken	Caner Budakoglu	Aaron Dudek	Jeffrey Greene	Amar Joshi
Jacobus Akkerhuis	Darrell Budic	Dmitriy Dudko	Richard Gregor	Javier Juan
Antonio Cuñat Alario	BugWorks	Andrew Dul	Martijn Groenleer	David Jump
William Allaire	Scott Burleigh	Joan Marc Riera	Geert Jan de Groot	Anders Marius Jørgensen
Nicola Altan	Chad Burnham	Duocastella	Ólafur Guðmundsson	Merike Kaeo
Shane Amante	Randy Bush	Pedro Duque	Christopher Guemez	Andrew Kaiser
Marcelo do Amaral	Colin Butcher	Holger Durer	Gulf Coast Shots	Vladislav Kalinovskiy
Matteo D'Ambrosio	Jon Harald Bøvre	Karlheinz Dölger	Sheryll de Guzman	Naoki Kambe
Selva Anandavel	Olivier Cahagne	Mark Eanes	Rex Hale	Akbar Kara
Jens Andersson	Antoine Camerlo	Andrew Edwards	Jason Hall	Christos Karayiannis
Danish Ansari	Tracy Camp	Peter Robert Egli	James Hamilton	Daniel Karrenberg
Finn Arildsen	Brian Candler	George Ehlers	Darow Han	David Kekar
Tim Armstrong	Fabio Caneparo	Peter Eisses	Handy Networks LLC	Stuart Kendrick
Richard Artes	Roberto Canonico	Torbjörn Eklöv	Stephen Hanna	Robert Kent
Michael Aschwanden	David Cardwell	Y Ertur	Martin Hannigan	Thomas Kernen
David Atkins	Richard Carrara	ERNW GmbH	John Hardin	Jithin Kesavan
Jac Backus	John Cavanaugh	ESdatCo	David Harper	Jubal Kessler
Jaime Badua	Lj Cemerar	Steve Esquivel	Edward Hauser	Shan Ali Khan
Bent Bagger	Dave Chapman	Jay Etchings	David Hauweele	Nabeel Khatri
Eric Baker	Stefanos Charchalakakis	Mikhail Evstiounin	Marilyn Hay	Dae Young Kim
Fred Baker	Molly Cheam	Bill Fenner	Headcrafts SRLS	William W. H. Kimandu
Santosh Balagopalan	Pierluigi Checchi	Paul Ferguson	Hidde van der Heide	John King
William Baltas	Greg Chisholm	Ricardo Ferreira	Johan Helsingius	Russell Kirk
David Bandinelli	David Chosrova	Kent Fichtner	Robert Hinden	Gary Klesk
A C Barber	Marcin Cieslak	Ulrich N Fierz	Michael Hippert	Anthony Klopp
Benjamin Barkin-Wilkins	Lauris Cikovskis	Armin Fisslthaler	Damien Holloway	Henry Kluge
Feras Batainah	Brad Clark	Michael Fiumano	Alain Van Hoof	Michael Kluk
Michael Bazarewsky	Narelle Clark	The Flirble Organisation	Edward Hotard	Andrew Koch
David Belson	Horst Clausen	Jean-Pierre Forcioli	Bill Huber	Ia Kochiashvili
Richard Bennett	James Cliver	Gary Ford	Hagen Hultzsich	Carsten Koempe
Matthew Best	Guido Coenders	Susan Forney	Kauto Huopio	Richard Koene
Hidde Beumer	Robert Collet	Christopher Forsyth	Asbjørn Højmark	Alexander Kogan
Pier Paolo Biagi	Joseph Connolly	Andrew Fox	Kevin Iddles	Matthijs Koot
Arturo Bianchi	Steve Corbató	Craig Fox	Mika Ilvesmaki	Antonin Kral
John Bigrow	Brian Courtney	Fausto Franceschini	Karsten Iwen	Robert Krejčí
Orvar Ari Bjarnason	Beth and Steve Crocker	Erik Fredriksson	Joseph Jackson	John Kristoff
Tyson Blanchard	Dave Crocker	Valerie Fronczak	David Jaffe	Terje Krogdahl
Axel Boeger	Kevin Croes	Tomislav Futivic	Ashford Jaggernaut	Bobby Krupczak
Keith Bogart	John Curran	Laurence Gagliani	Thomas Jalkanen	Murray Kuchera
Mirko Bonadei	André Danthine	Edward Gallagher	Jozef Janitor	Warren Kumari
Roberto Bonalumi	Morgan Davis	Andrew Gallo	Martijn Jansen	George Kuo
Lolke Boonstra	Jeff Day	Chris Gamboni	John Jarvis	Dirk Kurfuerst
Cente Cornelis Boot	Fernando Saldana Del	Xosé Bravo Garcia	Dennis Jennings	Mathias Körber
Julie Bottorff Photography	Castillo	Oswaldo Gazzaniga	Edward Jennings	Darrell Lack
Gerry Boudreaux	Rodolfo Delgado-Bueno	Kevin Gee	Aart Jochem	Andrew Lamb
Leen de Braal	Julien Dhallenne	Rodney Gehrke	Nils Johansson	Richard Lamb
Kevin Breit	Freek Dijkstra	Radu Cristian Gheorghiu	Brian Johnson	Yan Landriault
Thomas Bridge	Geert Van Dijk	Greg Giessow	Curtis Johnson	Edwin Lang
Ilia Bromberg	David Dillow	John Gilbert	Richard Johnson	Sig Lange
Lukasz Bromirski	Richard Dodsworth	Serge Van Ginderachter	Jim Johnston	Markus Langenmair

Fred Langham	Eduard Metz	Harald Pilz	James Schneider	Lorin J Thompson
Tracy LaQuey Parker	William Mills	Derrell Piper	Peter Schoo	Fabrizio Tivano
Christian de Larrinaga	David Millsom	Rob Pirnie	Dan Schrenk	Peter Tomsu Fine Art
Alex Latzko	Desiree Miloshevic	Jorge Ivan Pincay	Richard Schultz	Photography
Jose Antonio Lazaro	Joost van der Minnen	Ponce	Timothy Schwab	Joseph Toste
Lazaro	Thomas Mino	Marc Vives Piza	Roger Schwartz	Rey Tucker
Antonio Leding	Rob Minshall	Victoria Poncini	SeenThere	Sandro Tumini
Rick van Leeuwen	Wijnand Modderman-	Blahoslav Popela	Scott Seifel	Angelo Turetta
Simon Leinen	Lenstra	Andrew Potter	Paul Selkirk	Brian William Turnbow
Robert Lewis	Mohammad Moghaddas	Ian Potts	Andre Serralheiro	Michael Turzanski
Christian Liberale	Charles Monson	Eduard Llull Pou	Yury Shefer	Phil Tweedie
Martin Lillepuu	Andrea Montefusco	Tim Pozar	Yaron Sheffer	Steve Ulrich
Roger Lindholm	Fernando Montenegro	David Preston	Doron Shikmoni	Unitek Engineering AG
Link Light Networks	Roberto Montoya	David Raistrick	Tj Shumway	John Urbanek
Art de Llanos	Joel Moore	Priyan R Rajeevan	Jeffrey Sicuranza	Martin Urwaleck
Mike Lochocki	Joseph Moran	Balaji Rajendran	Thorsten Sideboard	Betsy Vanderpool
Chris and Janet Lonvick	John More	Paul Rathbone	Greipur Sigurdsson	Surendran Vangadasalam
Mario Lopez	Maurizio Moroni	William Rawlings	Fillipe Cajaiba da Silva	Ramnath Vasudha
Sergio Loreti	Brian Mort	Mujtiba Raza Rizvi	Andrew Simmons	Randy Veasley
Eric Louie	Soenke Mumm	Bill Reid	Pradeep Singh	Philip Venables
Adam Loveless	Tariq Mustafa	Petr Rejhon	Henry Sinnreich	Buddy Venne
Josh Lowe	Stuart Nadin	Robert Remenyi	Geoff Sisson	Alejandro Vennera
Guillermo a Loyola	Michel Nakhla	Rodrigo Ribeiro	John Sisson	Luca Ventura
Hannes Lubich	Mazdak Rajabi Nasab	Glenn Ricart	Helge Skrivervik	Scott Vermillion
Dan Lynch	Krishna Natarajan	Justin Richards	Terry Slattery	Tom Vest
David MacDuffie	Naveen Nathan	Rafael Riera	Darren Sleeth	Peter Villemoes
Sanya Madan	Darryl Newman	Mark Risinger	Richard Smit	Vista Global Coaching
Miroslav Madić	Mai Nguyen	Fernando Robayo	Bob Smith	& Consulting
Alexis Madriz	Thomas Nikolajsen	Michael Roberts	Courtney Smith	Dario Vitali
Carl Malamud	Paul Nikolich	Gregory Robinson	Eric Smith	Rüdiger Volk
Jonathan Maldonado	Travis Northrup	Ron Rockrohr	Mark Smith	Jeffrey Wagner
Michael Malik	Marijana Novakovic	Carlos Rodrigues	Tim Sneddon	Don Wahl
Tarmo Mammers	David Oates	Magnus Romedahl	Craig Snell	Michael L Wahrman
Yogesh Mangar	Ovidiu Obersterescu	Lex Van Roon	Job Snijders	Lakhinder Walia
John Mann	Jim Oplotnik	Marshall Rose	Ronald Solano	Laurence Walker
Bill Manning	Tim O'Brien	Alessandra Rosi	Asit Som	Randy Watts
Diego Mansilla	Mike O'Connor	David Ross	Ignacio Soto Campos	Andrew Webster
Harold March	Mike O'Dell	William Ross	Evandro Sousa	Jd Wegner
Vincent Marchand	John O'Neill	Boudhayan	Peter Spekreijse	Tim Weil
Normando Marcolongo	Carl Ötne	Roychowdhury	Thayumanavan Sridhar	Westmoreland
Gabriel Marroquin	Packet Consulting Limited	Carlos Rubio	Paul Stancik	Engineering Inc.
David Martin	Carlos Astor Araujo	Rainer Rudigier	Ralf Stempffer	Rick Wesson
Jim Martin	Palmeira	Timo Rüter	Matthew Stenberg	Peter Whimp
Ruben Tripiana Martin	Gordon Palmer	RustedMusic	Martin Štěpánek	Russ White
Timothy Martin	Alexis Panagopoulos	Babak Saberi	Adrian Stevens	Jurrien Wijlhuizen
Carles Mateu	Gaurav Panwar	George Sadowsky	Clinton Stevens	Joseph Williams
Juan Jose Marin Martinez	Chris Parker	Scott Sandefur	John Streck	Derick Winkworth
Ioan Maxim	Alex Parkinson	Sachin Sapkal	Martin Streule	Pindar Wong
David Mazel	Craig Partridge	Arturas Satkovskis	David Strom	Makarand Yerawadekar
Miles McCredie	Manuel Uruena Pascual	PS Saunders	Colin Strutt	Phillip Yialeloglou
Gavin McCullagh	Ricardo Patara	Richard Savoy	Viktor Sudakov	Janko Zavernik
Brian McCullough	Dipesh Patel	John Sayer	Edward-W. Suor	Bernd Zeimet
Joe McEachern	Dan Paynter	Phil Scarr	Vincent Surillo	Muhammad Ziad
Alexander McKenzie	Leif Eric Pedersen	Gianpaolo Scassellati	Terence Charles Sweetser	Ziauddin
Jay McMaster	Rui Sao Pedro	Elizabeth Scheid	T2Group	Tom Zingale
Mark Mc Nicholas	Juan Pena	Jeroen Van Ingen	Roman Tarasov	Matteo Zovi
Olaf Mehlberg	Luis Javier Perez	Schenau	David Theese	Jose Zumalave
Carsten Melberg	Chris Perkins	Carsten Scherb	Rabbi Rob and	Romeo Zwart
Kevin Menezes	Michael Petry	Ernest Schirmer	Lauren Thomas	廖明沂.
Bart Jan Menkveld	Alexander Peuchert	Benson Schliesser	Douglas Thompson	
Sean Mentzer	David Phelan	Philip Schneck	Kerry Thompson	

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at **ole@protocoljournal.org** or **olejacobsen@me.com**

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

Supporters



Diamond Sponsors

Your logo here!

Ruby Sponsors



Sapphire Sponsors

Your logo here!

Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

John Crain, Senior Vice President and Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Merike Kaeo, Founder and vCISO
Double Shot Security

Olaf Kolkman, Principal – Internet Technology, Policy, and Advocacy
The Internet Society

Dr. Jun Murai, Founder, WIDE Project
Distinguished Professor, Keio University
Co-Director, Keio University Cyber Civilization Research Center, Japan

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

October 2024

Volume 27, Number 3

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

FROM THE EDITOR

In This Issue

From the Editor	1
The Interop Shownet	2
“IPv6 Mostly” Experiment... ..	16
Fragments	22
Thank You.....	28
Call for Papers.....	30
Supporters and Sponsors	31

These days, it is common for Internet technology conferences to deploy a *temporary* network in convention centers or hotels to support Internet access for attendees and exhibitors. In some cases, these networks are used for special technology demonstrations. In this issue we will look at two examples of such networks.

The *TCP/IP Interoperability Conference*—later renamed *Interop*—began as a small workshop in August 1986. It quickly grew in scope to incorporate tutorials, and by 1988 an exhibition network connected 51 exhibitors to each other and to the global Internet. This network, or “Shownet,” was designed and deployed by a group of volunteers, and it became the proving ground for many emerging technologies. In 1994, Interop added Tokyo to its international venues, where 30 years later the conference and exhibition attracts more than 120,000 visitors. We will publish a separate article about the Japanese version of Shownet in a future edition. This time David Strom describes the history and evolution of the Interop Shownet. His article is dedicated to the memory of Daniel C. Lynch, the founder of Interop, who passed away earlier this year.

Work on the transition from IPv4 to IPv6 continues to be a major focus of several working groups in the *Internet Engineering Task Force* (IETF). As solutions are developed, technology events such as the *Asia Pacific Internet Conference on Operational Technologies* (APRICOT) provide end users an opportunity to experience IPv6 by simply selecting a designated Wi-Fi network on their devices. One such “IPv6 Mostly” experiment was conducted during APRICOT 2024 in Bangkok; Brian Candler describes it in our second article.

Ten years ago, *The Internet Protocol Journal* was relaunched with the help of numerous supporters, sponsors, and individual donors. Today we very much depend on this funding model, and once again we encourage you to make a donation or ask your organization to become a sponsor. We appreciate your feedback and suggestions. Please contact us via e-mail at: ipj@protocoljournal.org

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

The Interop Shownet

by David Strom

This is the story about how a group of very dedicated people came together at the dawn of the Internet era to build something special, something unique and memorable. It was called the *Interop Show 'n Tel-Net*, later known as *InteropNet*, or *Shownet*, and it was created in September 1988 at the third *TCP/IP Interoperability Conference* held in Santa Clara, California. This story tells how it evolved and specifically how the larger context of this network became a powerful tool that moved the Internet from a mostly government-sponsored research project to a network that would support commercial businesses and could be used by millions of ordinary people in their daily lives. But before we consider what happened then, we must turn back the clock a couple of years.

In August 1986, a few very motivated people decided to teach others how to implement the early Internet protocols. This first conference, called the *TCP/IP Vendors Workshop*, was held in Monterey, California, and was by invitation only (See Figure 1).

Figure 1: TCP/IP Vendors Workshop Agenda, August 1986.

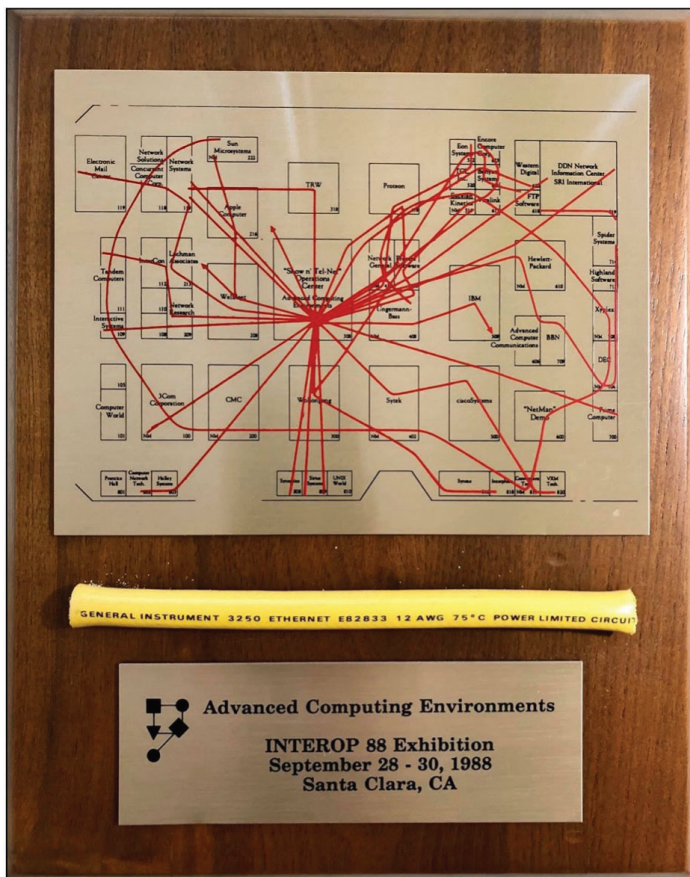
<p style="text-align: center;">TCP/IP Vendors Workshop Agenda Announcement</p> <p style="text-align: center;">25-27 August 1986 Doubletree Inn, Monterey, California</p> <p style="text-align: center;">Workshop Chair, Dan Lynch</p> <p>Sunday, August 24 5-7 PM -- Welcome Mixer at the Hotel</p> <p>Monday, August 25</p> <p>8:30 -- Coffee and Danish</p> <p>9:00 -- Welcome and Workshop Overview -- Dan Lynch, ACE</p> <p>9:10 -- DoD and ISO Protocol Coexistence Plan -- Mike Corrigan, OSD</p> <p>9:30 -- Future Directions for the Internet -- Jon Postel, USC-ISI</p> <p>9:45 -- NSFnet - Status and Plans -- Steve Wolff, NSF</p> <p>10:00 -- Certification Activity Status -- Bob Jones, SDC</p> <p>10:15 -- OSI Gateway Project -- John Heafner, NBS</p> <p>10:30 -- Break</p> <p>11:00 -- ULANA Specification -- Stan Ames, Mitre</p> <p>11:15 -- How to get information -- Elizabeth Feinler, SRI International</p> <p>11:30 -- Specification and Testing -- Dan Lynch (Testing and Certification, IMP (PSN) port availability)</p> <p>12:00 -- Lunch -- Find a place to eat nearby. -- Easy task.</p> <p>1:30 -- Network -- Vanilla Style -- Carl Sunshine, SDC (IP, ICMP, Net Management, Fragmentation, Security, Precedence (TOS))</p> <p>3:00 -- Break</p> <p>3:30 -- Network -- With Nuts -- Dave Mills, Linkabit (M/A-COM) Routing, ARP, RARP, BOOTP, Address Mapping Issues, Gateways (EGP, GGP, IGP), Subnetting, ISDN</p> <p>5:30 -- Break</p> <p>7:00 -- Meet at Monterey Bay Aquarium for a Private Viewing, Hors d'Oeuvre, Dinner, and a talk: Don't follow me, I'm a lost Internet Datagram -- Vint Cerf, NRI</p>	<p style="text-align: center;">TCP/IP Vendors Workshop Agenda Announcement</p> <p style="text-align: center;">Page 2</p> <p>Tuesday, August 26</p> <p>8:30 -- Coffee and Danish</p> <p>9:00 -- Files -- Jeff Mogul, Stanford (FTP, TFTP, Network File Servers (NFS, RFS, TOPS))</p> <p>10:15 -- Break</p> <p>10:45 -- Domain Name System Architecture -- Jon Postel</p> <p>11:15 -- Telnet -- Jon Postel (Telnet, Telnet Options, NACs)</p> <p>12:00 -- Lunch -- Find a place nearby.</p> <p>1:30 -- Mail -- Paul Mockapetris, USC-ISI (SMTP, User Mail Agents, Multi Media Mail, X.400)</p> <p>3:15 -- Break</p> <p>3:45 -- Transport -- Vint Cerf (TCP, Retransmission, UDP, Netbios on TCP)</p> <p>5:30 -- Done for the day...</p> <p>Wednesday, August 27</p> <p>8:30 -- Coffee and Danish</p> <p>9:00 -- Transactions -- Bob Braden, USC-ISI (Domain Names Implementation, Remote Procedure Calls, Transaction Services)</p> <p>10:30 -- Break</p> <p>11:00 -- Data Link and Physical-- Steve Holmgren, CMC (1822, X.25, Ethernet, Token Ring, Token Bus, Satellite, Packet Radio, Appletalk)</p> <p>12:00 -- Lunch -- As usual, find a spot on the wharf...</p> <p>1:30 -- Data Link and Physical, continued from morning</p> <p>3:00 -- Wrapup -- Dan Lynch (Identification of pressing issues for vendors, plans for follow on meetings.)</p>
--	--

Speakers included Vint Cerf, who was at MCI at the time, Jon Postel, *Request For Comments* (RFC) Editor and at ARPANET before playing a key role in Internet administration, and Paul Mockapetris and Bob Braden, both at the *University of Southern California's Information Sciences Institute* (USC-ISI).

Two subsequent conferences were held the following March in Monterey and then December in Crystal City, Virginia, and both were called the *TCP/IP Interoperability Conference*. All three were unusual events for several reasons: first, the presenters and instructors were the actual engineers that developed the earliest Internet protocols. They were also there to impart knowledge, rather than sell products—mainly because few commercial products were yet invented. One of the instructors, Douglas Comer of Purdue University, wrote the first and best-selling book on the topic: *Internetworking with TCP/IP, Volume 1, Principles, Protocols and Architecture*.

By September of 1988, the format of the conference changed, and expanded beyond lectures to a more practical proving ground. The event was renamed once again, and so *Interop*—and its show network—was born. The mission was still to teach Internet technologies and protocols, but for the first time the event was used to test and demonstrate various Internet communications devices on an active computer network. That show used a variety of Ethernet cables to connect 51 exhibitors together, with T1 links to the NASA Ames Research Center in Mountain View, California, and the NSFNet in Ann Arbor, Michigan (Figure 2).

Figure 2: Left: Interop 88 Exhibition show network. Right: NETWORLD+INTEROP 1996 advertisement.



THIS EVENT COULD REALLY GET ROUTY.

And the routing part is just for starters. After all, this is the world's top interoperability event.

At NetWorld+Interop® '96 Las Vegas, you'll get a chance to work with more than 50 hard-core networking experts like Wej, Doug, Robin and Steve. Together with other members of our Network Operations Center (NOC) team, they're in charge of operating our 6000 node, multi-vendor, multi-protocol network proving ground—the InteropNet™.

The NOC team designs and integrates classical networks with cutting-edge technologies.

So you get to see ATM running with existing networks. Or test 100base-T vs. 100VG, routing vs. switching, client/server interoperability, Internet applications and more.

The InteropNet also connects more than 550 vendors with the latest networking solutions, like LAN emulation, Fibre Channel, 155Mbps Wireless, FDDI, 100VG-AnyLAN, ATM, frame relay transports...the list goes on and on.

Call today to get a FREE VIP Exhibition Pass! Or access the Web at <http://www.interop.com> to register and get complete event information.

NETWORLD+INTEROP

Exhibition is April 2-4 • Conference is April 1-5 • Las Vegas Convention Center • 800-488-2883

InteropNet is a registered trademark of InteropNet, Inc. © 1996 InteropNet, Inc. All other names are the property of their respective holders.

The Interop conference quickly grew into a worldwide series of events^[1] with multiple shows held in different cities that were attended by tens of thousands of visitors with more than a thousand connected booths. In those early years, the largest shows were held in Tokyo, which began in 1994 and continued annually (with a pause because of the pandemic), with the latest show held in mid-2024. This year's show spanned over 500 vendors' booths and drew about 40,000 visitors each of its three days. The Tokyo Interop is also where the *ShowNet* (this is the chosen capitalization for the Tokyo show) not only has survived, but also has thrived, and continues to innovate and demonstrate Internet interoperability to this day. Many products had their world or Japanese debuts at various Tokyo Interop events, including Cisco's XDR and 8608 Router and NTT's Open APN.

My Own Interop Journey with Network Computing Magazine

Before I discuss the evolution of Interop and the role and history of its show network, I should first mention my own personal journey with Interop. In 1990, I was in the process of creating the first issue of *Network Computing* magazine for CMP Media. Our first issue was going to debut at the Interop show, the second time it was held that year in the San Jose, California, Convention Center.

The publisher and I both thought this place was the best one for our debut for several reasons. First, our magazine was designed for similar motivations—to demonstrate what worked in the new field of computer networking. We had designed our publication around a series of laboratories that had the same equipment found in a typical corporate office, including wide-area links and a mixture of PC MS-DOS, Apple Macintosh, and Unix devices and even a DEC minicomputer connected together. Second, we wanted to make a “big splash,” and our salespeople were already showing prototype issues ahead of the show to entice advertisers to sign up. Finally, *Network Computing's* booth would be connected to the Shownet and the greater Internet, just like many of the exhibitors who were trying out some product for the very first time.

One other feature about *Network Computing* that set us apart from other business trade magazines at the time: each bylined article would contain the email address of the author, so that readers could contact them with questions and comments. I wanted to use the domain **cmp.com** and set up an actual Internet presence, but alas I was overruled by management, so we ended up using a gateway maintained by one of the departments at UCLA where a couple of our editors were housed. While posting an author's email contact is now common, it was a radical notion at the time.

That 1990 Interop conference began my own personal journey of numerous shows around the world through the years, including speaking and teaching, as well as covering them for various business and networking publications that I would write for.

The Earliest Days of Interop

The Interop show in the late 1980s was a markedly different trade show from others of its era. At the time, trade shows with networked booths were non-existent. By way of perspective, up until that point in those early years, there were two kinds of conferences: one focused on the trade show with high-priced show floors and fancy exhibits. There, exhibitors were forced to “pay to play,” meaning if they bought booth space, they could secure a speaking slot at the associated conference. The other was a more staid affair that was a gathering of the engineers and actual implementers. Interop was a notable early example bridging the two: it looked like a trade show but was more of a conference, all in the guise of getting better commercial products out into the marketplace. It helped that it had its roots in those early TCP/IP conferences.

“You could see the Internet in a room, thanks to the Shownet, with hundreds of nodes talking to each other. That was unique for its time,” said Carl Malamud. “The Shownet was the most complex Internet installation you could do at any moment of time.”

Malamud would play several key roles in the development of various early Internet-based projects, including running the first Internet-based radio station, and he was a Shownet volunteer in 1991. In addition, Interop commissioned him to write the book *Exploring the Internet: A Technical Travelogue*.^[6]

That complexity has been true from the moment the Shownet was first conceived to the present day. Many of the Internet protocols—both in their earliest years and up to the present era—were debugged over the Shownet: volunteers recall testing NetBIOS over TCP/IP, 10BaseT Ethernet, SNA over TCP/IP, *Fiber Distributed Data Interface* (FDDI), *Simple Network Management Protocol* (SNMP), *Internet Protocol Version 6* (IPv6), various versions of segment routing, and numerous others. That extensive protocol catalog is a testament to the influence and effectiveness of the Shownet, and how enduring a concept it has been over the course of Internet history. Steve Hultquist, who was part of the early *Network Operations Center* (NOC) teams, remembers that the first version of 3Com’s 100BaseT switch—with “serial number 1”—was installed on the show network.

The force behind Interop was Dan Lynch, who passed away earlier this year. Lynch foresaw the commercial Internet and designed Interop to hasten its adoption. He based Interop on a series of efforts to bring together TCP/IP vendors, and the proto-Interop shows that were run in the middle 1980s that were more “Plugfests” or “Connectathons,” where vendors would try out their products. The main difference was those efforts deployed mostly proprietary protocols, whereas Interop ran on open source.

He told Sharon Fisher in November 1987: “There are millions of PCs out there and they’re starting to get networked in meaningful ways, not just in little printer-sharing networks.”

Part of his vision—and those that he recruited—was the notion of interoperability that could be used as a selling point and as an alternative to single-vendor proprietary networks from IBM, Digital Equipment Corporation, and others that were common in that era. Larry Lang, who worked for many years at Cisco, said, “The reassurance that it was okay to give up having ‘one throat to choke’ came from confidence that the equipment was interoperable. It is hard to remember a time when that was a worry, but it surely was.”

Part of Lynch’s vision was to ensure that proving interoperability was a very simple litmus test: did the product being exhibited work as advertised in a real-world situation? The answer to that question seems like common sense, but doing so in a trade show context was a relatively rare idea. And while it was a simple question, the answer was usually anything but simple, and sometimes the reasons why some product didn’t work—or didn’t work all the time—was what made the Shownet a powerful product improvement tool. That is just as true today as it was back then. The more realistic the Shownet, the more often it would expose these special circumstances that would bring out the bugs and other implementation problems.

It would prove to be a potent and enduring vision.

What Does Interoperability Mean?

The notion of interoperability seems so common sense now—and indeed it is the default position for most of the current networking world. However, in the early days of the Internet it was fraught with problems in terms of both larger-scale implementations and smaller issues that would prevent products from working reliably. One of Sharon Fisher’s articles in *Computerworld* in 1991^[2] speaks about TCP/IP this way: “The astounding thing is not how gracefully it performs but that it performs at all. TCP/IP is not for everyone.” Times certainly have changed in the 33 years since that was written. Today the notion of Internet connectivity, using TCP/IP protocols, is a given assumption in any computing product—from smart watches to the largest mainframe computers.

Those early implementation differences plagued both large and small vendors alike, and required a meeting of the minds where the protocol specifications weren’t exacting enough to ensure its success, or where bugs took time to resolve. Enter the Interop conference. As a reminder, in those early days the popular IP applications were based on the *File Transfer Protocol* (FTP), the *Simple Mail Transfer Protocol* (SMTP), and the *Simple Network Management Protocol* (SNMP). The web was still being invented and far away from being the de facto smash hit that it is today. Video conferencing and streaming didn’t exist. Telephones still ran on non-IP networks.

One of the early casualties was the *Open Systems Interconnection* (OSI) series of protocols promulgated by the *International Organization for Standardization* (ISO).

It was precisely because of the interoperability among TCP/IP products and “the failure of OSI to effectively demonstrate interoperability in the early 1990s that was the final nail in its coffin,” said Brian Lloyd, who worked at Telebit at the time. There are other stories of the defeat of OSI, such as this one in the *IEEE Spectrum*.^[3]

The Relationship Among the Shownet, the Conference Tutorials, and the NOC

To accomplish Lynch’s vision, Interop was not only the Shownet, but also its interaction with two other elements that became force multipliers in the quest for interoperability. These elements were the tutorials that were given before the opening of the show floor, and the NOC team that ran the network itself. All three had an important synergy to promote the actual practice of interoperability among different vendors’ products: not only in the demonstration of what worked with what but also in the discovery of protocol mismatches or programming errors so that new equipment could be made to interoperate.

Dave Crocker, who authored many RFCs and served on the Interop program committee that selected speakers in the 1990s, called out this tripartite structure of Shownet, NOC, and conference as a major strength of Interop. “Interop was able to contrast the technologies of the Internet with the interoperability of non-Internet technologies, such as IBM’s *Systems Network Architecture* (SNA). It had very pragmatic implications and wasn’t just promoting marketing speak.”

Many of the engineers who developed those early protocols and techniques and other pieces of Internet technology taught the tutorials (and as I mentioned, I taught a few of them during those early years, in addition to serving with Crocker on the program committee), so that others could learn how to best implement them. Here is where Interop contained its secret sauce: the people who taught the tutorials were the people who contributed to the underlying protocols and code, in some cases code so new it was changing over the duration of the show itself. “It was only after getting to Interop that we found out how few options were actually used by most implementations, and only then did we have access to the larger Internet and various versions of Unix computers,” said Brian Lloyd. “It was real bleeding edge stuff back then and the place to go for product testing and see how marketers and engineers would work together.”

And the NOC was a real one, like what could be found at large corporations, monitoring the network for anomalies and using it to debug various implementations leading up to the opening moments of the show. “It was unusual for its time,” said Fisher in another article in *Infoworld*.^[4] “The NOC team was infamous in the trade press for its tours and the time members took to explain things to us,” she said. Malamud recalls that the NOC had a strict “no suits” policy, meaning that its denizens were engineers that rolled up their sleeves and got stuff done.

All of this happened with very few paid staff: most of the people behind the Shownet and NOC were volunteers who came back, show after show, to work on setting things up and then taking them down after the show ended. That was, and to some extent still is, a very high-pressure environment: imagine wiring up a large convention center and connecting all of its conference rooms with a variety of network cabling. Several of the original Interop Shownet and NOC volunteers are continuing the tradition by helping to build and run the *Internet Engineering Task Force* (IETF) event networks at every IETF meeting around the world.

One of the more infamous moments of Interop was the *Internet Toaster*, created originally by John Romkey for the 1990 show^[5].

“I wanted to get people thinking of SNMP not just as getting variables, but for control applications, a wider vision. So we had an SNMP controlled toaster. If you put bread in the toaster, and set a variable in SNMP, the toaster would start toasting. A whole *Management Information Base* (MIB) was written for it, including how you wanted the toast, and whether it was a bagel or Wonder Bread. I ended up with lots and lots of bread in my garage. It got a lot of attention, but I don’t think that managing your kitchen through SNMP is very practical today.”

Dave Buerger, who was an early tech journalist at CMP, remembers Interop as having “a strange sense of awe unfolding for everyone as we glimpsed the possibilities of global connectivity. Exhibits on the show floor were more experiments in connecting their booths to the rest of the world.”

Construction of the Earliest Shownets (before 1993)

To say that Lynch was very persuasive is perhaps a big understatement. He convinced people who were quirky, unruly, or difficult to work with to spend lots of time pulling things together. “Dan allowed us to do stuff that the usual convention wouldn’t normally allow, and managed people that weren’t used to being managed,” said Malamud.

Peter de Vries was one of the earliest Shownet builders when he was working at Wollongong as one of the early Internet vendors. He ended up working for Interop for three years before opening the West Coast office for FTP Software. He remembers Lynch “dragging people kicking and screaming into using the Internet” back in the late 1980s and early 1990s. “But he was a fun guy to work for, and he had an unusual management style where he didn’t issue demands but convinced you to do something through more subtle suggestions, so by the time you did it you were convinced you had the original idea.”

These volunteers would essentially be working year-round, especially once the calendar was filled with multiple shows per year. Back in those early years, the convention centers didn’t care about cabling, and hadn’t yet figured out that having a more permanent physical networking plant could be used as an asset for attracting future meetings.

“We were often the first show to hang cables from the ceiling, and it wasn’t easy to do,” said Malamud, who chronicled the 1991 Shownet assembly^[6]. The first Interop shows used thick Ethernet cables that required a great deal of finesse to work with; de Vries recalls they had to pass wires through expansion joints and other existing holes in the walls and floors, wires that didn’t easily bend around corners.

“Each network tap took at least ten minutes of careful drilling to attach to this thick cable.” He has many fond memories of Lynch: “My goal was to try to get everyone to use TCP/IP, but Dan took it to the next step and showed that TCP could be a useful tool, something better than a fax. He was a real visionary.”

Ethernet—in all of its variations over the years, including early implementations of 10BaseT and 100-megabit speeds—wasn’t the only cabling choice for Interop; the show would expand to fiber and *Token Ring* cabling as part of its mission. Brian Chee was one of those volunteers who remembers having to re-terminate 150 different fiber strands across the high catwalks of the convention spaces. “We even had to terminate the fiber on the roof of the convention center to connect it to the Las Vegas Hilton across the street,” he said.

Getting all that cabling up in the air wasn’t an easy task either. Patrick Mahan worked on the San Francisco show in 1992 and recalls that he and other networking volunteers were paired with the union electrical workers on a series of scissor lifts. “We needed multiple lifts operating in tandem to raise them, and we had just started hanging a 100-foot length [cable] when a loud air-horn goes off and each lift immediately starts descending to the floor, because it was time for the mandatory union 15-minute break. It took about three minutes before the cable bundles started breaking apart and crashing to the cement floor. You could hear the glass in the fiber cables breaking!”

The Vegas climate made installations difficult, especially when its non-air-conditioned convention halls reached temperatures of 110 degrees Fahrenheit outside. “Many convention centers don’t turn on their air conditioning until the night before the show begins, so it was a particularly harsh environment,” said Glenn Evans, who worked as both a volunteer and an employee of Interop during the late 1990s and early 2000s. “Vegas in May is very dry, and static electricity is a big issue. We fried several switch ports inadvertently and spent long nights adding static filters to avoid it because some shows had more than 1,200 connections across their networks.” Evans emphasized that the install teams relied on “redneck engineering to come up with creative solutions, and it didn’t have to be perfect, [it] just had to work for five days.”

The cabling had to be laid out three times for each Shownet. The volunteers would have access to the convention center for a day months before any actual show. They would lay out the first cable segments and add connectors, then roll them up and store them in a warehouse. Then before the show there would be a “hot staging” event where the cables were connected to their equipment racks and tested.

Finally, several nights before the show began saw the real deployment at the convention center, which would span several 24-hour days before the actual opening. Those long nights were epic: de Vries recalls falling asleep in the middle of one night at the top of a 15-foot ladder, only to be gently awakened by the only other person in the convention center at the time. “Those installations nearly killed me!”

Many of the participants during those early years were motivated by a sense of common purpose, that their efforts were directly contributing to the Internet and its usefulness. “I loved that we could help the overall industry get stuff right,” said Hultquist. “They were some of the smartest people that I have ever worked with and were constantly pushing the envelope to try to deploy all sorts of emerging technologies.”

But the physical plant was just one issue; once the cabling was in place, the real world of getting equipment up and running across these networks was challenging. In those early Interops, equipment was often at the cutting edge, and engineers would make daily or even minute-by-minute changes to their protocol stacks and application code.

James van Bokkelen was the president of FTP Software then, and he recalls seeing the Shownet in 1988 crash while running BSD v4.3, thanks to a buggy version of one TCP/IP command. Turns out the bug was present in Cisco’s routers that were used on the Shownet. “It took a few minutes of scampering before everything was in place and we got Shownet back online,” he said. Scampering indeed: the volunteers had to compare notes, debug their code, and reboot equipment often located at different ends of the convention floor.

“We were getting alpha software releases during the show. This network created an environment where people had to fix things in real time in real production environments,” said Hultquist. “Wellfleet, 3Com, and Cisco were all sending us router firmware updates so their gear could interoperate with each other. I loved that we could help the overall industry get stuff right.”

At one of the 1991 Interop events, “FDDI completely melted down,” said Merike Kaeo, who at the time was working for Cisco in charge of their booth and volunteering in the NOC. “There was some obscure bug where a router reboot wasn’t enough, you had to reset the FDDI interface adapter separately. It didn’t take all that long to get things running, thankfully.”

Some of the problems were far more mundane, such as using equipment with NiCad batteries that had very short shelf life. Chee recalls that one Fluke engineering director got tired of trying to get these batteries replaced with Lithium-ion batteries. “He would send his team up to the rafters with network test equipment that had very short battery life; they were quickly replaced in their newer products.”

As more Shownets were brought up over the years, they had built-in redundant—and segregated—links. “We all played a part to make sure that after 1991, we would have a stable portion that would run reliably and put any untested equipment on another network that wouldn’t bring the working network down after the show started,” said Bill Kelly, who worked for Cisco in its early days.

de Vries said the first couple of Interop Shownets had less than 10 miles of cabling, which grew by 1991, according to Malamud, to having more than 35 miles of cabling, connecting a series of ribs, each one running down an aisle of the convention floor or some other well-defined geographic area. “Each rib had both Ethernet and Token Ring connected to an equipment rack with various routers,” said Malamud.

“There were two backbones that connected 50 different subnets, one based on FDDI and the other on Ethernet, which in turn were connected via T-1 lines to NASA Ames Research Center and Bay Area Internet points.”

Bill Kelly, who worked on the Shownet NOC while he was at Cisco, developed a three-stage model that covered a product lifecycle. “The first stage is using the IETF RFCs to try to make something work. Then the second stage is when a vendor is late to market and must figure out how to play nicely with the incumbents and the standards. The third stage is mostly commodity products, and everything works as advertised.”

The Middle Years (1994–1999)

The Internet—and Interop—were both growing quickly during this time. New Internet protocols and RFCs were being created frequently, and applications—and dot-com businesses—sprang up without any business plans, let alone initial paying customers. There were new venues each year in Europe, shows in Sydney, Australia, and Singapore, and Sao Paulo, Brazil. Some years had as many as seven or eight different shows, each with its own Shownet that needed to be customized for the exhibit halls in these cities.

Let’s return to 1986 for a moment. That year Novell began its own trade shows, called *NetWorld*, to explain its growing *Netware* community. By 1994, these shows had grown, and that is when Novell and Interop merged their shows, calling them *NetWorld+Interop*. This moniker held until 2004, when a series of technical media companies purchased Interop.

“Shownet didn’t change much after the Novell merger. We could accommodate their stuff at the edge, but it didn’t impact the core network,” said Hultquist. For the Shownet team, Netware was just another protocol to interoperate across. Despite Novell’s influence, during these years, TCP/IP became a networking standard. So did the cabling that made up the Shownet: “In the mid-1990s, a lot of the cable plant could be reused from show to show, with a standard set of 29-strand multimodal fiber with quick connectors and 48 strands of Category 5 copper cable for the ribs,” said Evans.

TCP/IP evolved too: by the end of the 1990s, the protocols and Ethernet hardware became commodities and were both factory-installed in millions of endpoint devices. “The Internet was becoming more standardized, and Interop became less of an experiment and more of a technology demonstration,” said Evans.

Nevertheless, vendors tried to differentiate themselves with quirky exhibits, pushing the envelope of connectivity. One stunt happened during the 1995 Interop at the Broadcom booth, which demonstrated Ethernet signals over barbed wire. “The wires were ugly and rusty and had nasty little barbs all over them,” according to one description written years later.^[6,11]

By 1999, the Shownet split into two separate parts: the live production network connecting the exhibit booths called *InteropNet*, and *InteropNet Labs* used for showcasing new technologies and products. Back then, these new technologies included VoIP, VPNs, and other “hot technologies,” according to a post by Tim Greene on CNN.^[7] Several market forces caused this situation. First, more and more conferences began promoting the idea of Internet connectivity for both attendees and vendor participants. “As that reality dawned on people, the Interop Shownet became an increasingly useless anachronism,” said Larry Lang, who was part of the team building Cisco’s support for FDDI at that time. “As our competition became Wellfleet rather than IBM, why would we want to participate in an expensive and time-consuming display that suggested complete equivalence among all the products?”

Hultquist was quoted in that CNN piece saying that attendees “won’t know whether a piece of equipment really worked because of the demands placed on them by more experimental or untested products.”

A second issue had to do with striking a balance between established vendors and newcomers. Kelly remembers the relationship between Cisco and Interop to be “complicated because we were the market leader and if we just donated equipment without any technical support, we ran the risk of outsiders misconfiguring the devices. Interop was also used to dealing with small engineering groups and not pesky marketing types that wanted to know the value of participating in the show.” Plus, long-running contributors to the original Shownets often got a jump on developing new gear and interacting with products that weren’t yet on the market.

By the end of the 1990s, the Shownet staging operation had also split into two. Prior to that moment, each Shownet would be staged in a Silicon Valley warehouse. But then the show runners for Tokyo decided to set up their own facilities to stage and constitute their own Shownet, and reformulated their NOC team from local talent, where they continue to build and demonstrate interoperability to the present day.

The New Millennium of Interop

Interop continued to grow in the new millennium. Two notable events affected the Shownet. The day the towers fell in New York in 2001 was also the day that the fall Interop show in Atlanta started. Many of the Shownet volunteers recall how quickly their network became the main delivery of news and video feeds to those attendees who were stuck in Atlanta, since all flights were grounded for the next several days. Brian Chee remembers that “within minutes of the disaster we maxed out the twin OC-12 WAN connections into the Shownet. We brought up streaming video of *CNN Headline News* over IP multicast, and that cut our wide-area traffic substantially, while at the same time it was an impressive demonstration of that technology.”

But then a few years later another event happened. “The day the Slammer virus hit, in 2003, we had just gone into production across the Shownet. That virus hurt our network throughput just enough that all our monitoring devices were useless,” said Chee. “But the NOC team was able to characterize the problem within a few minutes, and we were able to use air gapped consoles to reset routers and filter out the virus-infected packets.” That is as real world as it gets and is an example of how the Shownet proved its worth, time and time again.

But what is amazing is how enduring the legacy of Shownet continues to be. For example, during the 2019 Tokyo Interop, it played a critical role in demonstrating the interoperability among various segment routing vendors running over IPv6, resulting in a draft Internet document.^[8] I had an opportunity to review a draft report from the Tokyo team about the 2024 network that will be published in a future IPJ issue. “We faced varied challenges and considerations to achieve this while serving user traffic,” they wrote^[9].

Subsequent Tokyo shows—indeed the now sole survivors of the Interop legacy—would continue to draw on a huge talent pool of local talent. This year’s show had more than 650 volunteer engineers, including 30 alone to operate its NOC. “In 2024, we had 11 working groups leading the following fields: facilities, optical transport, external connectivity, backbone network, data center and cloud, wireless network, monitoring, security, testers, 5G, and media over IP,” said Takashi Tomine, who was part of the NOC team. The NOC occupies an impressive amount of show-floor real estate, where it continues to serve as a teaching and demonstration tool, as well as a working network nerve center.

It also is an opportunity for university students and junior staff to obtain hands-on experience in network operations and spend two weeks touching technology in ways that they might not have in their jobs or classrooms. The NOC team conducts walking tours, wherein guides describe what these teams have done in the many Interops held elsewhere down through the years. The 2024 show endures in another way, the “hot staging” model that was developed more than 30 years ago at the first Silicon Valley shows. The team has a total of eight days to assemble the network, and a few hours to take it apart after the show ends.

“First, we install every device in the right place on the racks, turn on the devices, and check their status. Checking device statuses is very important because some devices are transported directly from overseas to the venue, so it is necessary to ensure that they are not malfunctioning. We usually finish this process on the first day. On the second day, we start the network setup,” said Tomine.

His article will be published in a subsequent IPJ issue that goes into further details about the Tokyo show and how it grew over the years.

Figure 3: Interop Tokyo
2024 ShowNet Walking Tour



The 2024 Tokyo Interop showcased several new technologies, or technologies used in new and innovative ways. For example, the Shownet shared streaming video content with three geographically distributed TV broadcasting stations, all over IP networks. The team built a special media operations center to control these broadcasts and to demonstrate real-time video recording and editing of several conference sessions and demonstrations. In that respect, it was back to the future when the first multicast IP streams were broadcast years ago.

“The coolest thing we got out of working at Interop is that technology doesn’t happen without the people, and the people involved were some of the hardest-working and smartest people that you’ll ever meet. They checked their egos at the door, and solved problems jointly,” said Evans. “It was run like a democratic dictatorship, where everyone had a say.”

Acknowledgements

This article wouldn’t be possible without the help of numerous volunteers and staff at Interop events of the past, including Bill Alderson, Karl Auerbach, Dave Buerger, Brian Chee, Dave Crocker, Peter J.L. de Vries, Glenn Evans, Sharon Fisher, Connie Fleenor, Steve Hultquist, Merike Kaeo, Bill Kelly, Larry Lang, Brian Lloyd, Carl Malamud, Jim Martin, Naoki Matsuhira, Ryota Roy Motobayashi, Takashi Tomine, and James van Bokkelen.



Dedication

This article is dedicated to the memory of Daniel Courtney Lynch, August 16, 1941 – March 30, 2024, founder of the Interop events and whose vision gave us the Interop Shownet.

References and Further Reading

- [1] Ryota Motobayashi, “Report on INTEROP,”
<http://motobayashi.net/interop/>
- [2] Sharon Fisher, “TCP/IP: Effective, yes. Well supported, definitely. But if you’re looking for a graceful way to link PCs to hosts, you’d better wait a while,” *ComputerWorld*, October 7, 1991.
- [3] Andrew L. Russell, OSI: The Internet That Wasn’t,” *IEEE Spectrum*, July 29, 2013.
<https://spectrum.ieee.org/osi-the-internet-that-wasnt>
- [4] Sharon Fisher, “Products Help Users Manage TCP/IP-Based Networks,” *Infoworld*, October 10, 1988.
- [5] John Romkey, as told to Bernard Aboba, “How PC-IP Came to Be,” *The Internet Archive*, August 7, 2011.
- [6] Carl Malamud, “San Jose,” a chapter from *Exploring the Internet: A Technical Travelogue*, ISBN-13, 978-0132968980, January 1992.
- [7] Tim Greene, “Bleeding edge tech booted from show net,” CNN, May 11, 1999.
- [8] Ryo Nakamura, Yukito Ueno, and Teppei Kamata, “An experiment of SRv6 Service Chaining at Interop Tokyo 2019 Shownet,” Internet-Draft, Work in Progress, October 30, 2019.
- [9] Takashi Tomine, Ryo Nakamura, and Ryota Motobayashi, “ShowNet at Interop Tokyo: A Continuously Evolving Demonstration Network,” to be published in a future edition of *The Internet Protocol Journal*.
- [10] Karl Auerbach, “Our Scrapbook (The Ghost of Interop Past)”:
<https://cavebear.com/archive/interop/>
- [11] Dr. Howard Johnson, “So Good, It Works on Barbed Wire,” *EDN Magazine*, July 5, 2001.
- [12] “Building the INTEROPnet 93,” YouTube video.

DAVID STROM has been writing about technology since 1986 for various print and online publications, including *PC Week/eWeek*, *ComputerWorld*, *Infoworld*, *eeTimes*, *Tom’s Hardware*, *CSOonline*, *SiliconANGLE.com*, *Dark-Reading.com*, and dozens of others. He last wrote in IPJ Volume 23, No. 2 about his experiences selling his own class C block of IPv4 addresses back in 2020. He was the founding editor-in-chief of *Network Computing*, and he helped build numerous technical and editorial websites focusing on networking, security, and communications. He has written two books, *Internet Messaging* (1998, Prentice Hall) with Marshall T. Rose and *Home Networking Survival Guide* (2001, Osborne/McGraw Hill). He has two daughters and lives in St. Louis and can be found online at strom.com

Experimental IPv6-only Network at APRICOT 2024

by Brian Candler, NSRC

The Asia Pacific Regional Internet Conference on Operational Technologies (APRICOT)^[0] is the Internet Network Operators Summit for the Asia Pacific region. Every year, a wireless network conference is deployed to provide connectivity for hundreds of delegates, with a separate “IPv6-only” SSID as an alternative for people to try. This year, we decided to experiment with a new approach, by using some of the recent mechanisms designed for “IPv6-mostly” networks to build a better “IPv6-only” network.

What Is IPv6-mostly?

“IPv6-mostly” is a way to gracefully sunset *Internet Protocol Version 4* (IPv4) on dual-stack access networks. The work in this area has been driven in part by Google’s enterprise network, which has become so large that they ran out of RFC 1918^[1] private addresses. Technically, there are two pieces to this plan:

- A new “IPv6-only preferred” option for *Dynamic Host Configuration Protocol Version 4* (DHCPv4) (option 108, RFC 8925)^[2]. By requesting this option, a client declares that it is willing to run in a single-stack, IPv6-only mode. And by returning this option, the DHCPv4 server confirms that the network is happy to work this way too. The client then doesn’t configure itself with any IPv4 address.
- A new “PREF64” *Neighbor Discovery Option for Router Advertisements* (RFC 8781)^[3]. This option tells the client that a *Network Address Translator 64* (NAT64) is available, and what prefix to use.

If both of these conditions are true, the client configures itself with a *Customer-Side NAT46 Translator* (CLAT), with a hidden private IPv4 address. Any IPv4 application traffic is routed through this translator and carried across the network as IPv6 until it reaches the *Provider-Side NAT64* (PLAT), where it is converted back to IPv4. This whole mechanism is called *464XLAT*.

The end result is that you can interact with IPv4 resources—even using IPv4 literals, like `ping 8.8.8.8`—when running on an IPv6-only network. In effect, your IPv6 network doubles as a large block of private addresses behind a NAT. A big advantage of this approach is that there is no need to use DNS64 to generate fake AAAA records for IPv4-only destinations.

“IPv6-mostly” is supported by modern versions of macOS (13+), iOS, and Android. Any other clients will simply continue with regular dual-stack operation, but overall the usage of your DHCPv4 address pools will decrease.

Using IPv6-mostly Features for IPv6-only

For APRICOT, we wanted to build a pure IPv6-only network, not dual-stack. But we also wanted to enable the CLAT in client devices that support it in order to get maximum compatibility with IPv4. Here is how we enabled it: the pieces were all built inside an Ubuntu 22.04 *Virtual Machine* (VM) running on a compact *Next Unit of Computing* (NUC) computer. The NUC is a line of small-form-factor barebone computer kits designed by Intel.

First, we needed a DHCPv4 server that would respond to clients that requested option 108, granting them permission to run IPv6-only. Regular DHCP servers like *ISC DHCP* and *Kea DHCP* are quite happy to do that. However, we also did *not* want to respond to clients who *didn't* support option 108; if we did, we'd have to offer them an IPv4 address, and we'd be back to a dual-stack network.

I couldn't find an off-the-shelf DHCPv4 server that was capable of working this way, so I found a modular DHCP server in Go called *coredhcp*^[4] and created a new plugin^[5] to implement the desired behavior. This plugin has now been merged into the main codebase.

Second, I needed to send router advertisements with the PREF64 option. The conference routers were Arista Layer 3 switches, and although they have this feature in very recent firmware, it wasn't available in the version we were using.

Therefore, I used Linux's *radvd*^[6] to perform the router advertisements. This feature was not available in the latest released version, only git HEAD, so I had to compile *radvd* from source. Since it's not possible for one router to send advertisements on behalf of another, it meant that the VM where *radvd* was running also had to act as the gateway for the IPv6-only network, turning the VM into a router for IPv6 traffic.

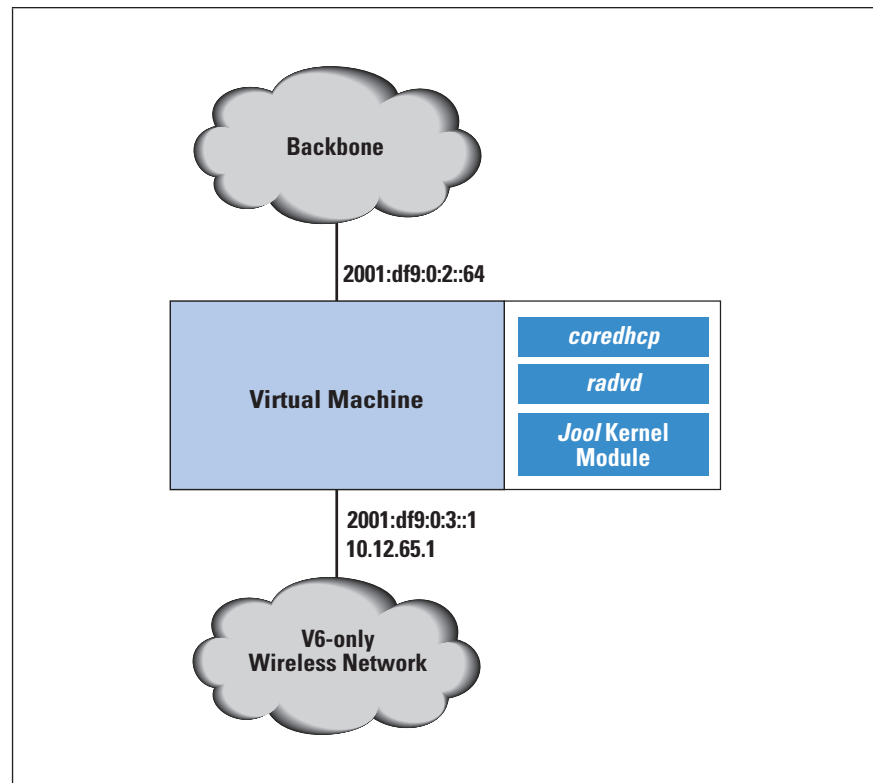
And finally, I needed a NAT64 translator somewhere on the network. This process could potentially have been done on a conference router, but since the IPv6 traffic had to pass through the VM anyway, I decided to implement NAT64 on the VM as well. I implemented it using *Jool*,^[7] a kernel module for NAT64 (refer to Figure 1 on the following page).

Configuration

We bundled all the components into a single VM running Ubuntu Linux 22.04. The VM acted as the router for the IPv6-only network, forwarding IPv6 packets in the kernel. It had an IPv6 address on the conference backbone, and an IPv6 address on a separate /64 subnet for the wireless network, with a static route on the upstream conference router.

We also configured the VM with an IPv4 address on the wireless-facing interface, but only so that the DHCPv4 server had an address to bind to. The VM offered no IPv4 addresses to clients, and enabled no IPv4 routing or NAT.

Figure 1: APRICOT 2024 Network Configuration for “IPv6 Mostly”



The other components running on the VM follow:

- The *Jool* kernel module, to perform the NAT64 translation;
- The *radvd* daemon, to announce the wireless prefix and the NAT64 prefix to clients; and
- The *coredhcp* server, to return DHCPv4 option 108 to those clients that requested it. It was also configured to act as a stateless DHCPv6 server, to give out DNS server addresses and domain search lists for any clients that didn't support Router Advertisement options *Recursive DNS Server Option* (RDNSS) and *Domain Name System Search List* (DNSSL).

Problems

The basic reachability to the IPv4 Internet via NAT64 seemed to work well. However, we were plagued by clients being repeatedly removed from the IPv6-only network by the Cisco *Wireless LAN Controller* (WLC), at seemingly random intervals. Strangely, this didn't seem to affect clients on the main conference SSID on the same access points. It turns out that several underlying issues were to blame.

First, the client private CLAT address **192.0.0.2** was leaking out as the source IP address on various multicast packets (such as multicast DNS, and Chrome doing service discovery on UDP port 1900). Logs showed that the wireless controller had a feature called “IP Theft or Reuse” that would add client MAC addresses to an exclusion list if it saw the same source IP address from multiple clients. We were able to turn that option off.

Second, some clients were being kicked off regularly every 10 minutes. The Cisco WLC had a setting where the wireless access points would perform dynamic channel reassignment every 10 minutes. We increased it to 24 hours.

Finally, when the conference was nearly over, we discovered another per-SSID setting, “IPv4 DHCP required,” which we also turned off. We believe that fixed the remaining problems.

There was one other major issue: clients would lose the ability to reach IPv4 destinations for a few minutes at a time, without being kicked off the wireless network—IPv6 connectivity continued to work. Using *tcpdump*, we saw that IPv6 neighbor discovery on the VM was forgetting about the IPv6 address of the CLAT. It turns out that the version of Jool in the Ubuntu 22.04 package repositories is old (v4.1.7), and this problem is known^[8]. It was straightforward to upgrade the package to the latest release, v4.1.11.

Apart from these problems, the CLAT/NAT64 mechanism worked very well for those devices that supported it. Remaining issues were minor: *traceroute* to an IPv4 destination showed only “*” for every hop, and the macOS ssh client didn’t work when given the “-4” flag (although it did work with an IPv4 literal address). Otherwise, it was just like being on a dual-stack network.

Conference Usage and Compatibility

From DHCP logs, I found that 142 unique devices had attempted to use the IPv6-only SSID, and of those, 115 (81%) supported DHCP option 108. That’s a surprisingly high proportion, representing a high usage of Apple laptops, iOS phones, and Android phones amongst delegates. Those devices should have gotten a good experience from the network, if it weren’t for the wireless disconnection issues.

The other 27 devices would have had a much worse experience. They got no DHCPv4 response, so they retried repeatedly, configured themselves with an IPv4 link-local address (**169.254.x.x**), and still would have been able to reach only Internet sites with IPv6 addresses.

We could have improved compatibility for these clients somewhat by providing a DNS64 service, which fakes AAAA records for DNS names that have only A records. However, these DNS settings would have applied to all hosts on the network, meaning that even those clients supporting option 108 would also have been exposed to fake DNS responses. I thought that the experiment was more useful without it, because the NAT64/DNS64 combination is already well-known and tested.

Post-conference Updates

Since the conference, I learned that it may no longer be necessary to provide any DHCPv4 server. Recent versions of MacOS will enable the CLAT even if the device has no IPv4 address at all.

However, an absent DHCPv4 server causes clients to keep sending DHCPDISCOVER messages, generating unnecessary load on both the clients and the network. Returning DHCPv4 option 108 stops the clients from doing this constant resending. There is also an older standard, “Auto-Config” DHCP option 116, in RFC 2563^[9], but DHCP logs from the IPv6-only conference network showed no clients using this option, so it appears to be obsolete.

I also discovered a problem about the choice of the NAT64 prefix. A *Well-Known Prefix* (WKP), **64:ff9b::/64**, is defined in RFC 6052^[10]. But if you use it, you will find that clients will be unable to connect via NAT64 to private addresses such as RFC 1918, because that is mandated by RFC 6052 section 3.1. If you want to use NAT64 on a typical home or enterprise network, and still be able to reach internal devices on private addresses, you will need to avoid the WKP. The conference network used a *Unique Local Address* (ULA) prefix (**fd64::/64**) instead.

Conclusion

IPv6-only using the IPv6-mostly mechanisms works surprisingly well, and is only going to improve over time as Windows^[11] and Linux add support for it.

Personally, I’d be quite happy to run this way at home, except that my Mikrotik router has no NAT64 capability. (RouterOS versions 7.8 and later do have the PREF64 router advertisement option though)^[12].

References and Further Reading

- [0] Asia Pacific Regional Internet Conference on Operational Technologies: <https://apricot.net/>
- [1] Yakov Rekhter, Robert G. Moskowitz, Daniel Karrenberg, Geert Jan de Groot, and Eliot Lear, “Address Allocation for Private Internets,” RFC 1918, February 1996.
- [2] Lorenzo Colitti, Jen Linkova, Michael C. Richardson, and Tomek Mrugalski, “IPv6-Only Preferred Option for DHCPv4,” RFC 8925, October 2020.
- [3] Lorenzo Colitti and Jen Linkova, “Discovering PREF64 in Router Advertisements,” RFC 8781, April 2020.
- [4] *coredhcp* server: <https://github.com/coredhcp/coredhcp>
- [5] Plugin for *coredhcp*:
<https://github.com/coredhcp/coredhcp/pull/1>. Option to Disable Stateless Auto-Configuration in IPv4 Clients,” RFC 2563, May 1999.
- [6] *radvd*: <https://github.com/radvd-project/radvd>
- [7] *Jool* kernel module: <https://github.com/NICMx/Jool>
- [8] “Jool when translating always drops the packet due to an error,” <https://github.com/NICMx/Jool/issues/382>
- [9] Ryan Troll, “DHCP Option to Disable Stateless Auto-Configuration in IPv4 Clients,” RFC 2563, May 1999.

- [10] Congxiao Bao, Christian Huitema, Marcelo Bagnulo, Mohamed Boucadair, and Xing Li, “IPv6 Addressing of IPv4/IPv6 Translators,” RFC 6052, October 2010.
- [11] Tommy Jensen, “Windows 11 Plans to Expand CLAT Support,” *Microsoft Tech Community Networking Blog*, March 7, 2024.
- [12] Mikrotik change log.
- [13] The video of my presentation at APRICOT 2024.
- [14] The slides from my presentation at APRICOT 2024.
- [15] Ondřej Caletka, “Deploying IPv6-mostly access networks,” RIPE Labs.
- [16] Jen Linkova, “Mission Possible: Turning off IPv4 in Google Enterprise Network,” presented at RIPE 87,

BRIAN CANDLER is a graduate of the University of Cambridge, UK. Having originally started in microelectronics design, he has developed systems and networking at several large UK-based ISPs. He delivers training courses for the Network Startup Resource Center in developing regions of the world, covering topics such as campus network design and network monitoring and management, and develops a virtual training platform used to provide hands-on lab exercises. He can be reached at: brian@nsrc.org

Brian Candler at APRICOT 2024



Check your Subscription Details!

Make sure that both your postal and e-mail addresses are up-to-date since these are the only methods by which we can contact you. If you see the words “Invalid E-mail” on your printed copy this means that we have been unable to contact you through the e-mail address on file. If this is the case, please contact us at ipj@protocoljournal.org with your new information. The subscription portal is located here: <https://www.ipjsubscription.org/>

IAB Workshop on Barriers to Internet Access of Services (BIAS)

The *Internet Architecture Board* (IAB) organizes workshops about topics of interest to the community that bring diverse experts together, raise awareness, and possibly identify the next steps that can be explored by the community. The IAB held its *Barriers for Internet Access of Services* (BIAS) fully online workshop during the week of January 15, 2024.

The Internet is a crucial component of our critical infrastructure that wields a significant influence on various aspects of society. It serves as a vital tool for advancing the United Nations' *Sustainable Development Goals* (SDGs)^[1] and upholding human rights on a global scale. Thus the absence of meaningful access to digital infrastructure and services amounts to a form of disenfranchisement. The barriers to meaningful access to Internet-based services and applications are increasing, posing challenges that persist even when Internet connectivity is available, thereby resulting in unequal information and service access.

The workshop solicited position papers about barriers to accessing content and services on the Internet, for example, based on filtering, blocking as well as due to general inequality of technological capabilities, like device or protocol limitations. 19 position papers were submitted to the workshop of which 12 papers were selected for publication^[2]. Two invited talks were also presented based on published papers. There were 40 participants in the workshop over three days.

This marked my first IAB workshop since joining the board. I am delighted to have collaborated with Mirja Kühlewind, Mallory Knodel, Tommy Pauly, and Christopher A. Wood in organizing this event. The themes of censorship, circumvention techniques, and the digital divide have surfaced in various IAB discussions lately. Our goal for this workshop was to present reports, expert opinions, and ignite discussions on these topics. Through this experience, I gained valuable insights and strongly believe that the IETF community must remain mindful of these crucial issues when designing protocols. It is imperative to ensure that we create the most secure, user-friendly protocols for all Internet users.

This article provides a short overview of the workshop discussion. However, if you would like to learn more you can also check out the initial draft version of the IAB workshop report^[3], or watch the entire thing on YouTube^[4]. The workshop was organized into three main themes across three days based on the submitted papers.

Community Networks

Community Networks are self-organized networks which are wholly owned by the community and thus provide an alternative mechanism to bring connectivity and Internet services to those places that lack commercial interest. Discussion ranged from highlighting the need for measuring *Quality of Experience* (QoE) for Community Networks, to the potential role a *Content Delivery Network* (CDN) can play in Community Networks, to the role of satellite networks, and finally, to the vital role of the spectrum in this space.

Digital Divide

The digital divide refers to disparities in access to the Internet and services. It signifies the gap between those who have effective and meaningful access to digital technologies and those who do not. Discussion recognized three key aspects of the digital divide: differences between population demographics in the provision of online resources by governments, inequality in the use of multilingual domains and email addresses, and increased costs for end-user downloads of contemporary websites' sizes. There was a general recognition that there may be more technical aspects of the digital divide that were not presented.

Censorship

Censorship is the legal control or suppression of what can be accessed, published, or viewed on the Internet. This discussion focused on reports of censorship as observed during recent years in different parts of the world, as well as on the use of and expectation for censorship circumvention tools, mainly the use of secure VPN services. This included censorship reports from India and Russia, where censorship has changed significantly recently, highlighting the legal frameworks and court acts that put obligations on regional network providers to block traffic. Further, measurements to validate the blocking, as well as analyses of how blocking is implemented were also discussed.

Next Steps

The discussion highlighted the need for the technical community to address the management gaps and document best practices for Community Networks including listing of manageability considerations explicitly for Community Networks. Further, the need to build consensus on solutions that have the most significant impact in fostering digital inclusion and the need to further promote them was discussed. We need to continue to work towards enhancing our protocols ensuring user privacy, develop further protocols that enable more transparency on filtering and new VPN-like services. Further discussion of these topics could happen in the *Global Access to the Internet for All* (GAIA)^[5], *Human Rights Protocol Considerations* (HRPC)^[6], *Privacy Enhancements and Assessments* (PEARG)^[7], and *Measurement and Analysis for Protocols* (MAPRG)^[8] research groups, based on the relevance to each group.

—Druv Dhody, IAB Member
dd@dhruvdhody.com

References

- [1] United Nations Department of Economic and Social Affairs, “Sustainable Development, The 17 Goals.”
- [2] IAB Workshop on Barriers to Internet Access of Services (BIAS) (biasws) Materials.
- [3] IAB Barriers to Internet Access of Services (BIAS) Workshop Report.
- [4] YouTube recordings from BIAS Workshop.
- [5] Global Access to the Internet for All (GAIA) Research Group.
- [6] Human Rights Protocol Considerations Research Group.
- [7] Privacy Enhancements and Assessments Research Group
- [8] Measurement and Analysis for Protocols Research Group.

IAB Statement on the Risks of Attestation on the Open Internet

While attestation of client software and hardware is a useful tool for preventing abuse or fraud on the Internet, the use of such attestation as a barrier to access otherwise open protocols and services would negatively impact the evolution of the Internet as a whole.

Openness and the empowerment of end users are core values of the IETF. RFC 3935^[1], Section 4.1, explains this as part of the IETF's mission statement:

“We want the Internet to be useful for communities that share our commitment to openness and fairness. We embrace technical concepts such as decentralized control, edge-user empowerment and sharing of resources, because those concepts resonate with the core values of the IETF community. These concepts have little to do with the technology that's possible, and much to do with the technology that we choose to create.”

The Internet is built upon the idea that anyone who implements the appropriate standards should be able to interoperate on the Internet. Many of the core services that run on the Internet, such as email and the web, are designed to be openly accessible in this way. Adding client attestation into otherwise open systems can significantly reduce openness for the Internet broadly. A recent “Web Environment Integrity”^[2] proposal has highlighted this risk, although such models pose a risk beyond just the web.

Attestation of client software and hardware is distinct from user authentication. User authentication verifies the identity of a user or a credential associated with a user, and is compatible with any implementation that supports the correct form of authentication. In contrast, attestation of client software and hardware places explicit restrictions on the implementations that are allowed to participate in the protocol. For services that have intentionally restricted access, such client attestation (as described in *Remote ATtestation procedureS* (RATS), RFC 9334^[3]) is a valuable security measure, particularly when used in conjunction with user authentication. However, this approach is not appropriate for openly accessible services.

Allowing clients to use a variety of software as long as it is protocol-compliant is an essential part of the IETF development process and the openness of the Internet. Although customized or open-source software can also be used to circumvent client-side security measures, the continuing viability of open software is required for continued innovation. Restricting access via attestation of software or hardware would limit the development of new protocols and extensions to existing protocols, lock users into a limited ecosystem of applications, and hamper the ability to audit implementations, conduct measurements, or perform essential security research.

If client attestation signals are used in open services to mitigate fraud or abuse, they should be designed to only signal the authenticity of a user or client without imposing strict software or hardware requirements. They should also be designed such that attestation is not required, but has a clear backup behavior when attestation is not possible. IETF-based protocols such as *Privacy Pass* [RFC 9576] attempt to provide a protocol that can be deployed in ways that promote user privacy without exposing detailed identifiers about the client systems that are being used. Fundamentally, attesting specific properties about a networking client (for example, there is some human user involved in this interaction) maintains the openness of the Internet, whereas attesting that a specific piece of software is in use does not and should be avoided.

The IAB invites those in the industry and standards community working on client attestation in open services to engage with the relevant IETF working groups (in particular, Privacy Pass^[4] and RATS^[5]), and encourages those groups to focus on defining safe deployment models for attestation and abuse prevention that will not put the openness of the Internet at risk.

References

- [1] Harald Alvestrand, “A Mission Statement for the IETF,” RFC 3935, October 2004.
- [2] Web Environment Integrity API:
<https://github.com/explainers-by-googlers/Web-Environment-Integrity/>
- [3] Henk Birkholz, Dave Thaler, Michael Richardson, Ned Smith, and Wei Pan, “Remote ATtestation procedureS (RATS) Architecture,” RFC 9334, January 2023.
- [4] Privacy Pass Working Group:
<https://datatracker.ietf.org/wg/privacypass/about/>
- [5] Remote ATtestation ProcedureS Working Group:
<https://datatracker.ietf.org/wg/rats/about/>

Internet Community Encouraged to Submit Event Proposals for UA Day 2025

Event proposals are now being accepted for the third annual *Universal Acceptance* (UA) Day, to be held between 1 March and 30 May 2025. UA Day, held annually, is an opportunity to rally local, national, regional, and global communities and organizations to spread UA awareness and to promote UA adoption with key stakeholders.

UA is a technical best practice that ensures all valid domain names and email addresses, regardless of script, language or character length, can be equally used by all Internet-enabled applications, devices, and systems. Co-organized by the *Universal Acceptance Steering Group* (UASG) and the *Internet Corporation for Assigned Names and Numbers* (ICANN), UA Day 2025 will consist of various virtual and in-person events held by the UASG, ICANN, global partners, and regional and local organizations.

UA Day 2025 builds on the success of *UA Day 2024* and the inaugural *UA Day* in 2023. Together, the 2023 and 2024 events attracted approximately 15,000 participants worldwide across dozens of countries. These milestone events have helped mobilize technical and language communities, companies, governments and *Domain Name System* (DNS) industry stakeholders to champion UA on a global scale.

Those interested in organizing a UA Day event must complete the *UA Day Event Proposal Form*^[1] by 11 October 2024. ICANN will provide limited support for proposed UA Day events based on this group's recommendations. The following types of events are eligible for support:

- *UA Awareness*: Provide a high-level overview of UA and *Email Address Internationalization* (EAI), the benefits of being UA-ready, basic technical concepts related to UA and next steps for becoming UA-ready.
- *UA Technical Training*: Provide in-depth training on becoming EAI-ready for email system administrators and on becoming UA-ready for software developers.
- *UA Academic Curricula*: Work with academic faculty members and experts to integrate *Internationalized Domain Names* (IDNs) and UA-related topics into existing technical curricula and design a roadmap.
- *UA Adoption*: Conduct a UA adoption exercise and share challenges and solutions to becoming UA-ready, and document your experience. Please note that advance preparatory work and review are required in order to qualify for a UA Adoption event.
- *UA Regional Strategy*: Discuss appropriate mechanisms for promoting UA adoption at the local, regional, and national levels.

Proposals will be considered from all relevant organizations, including international, regional and local organizations, technology organizations and companies, open-source communities, standards bodies, email service providers, academia, industry groups, and language communities.

UA is considered a foundational requirement for the continued expansion of the Internet. Since 2009, the landscape for domain names has changed markedly—in overall number of *Top-Level Domain Names* (TLDs) available, TLD character length and scripts available. However, the checks used by many software applications to validate domain names and email addresses often use rules that do not fully support Universal Acceptance. Achieving UA ensures everybody has the ability to experience the full social and economic power of the Internet using their chosen domain name and email address that best aligns with their interests, business, culture, language, and script.

Questions can be directed to **UAProgram@icann.org**. A full UA Day event calendar will be published in due course. In the meantime, interact with the UASG on social media (X, Facebook and LinkedIn) using the hashtag **#Internet4All**.

The UASG is a community-led initiative that was formed in 2015 and funded by ICANN. It consists of volunteers from many companies, governments, and community groups. The UASG works to raise awareness of the importance of UA globally, provide free resources to organizations to help them become UA-ready, and measure the progress of UA adoption. To learn more, visit <https://uasg.tech/>

ICANN's mission is to help ensure a stable, secure, and unified global Internet. To reach another person on the Internet, you need to type an address—a name or a number—into your computer or other device. That address must be unique so computers know where to find each other. ICANN helps coordinate and support these unique identifiers across the world. ICANN was formed in 1998 as a nonprofit public benefit corporation with a community of participants from all over the world. To learn more, visit <https://icann.org>

References

- [1] UA Day Event Proposal Form:
<https://tinyurl.com/UADayForm>

Our Privacy Policy

The *General Data Protection Regulation* (GDPR) is a regulation for data protection and privacy for all individual citizens of the *European Union* (EU) and the *European Economic Area* (EEA). Its implementation in May 2018 led many organizations worldwide to post or update privacy statements regarding how they handle information collected in the course of business. Such statements tend to be long and include carefully crafted legal language. We realize that we may need to provide similar language on our website and in the printed edition, but until such a statement has been developed here is an explanation of how we use any information you have supplied relating to your subscription:

- The mailing list for *The Internet Protocol Journal* (IPJ) is entirely “opt in.” We never have and never will use mailing lists from other organizations for any purpose.
- You may unsubscribe at any time using our online subscription system or by contacting us via e-mail. We will honor any request to remove your name and contact information from our database.
- We will use your contact information only to communicate with you about your subscription; for example, to inform you that a new issue is available, that your subscription needs to be renewed, or that your printed copy has been returned to us as undeliverable by the postal authorities.
- We will never use your contact information for any other purpose or provide the subscription list to any third party other than for the purpose of distributing IPJ by post or by electronic means.
- If you make a donation in support of the journal, your name will be listed on our website and in print unless you tell us otherwise.

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Kjetil Aas	Ilia Bromberg	Geert Van Dijk	Greg Giessow	Curtis Johnson
Fabrizio Accatino	Lukasz Bromirski	David Dillow	John Gilbert	Richard Johnson
Michael Achola	Václav Brožík	Richard Dodsworth	Serge Van Ginderachter	Jim Johnston
Martin Adkins	Christophe Brun	Ernesto Doelling	Greg Goddard	Jose Enrique Diaz Jolly
Melchior Aelmans	Gareth Bryan	Michael Dolan	Tiago Goncalves	Jonatan Jonasson
Christopher Affleck	Ron Buchalski	Eugene Doroniuk	Ron Goodheart	Daniel Jones
Scott Aitken	Paul Buchanan	Michael Dragone	Octavio Alfageme	Gary Jones
Jacobus Akkerhuis	Stefan Buckmann	Joshua Dreier	Gorostiaga	Jerry Jones
Antonio Cuñat Alario	Caner Budakoglu	Lutz Drink	Barry Greene	Michael Jones
William Allaire	Darrell Budic	Aaron Dudek	Jeffrey Greene	Amar Joshi
Nicola Altan	BugWorks	Dmitriy Dudko	Richard Gregor	Javier Juan
Shane Amante	Scott Burleigh	Andrew Dul	Martijn Groenleer	David Jump
Marcelo do Amaral	Chad Burnham	Joan Marc Riera	Geert Jan de Groot	Anders Marius Jørgensen
Matteo D'Ambrosio	Randy Bush	Duocastella	Ólafur Guðmundsson	Merike Kao
Selva Anandavel	Colin Butcher	Pedro Duque	Christopher Guemez	Andrew Kaiser
Jens Andersson	Jon Harald Bøvre	Holger Durer	Gulf Coast Shots	Vladislav Kalinovsky
Danish Ansari	Olivier Cahagne	Karlheinz Dölger	Sheryll de Guzman	Naoki Kambe
Finn Arildsen	Antoine Camerlo	Mark Eanes	Rex Hale	Akbar Kara
Tim Armstrong	Tracy Camp	Andrew Edwards	Jason Hall	Christos Karayiannis
Richard Artes	Brian Candler	Peter Robert Egli	James Hamilton	Daniel Karrenberg
Michael Aschwanden	Fabio Caneparo	George Ehlers	Darow Han	David Kekar
David Atkins	Roberto Canonico	Peter Eisses	Handy Networks LLC	Stuart Kendrick
Jac Backus	David Cardwell	Torbjörn Eklöv	Stephen Hanna	Robert Kent
Jaime Badua	Richard Carrara	Y Ertur	Martin Hannigan	Thomas Kernen
Bent Bagger	John Cavanaugh	ERNW GmbH	John Hardin	Jithin Kesavan
Eric Baker	Lj Cemerar	ESdatCo	David Harper	Jubal Kessler
Fred Baker	Dave Chapman	Steve Esquivel	Edward Hauser	Shan Ali Khan
Santosh Balagopalan	Stefanos Charchalakias	Jay Etchings	David Hauweele	Nabeel Khatri
William Baltas	Molly Cheam	Mikhail Evstiounin	Marilyn Hay	Dae Young Kim
David Bandinelli	Christof Chen	Bill Fenner	Headcrafts SRLS	William W. H. Kimandu
A C Barber	Pierluigi Checchi	Paul Ferguson	Hidde van der Heide	John King
Benjamin Barkin-Wilkins	Greg Chisholm	Ricardo Ferreira	Johan Helsingius	Russell Kirk
Ryan Barnes	David Chosrova	Kent Fichtner	Robert Hinden	Gary Klesk
Feras Batainah	Marcin Cieslak	Ulrich N Fierz	Michael Hippert	Anthony Klopp
Michael Bazarewsky	Lauris Cikovskis	Armin Fisslthaler	Damien Holloway	Henry Kluge
David Belson	Brad Clark	Michael Fiumano	Alain Van Hoof	Michael Kluk
Richard Bennett	Narelle Clark	The Flirble Organisation	Edward Hotard	Andrew Koch
Matthew Best	Horst Clausen	Jean-Pierre Forcioli	Bill Huber	Ia Kochiashvili
Hidde Beumer	James Cliver	Gary Ford	Hagen Hultzs	Carsten Koempe
Pier Paolo Biagi	Guido Coenders	Susan Forney	Kauto Huopio	Richard Koene
Arturo Bianchi	Robert Collet	Christopher Forsyth	Asbjørn Højmark	Alexader Kogan
John Bigrow	Joseph Connolly	Andrew Fox	Kevin Iddles	Matthijs Koot
Orvar Ari Bjarnason	Steve Corbató	Craig Fox	Mika Ilvesmaki	Antonin Kral
Tyson Blanchard	Brian Courtney	Fausto Franceschini	Karsten Iwen	Robert Krejčí
Axel Boeger	Beth and Steve Crocker	Erik Fredriksson	Joseph Jackson	John Kristoff
Keith Bogart	Dave Crocker	Valerie Fronczak	David Jaffe	Terje Krogdahl
Mirko Bonadei	Kevin Croes	Tomislav Futivic	Ashford Jaggernaut	Bobby Krupczak
Roberto Bonalumi	John Curran	Laurence Gagliani	Thomas Jalkanen	Murray Kucherawy
Lolke Boonstra	André Danthine	Edward Gallagher	Jozef Janitor	Warren Kumari
Cente Cornelis Boot	Morgan Davis	Andrew Gallo	Martijn Jansen	George Kuo
Julie Bottorff Photography	Jeff Day	Chris Gamboni	John Jarvis	Dirk Kurfuerst
Gerry Boudreaux	Fernando Saldana Del	Xosé Bravo Garcia	Dennis Jennings	Mathias Körber
Leen de Braal	Castillo	Oswaldo Gazzaniga	Edward Jennings	Darrell Lack
Stephen Bradley	Rodolfo Delgado-Bueno	Kevin Gee	Aart Jochem	Andrew Lamb
Kevin Breit	Julien Dhallenne	Rodney Gehrke	Nils Johansson	Richard Lamb
Thomas Bridge	Freek Dijkstra	Radu Cristian Gheorghiu	Brian Johnson	Yan Landriault

Edwin Lang	Kevin Menezes	Alexander Peuchert	Philip Schneck	Lorin J Thompson
Sig Lange	Bart Jan Menkveld	David Phelan	James Schneider	Jerome Tissieres
Markus Langenmair	Sean Mentzer	Harald Pilz	Peter Schoo	Fabrizio Tivano
Fred Langham	Eduard Metz	Derrell Piper	Dan Schrenk	Peter Tomsu Fine Art
Tracy LaQuey Parker	William Mills	Rob Pirnie	Richard Schultz	Photography
Christian de Larrinaga	David Millsom	Jorge Ivan Pincay	Timothy Schwab	Joseph Toste
Alex Latzko	Desiree Miloshevic	Ponce	Roger Schwartz	Rey Tucker
Jose Antonio Lazaro	Joost van der Minnen	Marc Vives Piza	SeenThere	Sandro Tumini
Lazaro	Thomas Mino	Victoria Poncini	Scott Seifel	Angelo Turetta
Antonio Leding	Rob Minshall	Blahoslav Popela	Paul Selkirk	Brian William Turnbow
Rick van Leeuwen	Wijnand Modderman-	Andrew Potter	Andre Serralheiro	Michael Turzanski
Simon Leinen	Lenstra	Ian Potts	Yury Shefer	Phil Tweedie
Anton van der Leun	Mohammad Moghaddas	Eduard Llull Pou	Yaron Sheffer	Steve Ulrich
Robert Lewis	Charles Monson	Tim Pozar	Doron Shikmoni	Unitek Engineering AG
Christian Liberale	Andrea Montefusco	David Preston	Tj Shumway	John Urbanek
Martin Lillepuu	Fernando Montenegro	David Raistrick	Jeffrey Sicuranza	Martin Urwaleck
Roger Lindholm	Roberto Montoya	Priyan R Rajeevan	Thorsten Sideboard	Betsy Vanderpool
Link Light Networks	Joel Moore	Balaji Rajendran	Greipur Sigurdsson	Surendran Vangadasalam
Art de Llanos	Joseph Moran	Paul Rathbone	Fillipe Cajaiba da Silva	Ramnath Vasudha
Mike Lochocki	John More	William Rawlings	Andrew Simmons	Randy Veasley
Chris and Janet Lonvick	Maurizio Moroni	Mujtiba Raza Rizvi	Pradeep Singh	Philip Venables
Mario Lopez	Brian Mort	Bill Reid	Henry Sinnreich	Buddy Venne
Sergio Loreti	Soenke Mumm	Petr Rejhon	Geoff Sisson	Alejandro Vennera
Eric Louie	Tariq Mustafa	Robert Remenyi	John Sisson	Luca Ventura
Adam Loveless	Stuart Nadin	Rodrigo Ribeiro	Helge Skrivervik	Scott Vermillion
Josh Lowe	Michel Nakhla	Glenn Ricart	Terry Slattery	Tom Vest
Guillermo a Loyola	Mazdak Rajabi Nasab	Justin Richards	Darren Sleeth	Peter Villemoes
Hannes Lubich	Krishna Natarajan	Rafael Riera	Richard Smit	Vista Global Coaching
Dan Lynch	Naveen Nathan	Mark Risinger	Bob Smith	& Consulting
David MacDuffie	Darryl Newman	Fernando Robayo	Courtney Smith	Dario Vitali
Sanya Madan	Mai Nguyen	Michael Roberts	Eric Smith	Rüdiger Volk
Miroslav Madić	Thomas Nikolajsen	Gregory Robinson	Mark Smith	Jeffrey Wagner
Alexis Madriz	Paul Nikolich	Ron Rockrohr	Tim Sneddon	Don Wahl
Carl Malamud	Travis Northrup	Carlos Rodrigues	Craig Snell	Michael L Wahrman
Jonathan Maldonado	Marijana Novakovic	Magnus Romedahl	Craig Snell	Lakhinder Walia
Michael Malik	David Oates	Lex Van Roon	Job Snijders	Laurence Walker
Tarmo Mamers	Ovidiu Obersterescu	Marshall Rose	Ronald Solano	Randy Watts
Yogesh Mangar	Jim Oplotnik	Alessandra Rosi	Asit Som	Andrew Webster
John Mann	Tim O'Brien	David Ross	Ignacio Soto Campos	Jd Wegner
Bill Manning	Mike O'Connor	William Ross	Evandro Sousa	Tim Weil
Diego Mansilla	Mike O'Dell	Boudhayan	Peter Spekrijse	Westmoreland
Harold March	John O'Neill	Roychowdhury	Thayumanavan Sridhar	Engineering Inc.
Vincent Marchand	Carl Ötne	Carlos Rubio	Paul Stancik	Rick Wesson
Normando Marcolongo	Packet Consulting Limited	Rainer Rudigier	Ralf Stempfer	Peter Whimp
Gabriel Marroquin	Carlos Astor Araujo	Timo Ruiter	Matthew Stenberg	Russ White
David Martin	Palmeira	RustedMusic	Martin Štěpánek	Jurrien Wijnhuizen
Jim Martin	Gordon Palmer	Babak Saberi	Adrian Stevens	Joseph Williams
Ruben Tripiana Martin	Alexis Panagopoulos	George Sadowsky	Clinton Stevens	Derick Winkworth
Timothy Martin	Gaurav Panwar	Scott Sandefur	John Streck	Pindar Wong
Carles Mateu	Chris Parker	Sachin Sapkal	Martin Streule	Brian Woods
Juan Jose Marin Martinez	Alex Parkinson	Arturas Satkovskis	David Strom	Makarand Yerawadekar
Ioan Maxim	Craig Partridge	PS Saunders	Colin Strutt	Phillip Yialeloglou
David Mazel	Manuel Uruena Pascual	Richard Savoy	Viktor Sudakov	Janko Zavernik
Miles McCredie	Ricardo Patara	John Sayer	Edward-W. Suor	Bernd Zeimet
Gavin McCullagh	Dipesh Patel	Phil Scarr	Vincent Surillo	Muhammad Ziad
Brian McCullough	Dan Paynter	Gianpaolo Scassellati	Terence Charles Sweetser	Ziayuddin
Joe McEachern	Leif Eric Pedersen	Elizabeth Scheid	T2Group	Tom Zingale
Alexander McKenzie	Rui Sao Pedro	Jeroen Van Ingen	Roman Tarasov	Matteo Zovi
Jay McMaster	Juan Pena	Schenau	David Theese	Jose Zumalave
Mark Mc Nicholas	Luis Javier Perez	Carsten Scherb	Rabbi Rob and	Romeo Zwart
Olaf Mehlberg	Chris Perkins	Ernest Schirmer	Lauren Thomas	廖明沂.
Carsten Melberg	Michael Petry	Benson Schliesser	Douglas Thompson	
			Kerry Thompson	

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Supporters and Sponsors

<p><i>Supporters</i></p> <div>   </div>	<p><i>Diamond Sponsors</i></p> <p>Your logo here!</p>
<p><i>Ruby Sponsors</i></p> <div>   </div>	<p><i>Sapphire Sponsors</i></p> <p>Your logo here!</p>
<p><i>Emerald Sponsors</i></p> <div>      </div> <div>      </div> <div>      </div> <div>      </div>	
<p><i>Corporate Subscriptions</i></p> <div>      </div> <div>     </div>	

For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

John Crain, Senior Vice President and Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Merike Kaeo, Founder and vCISO
Double Shot Security

Olaf Kolkman, Principal – Internet Technology, Policy, and Advocacy
The Internet Society

Dr. Jun Murai, Founder, WIDE Project
Distinguished Professor, Keio University
Co-Director, Keio University Cyber Civilization Research Center, Japan

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

May 2025

Volume 28, Number 1

*A Quarterly Technical Publication for
Internet and Intranet Professionals*

FROM THE EDITOR

In This Issue

From the Editor	1
ShowNet at Interop Tokyo	2
The IPv6 Transition	13
Book Review	33
Fragments	35
Thank You!	36
Call for Papers	38
Supporters and Sponsors	39

The *TCP/IP Interoperability Conference*—later renamed *Interop*—began as a small workshop in August 1986. It quickly grew in scope to incorporate tutorials, and by 1988 an exhibition network connected 51 exhibitors to each other and to the global Internet. This network was designed and deployed by a group of volunteers, and it became the proving ground for many emerging technologies. In 1994, Interop added Tokyo to its international venues, where 30 years later the conference and exhibition attracts more than 120,000 visitors annually. Following an article by David Strom describing the history and evolution of the Interop show network in our previous issue, we now bring you the first installment of an article that describes how this network continues to evolve at the Tokyo Interop show. The article is by Takashi Tomine, Ryo Nakamura, and Ryota Motobayashi—all members of the team that designs and deploys their version called *ShowNet*. A second installment detailing the technologies demonstrated in the 2024 ShowNet will be published in a future issue.

Our previous issue also contained an article about the “IPv6 Mostly” experiment that was conducted during APRICOT 2024 in Bangkok. It is perhaps surprising that we are still very much living in an Internet that is heavily dependent on *IP Version 4* (IPv4) given the amount of time that has passed since the initial *IP Version 6* (IPv6) specifications were published. In our second article, Geoff Huston provides an in-depth analysis on the topic of IPv6 Transition and suggests that perhaps changes in Internet Architecture and technological developments will have us waiting a very long time before IPv4 addressing becomes obsolete.

Book reviews used to be a fairly regular feature in this journal, but it has been quite a long time since we have published any reviews. We asked Craig Partridge to review the book *The Real Internet Architecture: Past, Present, and Future Evolution*, and we hope this latest review will encourage you to send us suggestions for other books on networking and related topics. As always, you can contact us with your feedback by sending an e-mail to: ipj@protocoljournal.org

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

ShowNet at Interop Tokyo: *A Continuously Evolving Demonstration Network*

by Takashi Tomine, Ryo Nakamura, and Ryota Motobayashi

Interop Tokyo, which inherits the objectives of the Interop conference series, is the largest annual exhibition of Internet technologies in Japan. It is held yearly for three days in June. Over 500 exhibitors showcase their products and services at individual booths, and more than 120,000 people visit the venue during the exhibition, as shown in Figure 1. Moreover, a co-located conference offers several dozen sessions and keynote talks.

Figure 1: A view of the Interop exhibition in 2024.



An essential part of Interop Tokyo is *ShowNet*, the largest demonstration network built at Interop Tokyo exhibitions. ShowNet provides network connectivity for Interop exhibitors and attendees but is not limited to this service. Since Interop originates from the word “interoperability,” ShowNet conducts various interoperability tests, experiments, and demonstrations of new networking technologies. For example, in 2019, we deployed service function chaining using *Segment Routing over IPv6* (SRv6)^[1] with four SRv6-capable nodes and five SRv6 proxies^[2] from different vendors.

At that time, SRv6 was an emerging packet-forwarding paradigm, and we faced varied challenges and considerations to achieve this archetype while serving user traffic.

The knowledge we gained through the deployment was published as an Internet Draft^[3]. We have deployed and demonstrated not only routing techniques but also broader technologies, including facilities, optical transport, wireless, security, monitoring, testers, and emerging technologies, such as 5G and media over IP in recent years.

The 2024 ShowNet, featuring comprehensive technical demonstrations, consisted of more than 20 full-height racks, allowing attendees to see the devices running in production. Figure 2 shows a picture of such a ShowNet booth in 2024.

Figure 2: The ShowNet booth in 2024. Whiteboards were mounted on the side walls of each rack, where the NOC members wrote explanations about the devices, design, and technologies. Attendees could see the running devices with those explanations.



History

Interop Tokyo celebrated its 30th event in 2023. In other words, thirty years have passed since an IT trade show, *NetWorld+Interop*, landed in Japan with cutting-edge technologies such as Ethernet 10BASE-T, *Fiber Distributed Data Interface* (FDDI), *Asynchronous Transfer Mode* (ATM) with IP, *Xerox Internetwork Packet Exchange* (IPX), and Apple's *AppleTalk*. Because of a consolidation among other strong Informa sister brands, no Interop event has been held in the United States since 2023. However, Interop is still alive and well in Tokyo, and it maintains its original mission: establishment of multi-vendor interoperability.

In the early 1990s, fathers of the Internet in Japan who visited an Interop event in the US were impressed by its effectiveness—a practical display of interoperability among multi-vendor networking equipment. The groundwork to adapt the event to a Japanese audience began then, and in 1994 Tokyo became one of the host cities of *NetWorld+Interop*, with Las Vegas, Berlin, Atlanta, and Paris (Interop was merged with Novell's *NetWorld*, a similar event that occurred from 1994 to 2004).

As an essential part of the event, Interop ordinarily deploys a temporary show-floor network *InteropNet*—or *ShowNet*, varied by years or venues—to provide Internet connectivity for each exhibitor. This concept was naturally introduced for the Tokyo show too. Until 1997, the design and equipment of InteropNet was basically shared through every show during the annual world tour. Actually, for the first Tokyo show, the construction and verification work of InteropNet (for Tokyo) was held not in Japan but at the Ziff-Davis's *Hot Stage Test Facility* in Sunnyvale, California. The persons of talent for the latest network construction and operation—the *Network Operations Center* (NOC) team—were also invited globally during the initial Tokyo era.

In 1998, following the event organizer's business operations review, Tokyo decided to set up a show-floor network with a local focus. The Tokyo NOC team has focused solely on the Japanese events since that time.

After making that decision in 1998, Tokyo has refrained from using the original *InteropNet* name, and now calls its own network *ShowNet*. ShowNet has since run every year except 2020, when the COVID-19 pandemic prohibited such public gatherings. In addition to Japan's Internet technology community, diverse members from industry and academia now gather every year to continue building and demonstrating ShowNet.

Volunteers are indispensable to achieving such complex networking. Initially named *InteropNet Team Members* (ITMs), volunteers are currently called *ShowNet Team Members* (STMs) in Tokyo. This program, which includes an educational aspect for young students and engineers, continues to be an essential component of ShowNet.

Who Makes ShowNet?

Building ShowNet at the Interop Tokyo exhibition is not an easy feat, so more than 650 engineers with diverse backgrounds are now involved in the project. The NOC team includes the core members, who design the ShowNet network and conduct broad experiments and demonstrations. In recent years, the NOC team has consisted of around 30 expert volunteers from academia, carriers, vendors, etc. They use their areas of expertise and skills to manage the project. Two leaders supervise the ShowNet project; they choose the NOC team members yearly. Teams are either selected by invitation from personal and professional connections or transitioned from STMs or Contributor Members.

The STM program offers a unique opportunity for university students and junior staff from companies to obtain hands-on experience in network operations. Participants in the STM program, who are relatively young and novice network engineers, engage in building ShowNet at the venue as volunteers, and they have the opportunity to touch, configure, operate, and debug various devices and cutting-edge technologies.

This valuable experience is difficult to gain in universities or in their regular workplace. The program also allows young engineers to build and foster relationships and learn from each other by spending two weeks building ShowNet together. In recent years, around 30 slots have been available for participants in this program, but we receive more than twice as many applications each year, so the NOC team members are responsible for the selection process. Figure 3 shows the 2024 STMs.

Figure 3:
The participants in the STM program in 2024.



ShowNet Team Members engaged in building ShowNet.



The third category of engineers involved in ShowNet is the *Contributor Members* who showcase products at Interop Tokyo each year. The contributor vendors make their products and services available to ShowNet and demonstrate them on the live network during the exhibition. These members are skilled engineers from those vendors, and they help build ShowNet with their expertise in the products. The presence of the contributor members is also indispensable for building ShowNet.

A Timeline in a Year

This section briefly introduces a timeline of ShowNet in a year. We, the people involved in ShowNet, put a long-term effort into accomplishing the ShowNet project every year.

Planning

ShowNet covers broad aspects of today's networking technologies. To manage this complexity, we organize the project into working groups, each focusing on a specific field. In 2024, we had 11 working groups leading the following fields: facilities, optical transport, external connectivity, backbone network, data center and cloud, wireless network, monitoring, security, testers, 5G, and media over IP. The NOC team consisted of approximately 30 members, with each working group led by two to four NOC team members who have expertise in that group's area of interest.

Preparation for ShowNet starts in October, the year before the Interop exhibition in June. First, the two leaders gather and organize the NOC team members and begin to discuss topics and technologies that ShowNet will address in the next Interop Tokyo. Then the leaders meet monthly with all NOC team members to share and discuss the overall structure and design of demonstrations. Additionally, each working group holds meetings at least once a month, as needed.

The Contributor Members—vendors providing their products to ShowNet—join the discussion in December. The NOC team members introduce the concept for the next ShowNet and the technologies they want to adopt to the contributor members per working group. Also, the contributors propose their products and use cases they wish to showcase. The NOC team members receive these requests and integrate them into demonstrations. From then until the end of May, the demonstration contents are continuously refined, and the NOC team members consolidate everything into a concrete network design.

Hot Stage

Two weeks before the Interop Tokyo exhibition, we start building ShowNet at the Makuhari Messe exhibition hall. Building ShowNet has two phases: *Hot Stage* and *Deployment*. During hot stage, we build and test all the designs and conduct planned interoperability tests and experiments. In recent years, we have allotted eight days for the hot stage, and we hold two all-hands meetings every day, one in the morning and one in the evening, to share progress as we continue the construction of ShowNet.

When the hot stage begins, all members of the NOC team, the STM, and the contributor members gather at the venue and start building the network. First, we install every device in the right place on the racks, turn on the devices, and check their status. Checking device status is very important because some devices are transported directly from overseas to the venue, so it is necessary to ensure that they are not malfunctioning. We usually finish this process on the first day.

On the second day, we start the network setup: connect appropriate links between devices with patch cables as designed. After the physical network connections are completed, NOC members in charge of the backbone network start configuring the backbone routers with the help of the ShowNet team members. In the early days, ShowNet backbone was a simple Layer-3 network with a single *Interior Gateway Protocol* (IGP) instance. But now, ShowNet adopts several overlay technologies such as *Multiprotocol Label Switching* (MPLS), SRv6, and *Virtual eXtensible Local Area Network* (VXLAN), so we have to configure more overlays after the Layer-3 routing configuration.

The working groups other than the backbone network group prepare their demonstrations in parallel. Every part of ShowNet is built with multi-vendor equipment, so we have to check interoperability everywhere. The working groups also conduct several interoperability tests during the hot stage. These interoperability tests are beneficial for finding bugs or slight differences in implementation.

Sometimes, these bugs or differences are critical to building ShowNet, so contributor members from vendors try to fix them with their development teams.

Testing the network is always essential, even in an event network. In ShowNet, we conduct failover tests in the latter part of the hot stage—stop and resume each backbone router sequentially and confirm that routing redundancy works as expected. If troubles arise during the test, the backbone network group of the NOC team and the ShowNet team members troubleshoot and debug the problems together. This collaborative troubleshooting process is also a good hands-on experience for the junior network engineers of STM. Figure 4 shows a snapshot of a failover test.

Figure 4: A snapshot of the failover test in 2024. Red lines in the display indicate that some user segments have unexpectedly lost the connectivity, and the NOC team and ShowNet team members start to troubleshoot. The software used here is deadman^[5], which was designed and implemented for ShowNet.



Deployment

In 2024, after we finished the hot stage, we started to deploy the ShowNet network in the whole Makuhari Messe venue four days before the Interop Tokyo exhibition began. Interop Tokyo used five halls in Makuhari Messe. The ShowNet network spread to each hall with optical transport from the ShowNet booth. Every hall had a small booth on which access switches were installed for ShowNet to extend the network to all the exhibitors' booths. Electrical construction members deployed optical fibers from the ShowNet booth to the small booths in each hall and copper cables from the switches on the booths to the exhibitors' booths. After spreading the cables, the ShowNet team members connected these cables to the access switches of the ShowNet backbone and checked the correctness of Layer-1 to Layer-7 connectivity. Figure 5 shows a snapshot of such a scene.

Figure 5: Three ShowNet Team members and a NOC team member are checking deployed cables for exhibitor booths in an exhibition hall.

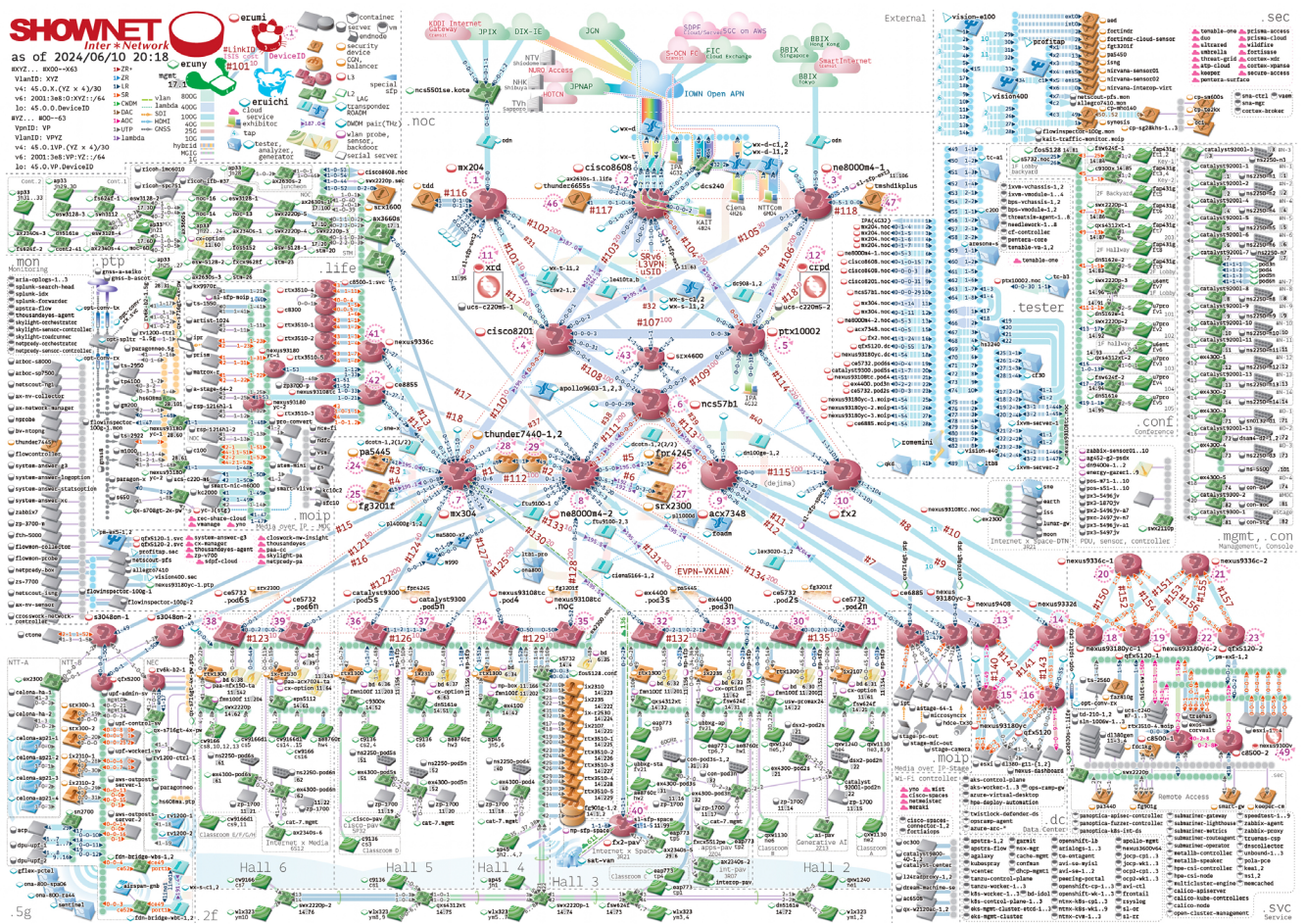


As the deployment phase begins, exhibitors of Interop Tokyo also arrive at the venue and start preparing their booths. ShowNet provides Internet connectivity for the demonstrations held in their booths. If they have any problems on the ShowNet network, they go to the *ShowNet Service Counter* and describe their problem, and we start to identify and resolve it. Usually, problems in this phase are caused by physical things like a cut cable or mis-connection of cables or access switches, but sometimes some logical bugs cause critical challenges. Such bugs are sometimes difficult to solve because we must fix them before the exhibition starts.

It is also crucial to ensure the visibility and presentation of the equipment and services contributed to ShowNet. After the whole ShowNet network is built, we tidy up the ShowNet booth. We try to ensure that every piece of equipment inside the racks is presented well, because it is not only a device but also an exhibit. We also post captions for all equipment and prepare description slides for attendees.

After we finish all the processes for building the ShowNet, we complete the network diagram. Figure 6 shows the network diagram of ShowNet 2024. You can see all the devices, links, services, and designs of ShowNet 2024 on this diagram. One of the NOC members creates the diagram. From the hot stage onward, the same member continuously monitors all the design and configuration changes and keeps the diagram up-to-date. Eventually, the diagram captures all of the ShowNet network on a single sheet. This diagram is an essential tool for ShowNet: engineers use it to grasp the overall network design, communicate and share changes, and troubleshoot problems. The diagram is practical and functional, especially when building the network and troubleshooting.

Figure 6: Diagram of the ShowNet network at Interop Tokyo 2024. The full-size version is available^[6]. Green icons indicate Layer-2 devices, red icons are Layer-3, and yellow icons are Layer-4 and above devices. The icons in the diagram are also available under the Creative Commons license^[7].



Exhibition

The Interop Tokyo exhibition usually starts on a Wednesday in mid-June and lasts three days. When it starts, we begin operating the ShowNet network and also offer some presentations for attendees. ShowNet is one of the main parts of Interop Tokyo, so many attendees come to the ShowNet booth. Attendees can see the demonstrations, learn about the technologies, observe the devices running in the racks, and see and feel the functions that run on ShowNet.

During the exhibition, ShowNet members, especially the NOC team members, work on showcasing the ShowNet network. The NOC room inside the ShowNet booth is equipped with large screens that show several tools monitoring the network status and security, as shown in Figure 7. The NOC team members operate ShowNet from there, and attendees can observe them; actually, the NOC room is also a part of the show. In addition, the NOC team members conduct tours, called *ShowNet Walking Tours*, during the exhibition. Participants on the tours can walk around each ShowNet rack with a NOC member and observe the running devices, and the NOC member will answer any questions. Figure 8 captures a scene from one of the tours. The NOC team members also give several presentations at the ShowNet booth and in session rooms at the exhibition.

Figure 7: The NOC room on the exhibition floor, visible to attendees.



Tear-Down

After the 2024 exhibition ended at 5 p.m. on Friday, we started to tear down the ShowNet network. Our contract required that we vacate the halls by midnight. First, the NOC team and contributor members shut down devices, a requirement before they could be powered off. Next, we shut off all power supplies, unplugged the patch cables, unmounted the devices, and returned them to the contributors. After that, the NOC and ShowNet team members wound up all cables and cleaned up all racks for use next year. Figure 9 shows the racks during tear-down.

Conclusion

The ShowNet network is a unique environment. It not only provides Internet connectivity for exhibitors and attendees, it also displays a large-scale ephemeral event network that will demonstrate emerging and cutting-edge technologies. ShowNet conducts various interoperability tests, experiments, and demonstrations with numerous devices contributed by multiple vendors. Furthermore, ShowNet offers an invaluable opportunity for engineers to collaborate with diverse engineers from different fields. We believe that the connections and relationships among them established through ShowNet have contributed to revitalizing network communities. In addition, it has more than 30 years of experience and adds to our knowledge to handle these cutting-edge trials. Interop Tokyo will continue this work inherited from US Interop.

Acknowledgments

We would like to express our gratitude to everyone around the world who has been involved in all Interop events, from the past to the present.

Figure 8: ShowNet Walking Tour: NOC team members explain the network to attendees in front of each rack.

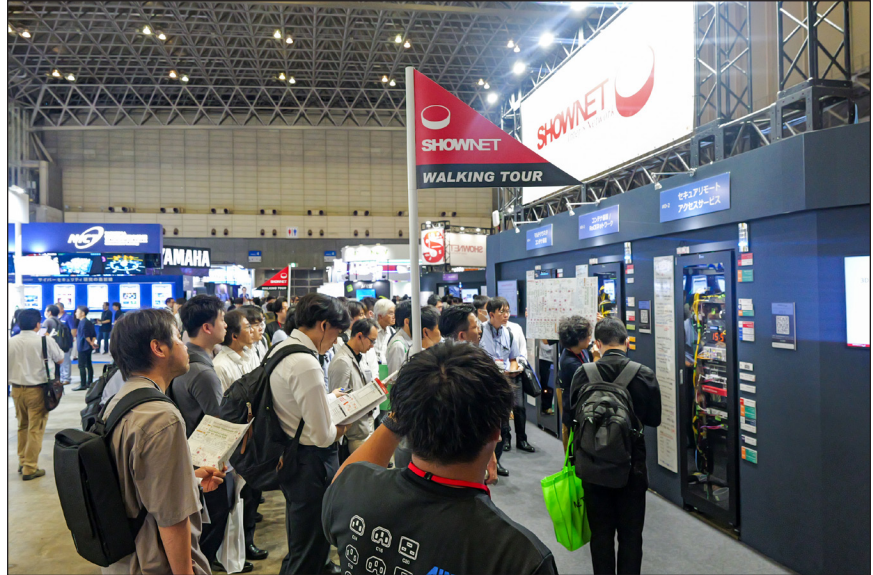


Figure 9: A scene of the tear-down process in 2024. All devices were unmounted from the racks.



References and Further Reading

- [0] David Strom, “The Interop Shownet,” *The Internet Protocol Journal*, Volume 27, No. 3, October 2024.
- [1] Clarence Filsfil, Pablo Camarillo, John Leddy, Daniel Voyer, Satoru Matsushima, and Zhenbin Li, “Segment Routing over IPv6 (SRv6) Network Programming,” RFC 8986, February 2021.

- [2] Francois Clad, Xiaohu Xu, Clarence Filsfils, Daniel Bernier, Cheng Li, Bruno Decraene, Shaowen Ma, Chaitanya Yadlapalli, Wim Henderickx, and Stefano Salsano, “Service Programming with Segment Routing,” Internet-Draft, Work in Progress, February 2025.

`draft-ietf-spring-sr-service-programming-09`

- [3] Ryo Nakamura, Yukito Ueno, and Teppei Kamata, “An experiment of SRv6 Service Chaining at Interop Tokyo 2019 Shownet,” Internet-Draft, Work in Progress, October 30, 2019.

`draft-upa-srv6-service-chaining-exp-00`

- [4] Glenn Evans, “Inside InteropNet’s Hot Stage,” *Network Computing*, April 2013.

- [5] *deadman*: <https://github.com/upa/deadman>

- [6] Interop 2024 ShowNet map:

`https://www.interop.jp/2024/assets/file/e-web.pdf`

- [7] ShowNet map icons:

`https://github.com/interop-tokyo-shownet/shownet-icons`

TAKASHI TOMINE received a Master’s degree from Keio University, Japan, and finished his Ph.D. program without a dissertation at Keio University. He is now an Associate Senior Research Engineer at the National Astronomical Observatory of Japan. He has been an Interop Tokyo ShowNet NOC team generalist since 2013. His research interests include network operation, international research educational networks, and cybersecurity. He can be reached at: **`tomine@interop-tokyo.net`**

RYO NAKAMURA received his Ph.D. degree in Information Science and Technology from the University of Tokyo, Tokyo, Japan, in 2017. He is currently an Associate Professor at the Information Technology Center, the University of Tokyo, where he operates the university’s campus network. His research interests include networking in operating systems, network virtualization, and network operations. Since 2009, he has been involved in Interop Tokyo ShowNet, as a ShowNet team member until 2011, and as a member of the NOC team from 2012 to the present. He has been primarily responsible for the backbone network of ShowNet, and led demonstrations of SDN-related technologies from 2013 to 2017. He can be reached at: **`ryo@interop-tokyo.net`**

RYOTA “ROY” MOTOBAYASHI holds a Bachelor of Engineering from Shinshu University and is qualified as CISSP, Japan’s Registered Information Security Specialist and Information Technology Engineer (Class I and Network Specialist). He went through various networking-related projects, from hardware design to corporate strategy planning for NEC Corporation 1988-2023. Since 2024, he has worked for Telecom Engineering Center, a certification body in Japan. His long-term contributions to Interop Tokyo are NOC 1994–1996, NOC Advisory 2006–2017, and Program Committee 2004–2023. He can be reached at: **`jj1wt1@jar1.com`**

The IPv6 Transition

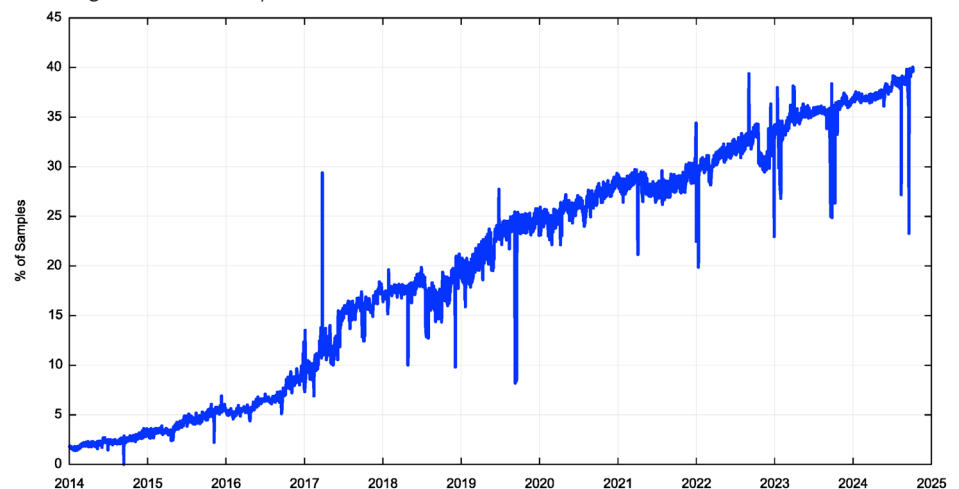
by Geoff Huston

The state of the transition to IPv6 within the public Internet continues to confound us. RFC 2460^[1], the first complete effort at a specification of the IPv6 protocol, was published in December 1998, more than 25 years ago. The entire point of IPv6 was to specify a successor protocol to IPv4 because of the prospect of running out of IPv4 addresses. Yet while the public Internet ran out of IPv4 addresses more than a decade ago, the contrary observation is that this network platform is still largely sustained through its use of IPv4. The transition of the public Internet to the IPv6 protocol has been going on for 25 years now, and if there were any urgency to be instilled in the transition effort by the prospect, and then the reality, of IPv4 address exhaustion, then we’ve been living with exhaustion a very long time now, and we’re largely inured to it. It’s probably time to ask the question again: How much longer will this transition to IPv6 take?

At APNIC Labs^[0] we’ve been measuring the uptake of IPv6 for more than a decade now. We use a measurement approach that looks at the network from the perspective of the Internet user base. What we measure is the proportion of users who can reach a published service when the only means to do so is by using IPv6. The data is gathered using a measurement script embedded in an online ad, and the ad placements are configured to sample a diverse collection of end users continually.

Figure 1 displays the IPv6 adoption report showing our measurements of IPv6 adoption across the Internet user base from 2014 to 2024.

Figure 1: IPv6 Adoption – 2014 to 2024. (APNIC Labs Data)



On the one hand, Figure 1 is one of those classic “up and to the right” Internet curves that shows continues growth in the adoption of IPv6. The problem is in the values in the scale of the Y-axis. The issue here is that in 2024 we were at a level where only a little more than one-third of the Internet user base could access an IPv6-only service. Everyone else, now in 2025, is still in an IPv4-only Internet.

This situation appears to be completely anomalous. It's been more than a decade since the supply of "new" IPv4 addresses was exhausted, and the Internet has not only been running on empty, but also is now tasked to span an ever-increasing collection of connected devices—and it has achieved this feat without collapsing. In late 2024 it is variously estimated (or guessed!) that some 20 billion devices used the Internet, yet the Internet IPv4 routing table encompasses only some 3.03 billion unique IPv4 addresses. The original "end-to-end" architecture of the Internet assumed that every device was uniquely addressed with its own IP address, yet the Internet is now sharing each individual IPv4 address across an average of 6 devices, and apparently it all seems to be working! If "end-to-end" was the sustaining principle of the Internet architecture in the 1980's, then as far as the current users of IPv4-based access and services across the public Internet are concerned, it's all over!

IPv6 was meant to address these issues, and the 128-bit wide address fields in the protocol have sufficient address space to allow every connected device to use its own unique address. The design of IPv6 was intentionally very conservative. To a first level of approximation IPv6 is simply "IPv4 with bigger addresses." There are also some changes to fragmentation controls, the address acquisition protocols [*Address Resolution Protocol* (ARP) vs. *Neighbour Discovery*], and the *IP Options* fields, but the upper-level transport protocols are unchanged. IPv6 was intended to be a largely invisible change to a single level in the protocol stack, and definitely not intended to be a massive shift to an entirely novel networking paradigm.

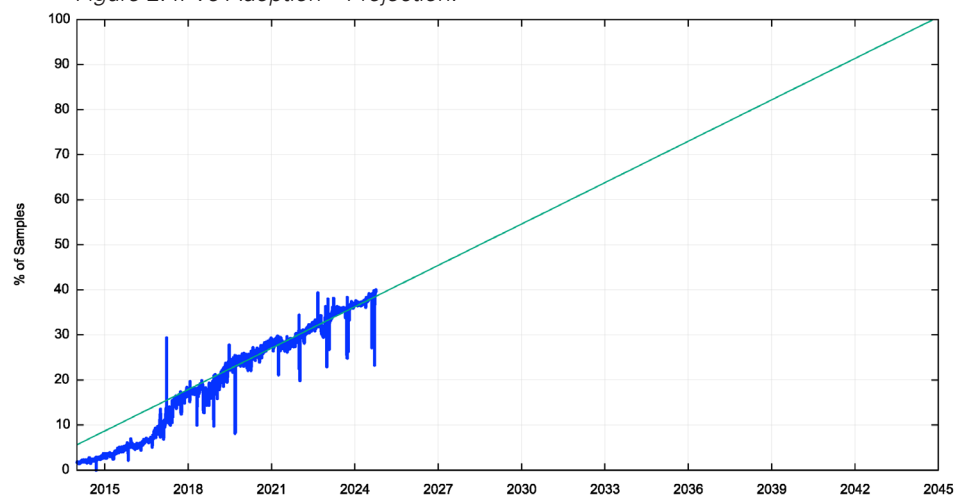
In the sense of representing a very modest incremental change to IPv4, the IPv6 design achieved its objective, but in so doing it necessarily provided little in the way of any marginal improvement in protocol use and performance. IPv6 was no faster, no more versatile, no more secure than IPv4. The major benefit of IPv6 was to mitigate the future risk of IPv4 address exhaustion. In terms of conventional market operations, many markets, including that of the Internet, apply a hefty discount factor to future risk. The result is that the level of motivation to undertake this transition is highly variable given that the expenditure to deploy this second protocol does not immediately realize tangible benefits in terms of lower cost, greater revenue, or greater market share. In a networking context where market-based coordination of individual actions is essential, a level of diversity of views of the net value of running a dual-stack network often leads to reluctance on the part of individual actors and sluggish progress of the common outcome of the transition. As a result, there is no common sense of urgency.

To illustrate this fact, we can look at the time series shown in Figure 1 and ask the question: "If the growth trend of IPv6 adoption continues at its current rate, how long will it take for every device to be IPv6-capable?"

Asking this question is the same as looking at a linear trend line placed over the data series used in Figure 1 for the date when this trend line reaches 100%. Using a least-squares best fit for this data set from January 2020 to the present day, and using a linear trend line, we can come up with Figure 2.

This exercise predicts that we'll see completion of this transition in late 2045, or some 20 years into the future. It must be noted that there is no deep modelling of the actions of various service providers, consumers, and network entities behind this prediction. The only assumption that drives this prediction is that the forces that shaped the immediate recent past are unaltered when looking into the future. In other words, this exercise simply assumes that "tomorrow is going to be a lot like today."

Figure 2: IPv6 Adoption – Projection.



The projected date in Figure 2 is less of a concern than the observation that this model predicts a continuation of this transition for a further two decades. If the entire concept of IPv6 was to restore a coherent address plan across the collection of Internet-connected devices, then placing this model of coherent unique device addressing in abeyance for some 30 years, from around 2015 through to 2045, leads to questioning the role and value of such a unique device addressing framework in the first place! If we can operate a fully functional Internet without such a coherent end-device address architecture for three decades, why would we feel the need to restore address coherence at some point in the future? What's the point of IPv6 if it's not address coherence?

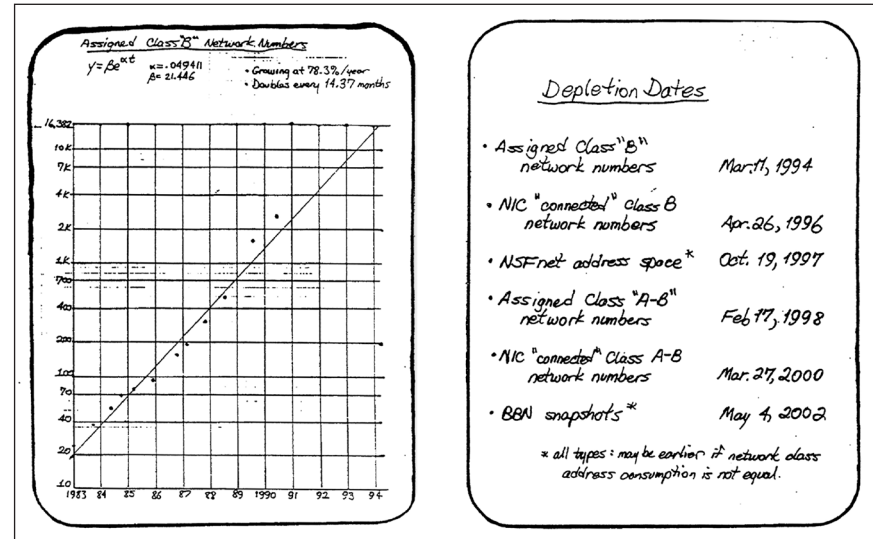
Something has gone very wrong with this IPv6 transition, and that's what I will examine in this article.

A Little Bit of History

By 1990 it was clear that IP had a problem. The Internet was still tiny at the time, but the growth patterns were exponential, doubling in size every 12 months. We were stressing out the pool of Class B IPv4 addresses, and in the absence of any corrective measures this address pool would be fully depleted in 1994 (Figure 3).

Frank Solensky presented predictions at the 18th meeting of the *Internet Engineering Task Force (IETF)*^[10].

Figure 3: IPv4 Depletion Predictions, Frank Solensky.



We were also placing pressure on the routing system at the time. The deployed routers in 1992 had only enough memory to support a further 12 to 18 months of routing growth. The combination of these routing and addressing pressures was collectively addressed in the IETF at the time under the umbrella of the ROAD effort, as described in RFC 1380^[2].

There was a collection of short-, medium- and longer-term responses that were adopted in the IETF to address the problem. In the short term, the IETF dispensed with the class-based IPv4 address plan and instead adopted a variably sized address prefix model. Routing protocols, including the *Border Gateway Protocol (BGP)*, were quickly modified to support these classless address prefixes. Variably sized address prefixes added additional burdens to the address-allocation process, and in the medium term the Internet community adopted the organisational measure of the *Regional Internet Registry (RIR)* structure to allow each region to resource the increasingly detailed operation of address-allocation and registry functions for their region. These measures increased the specificity of address allocations and provided the allocation process with a more exact alignment to determine adequate resource allocations that permitted a more diligent application of relatively conservative address-allocation practices. These measures realized a significant increase in address usage efficiency. The concept of "address sharing" using *Network Address Translation (NAT)*^[3] also gained some traction in the *Internet Service Provider (ISP)* world. Not only did NATs dramatically simplify the address administration processes in ISPs, they also played a major role in reducing the pressures on overall address consumption.

The adoption of these measures across the early 1990's pushed a 2-year imminent crisis into a more manageable decade-long scenario of depletion. However, they were not considered to be a stable long-term response. It was thought at the time that an effective long-term response really needed to extend the 32-bit address field used in IPv4. Then the transition from mainframe to laptop was well underway in the computing work, and the prospect of further reductions in size and expansion of deployment in smaller embedded devices was clear. An address space of 4 billion was just not large enough for what was likely to occur in the coming years in the computing world.

But in looking at a new network protocol with a vastly increased address space, the IETF realized that any such change would not be backward-compatible with the installed base of IPv4 systems. As a result, there were a few divergent schools of thought as to what to do. One approach was to jump streams and switch over to use the Connectionless Transport profile of the *Open Systems Interconnection* (OSI) Protocol Suite and adopt OSI *Network Service Access Point Address* (NSAP) addresses along the way. Another was to change as little as possible in IP except the size of the address fields. And numerous ideas were thrown about in the area of proposing significant changes to the IP model.

By 1994 the IETF had managed to settle on the minimal change approach, which was IPv6. The address field was expanded to 128 bits, a *Flow ID* field was introduced, fragmentation behaviour was altered and pushed into an optional header, and ARP was replaced with *multicast*.

The main thing to note was that IPv6 did not offer any new functionality that was not already present in IPv4. It did not introduce any significant changes to the operation of IP. It was just IP with larger addresses.

Transition

While the design of IPv6 consumed a lot of attention at the time, the concept of transition of the network from IPv4 to IPv6 did not.

Given the runaway adoption of IPv4, there was a naive expectation at that time that IPv6 would similarly just take off, and there was no need to give the transition much thought. In the first phase, we would expect to see applications, hosts, and networks adding support for IPv6 in addition to IPv4, transforming the Internet into a dual-stack environment. In the second phase we could then phase out support for IPv4. The expectation was that the process would take a few years.

This plan had numerous problems. Perhaps the most serious one was a resource-allocation problem. The Internet was growing extremely quickly, and most of our effort was devoted to keeping pace with demand. More users, more capacity, larger servers, more content and services, more responsive services, more security, better defence. All of these factors shared a common theme: *scale*.

We could either concentrate our resources on meeting the incessant demands of scaling, or we could work on IPv6 deployment. The short- and medium-term measures that we had already taken had addressed the immediacy of the problems of address depletion, so in terms of priority, scaling was a far more important priority for the industry than IPv6 transition. Through the decade from 1995 to 2005 the case for IPv6 quietly slumbered in terms of mainstream industry attention.

IPv4 addresses were still available, and the use of *Classless Inter-Domain Routing* (CIDR) and far more conservative address-allocation practices had pushed the prospect of IPv4 address depletion out by more than a couple of decades. Many more pressing operational and policy issues for the Internet absorbed the industry's collective attention in those days.

However, this period of respite was brief. The scaling problem accelerated by a whole new order of magnitude in the mid 2000's with the introduction of the iPhone and its brethren^[4]. Suddenly this scale problem was not just of the order of tens or even hundreds of millions of households and enterprises, it transformed into a problem of billions of individuals and their personal devices, and it added mobility into the mix. As a taste of a near-term future, the production scale of these "smart" devices quickly ramped up into annual volumes of hundreds of millions of units. The entire reason why IPv6 was a necessity was coming into fruition, but at this stage we were just not ready to deploy IPv6 in response. Instead, we rapidly increased our consumption of the remaining pools of IPv4 addresses and we supported the first wave of large-scale mobile services with IPv4. Dual stack was not even an option in the mobile world at the time. The rather bizarre economics of financing 3G infrastructure meant that dual-stack infrastructure in a 3G platform was impractical, so IPv4 was used to support the first wave of mobile services. This situation quickly turned to IPv4 and NATs as the uptake of mobile services gathered momentum.

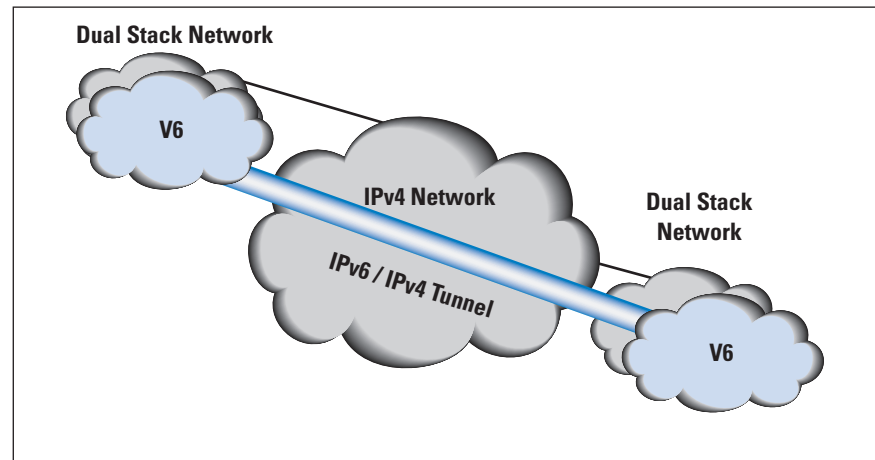
At the same time the decentralised nature of the Internet was hampering IPv6 transition efforts. What point was there in developing application support for IPv6 services if no host had integrated IPv6 into its network stack? What point was there in adding IPv6 to a host networking stack if no ISP was providing IPv6 support? And what point was there for an ISP to deploy IPv6 if no hosts and no applications would use it? In terms of IPv6 at this time, nothing happened.

The operating-system sector made the first efforts to try to break this impasse of mutual dependence, and fully functional IPv6 stacks were added to the various flavours of Linux, Windows, and MacOS, as well as in the mobile host stacks of iOS and Android.

But even these implementations were not enough to allow a transition to achieve critical momentum. It could be argued that this situation made the IPv6 situation worse and set back the transition by some years.

The problem was that with IPv6-enabled hosts there was some desire to use IPv6, but these hosts were isolated “islands” of IPv6 sitting in an ocean of IPv4. The concentration of the transition effort then fixated on various tunnelling methods to tunnel IPv6 packets through the IPv4 networks (Figure 4). While you can perform this tunnelling manually when you have control over both tunnel endpoints, this approach was not that useful. What we wanted was an automated tunnelling mechanism that took care of all these details.

Figure 4: Phase 1 of the IPv6 Transition.



The first such approach that gathered some momentum was 6to4^[5]. The first problem with 6to4 was that it required public IPv4 addresses, so it could not provide services to IPv6 hosts that were behind a NAT. The more critical problem was that firewalls had no idea how to handle these 6to4 packets, and the default action when in doubt is to deny access. So 6to4 connections encountered an average of a 20 to 30% failure rate in the public Internet, making it all but unusable as a mainstream service. The NAT traversal issue was also a problem, so a second auto-tunnel mechanism was devised that performed NAT sensing and traversal. This mechanism, *Teredo*^[6], was even worse in terms of failure rates, and some 40% of Teredo connection attempts were observed to fail^[7].

Not only were these Phase 1 IPv6 transition tools extremely poor performers, as they were so unreliable, but even when they worked the connection was both fragile and slower than IPv4. The result was perhaps predictable, even if unfair. It was not just the transition mechanisms that were viewed with disfavour, but IPv6 itself also attracted some opprobrium.

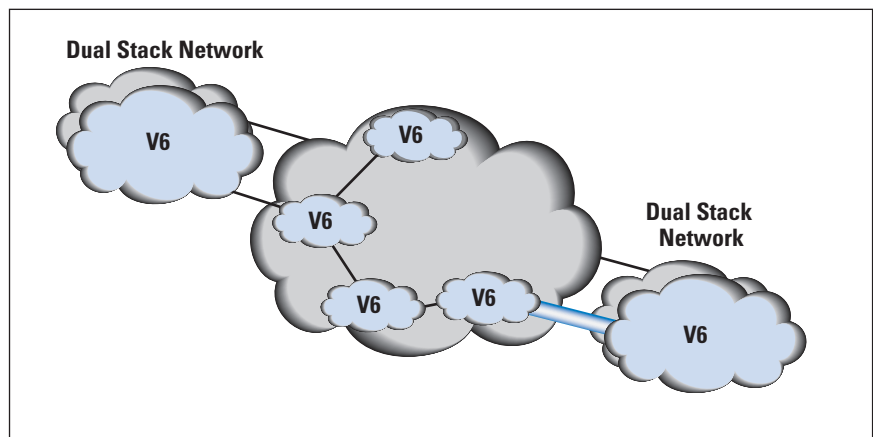
Up until around 2011 IPv6 was largely ignored as a result in the mainstream of the public Internet. A small number of service providers tried to deploy IPv6, but in each case they found themselves with a unique set of challenges that they and their vendors had to solve, and without a rich set of content and services on IPv6, the value of the entire exercise was highly dubious! So, nothing much happened.

Movement at Last!

It wasn't until the central IPv4 address pool that the *Internet Assigned Numbers Authority* (IANA) managed was depleted at the start of 2011, and the first RIR, APNIC, ran down on its general allocation pool in April of that year, that the ISP industry started to pay more focussed attention to this IPv6 transition.

At around the same time, the mobile industry commenced its transition into 4G services. The essential difference between 3G and 4G was the removal of the *Point-to-Point Protocol* (PPP) tunnel through the radio access network from the gateway to the device and its replacement by an IP environment. This solution allowed a 4G mobile operator to support a dual-stack environment without an additional cost component, and it was a major enabler for IPv6. Mapping IPv4 into IPv6 (or the reverse) is fragile and inefficient for service providers as compared to native dual stack. In the 6-year period from 2012 to the start of 2018, the level of IPv6 deployment rose from 0.5 to 17.4%. At this stage IPv6 was no longer predominately tunnelled, as many networks supported IPv6 in native mode (Figure 5).

Figure 5: Phase 2 of the IPv6 Transition.



The problem here was that we were late with this phase of the transition. The intention of this transition was to complete the work and equip every network and host with IPv6 before we ran out of IPv4 addresses (Figure 6).

The position we had arrived at by 2012 was far more challenging. The pools of available IPv4 address space were rapidly depleting, and the regional address policy communities were introducing highly conservative address-allocation practices to eke out the remaining address pools. At the same time the amount of IPv6 uptake was minimal. The transition plan for IPv6 was largely broken (Figure 7).

Figure 6: The IPv6 Transition Plan.

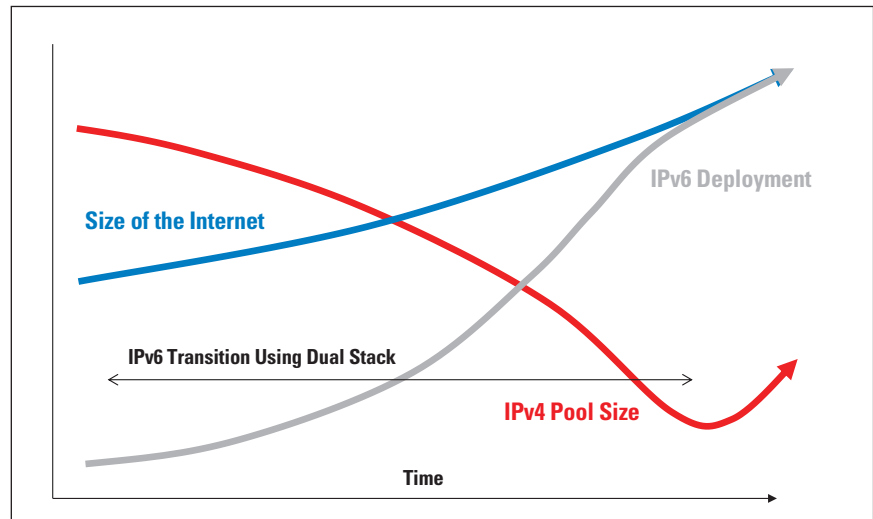
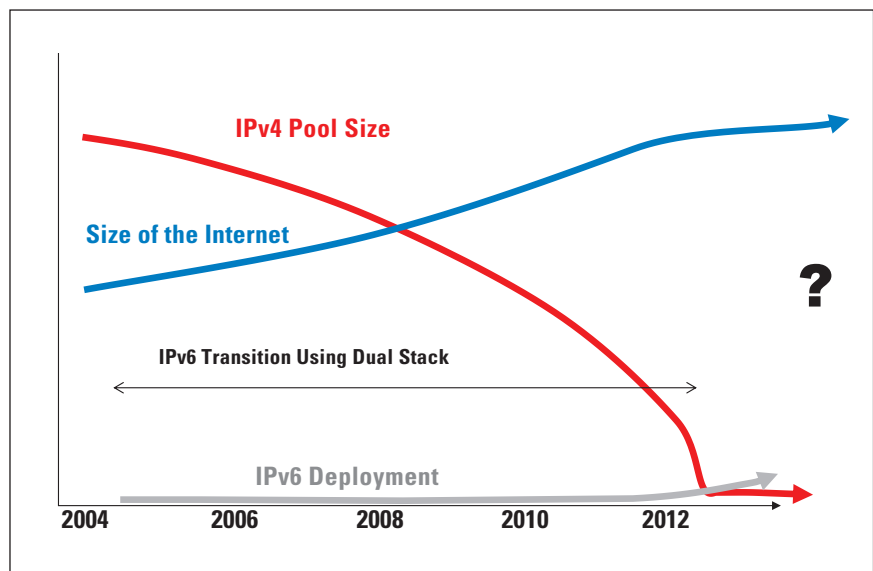


Figure 7: The IPv6 Transition Plan in 2012.



NATs and Address Scarcity Pressures

At this point there was no choice for the Internet, and to sustain growth in the IPv4 network while we were waiting for IPv6 to gather momentum we turned to NATs. NATs were a challenging subject for the IETF. The entire concept of coherent end-to-end communications was to eschew active middleware such as NATs in the network. NATs created a point of disruption in this model, thereby causing a critical dependency upon network elements. They removed elements of network flexibility from the network and at the same time reduced the set of transport options to the *Transmission Control Protocol* (TCP) and *User Datagram Protocol* (UDP).

The IETF resisted any efforts to standardise the behaviour of NATs, fearing perhaps that standard specifications of NAT behaviour would bestow a legitimacy on the use of NATs, an outcome that many IETF participants were very keen to avoid.

This aversion did not reduce the level of impetus behind NAT deployment. We had run out of IPv4 addresses and IPv6 was still a distant prospect, so NATs were the most convenient solution. What this action did achieve was to create a large variance of NAT behaviours^[15] in various implementations, particularly with respect to UDP behaviours. This situation has exacted a cost in software complexity where an application needs to dynamically discover the type of NAT (or NATs) in the network path if it wants to perform anything more complex than a simple two-party TCP connection.

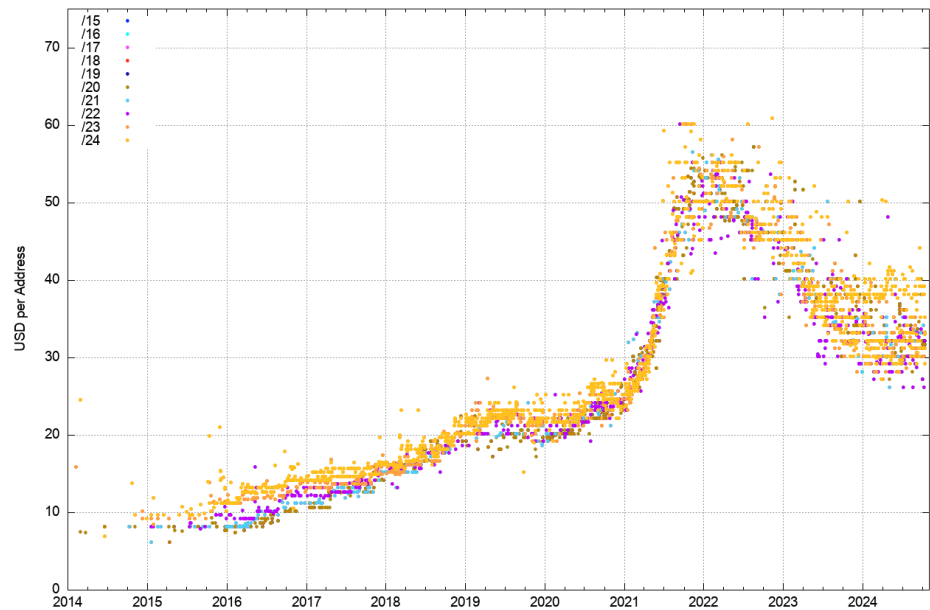
Despite these issues NATs were a low-friction response to IPv4 address depletion, where individual deployment could be undertaken without incurring external dependencies. On the other hand, deployment of IPv6 was dependant on other networks and servers also deploying IPv6. NATs made highly efficient use of address space for clients, as not only could a NAT use the 16-bit source port field, but by time-sharing the NAT binding, NATs achieved an even greater level of address efficiency. A major reason why we've been able to sustain an Internet with tens of billions of connected devices is through the widespread use of NATs.

Server architectures were also changing. The introduction of *Transport Layer Security* (TLS)^[8] into the web-server world included a point in TLS session establishment where the client informs the server platform the name of the service that it intends to connect to. Not only did this information allow TLS to validate the authenticity of the service point, but it also allowed a server platform to host an extremely large collection of services from a single platform (and a single platform IP address) and perform individual service selection via this *TLS Server Name Indication* (SNI). The result is that server platforms perform service selection by name-based distinguishers [*Domain Name System* (DNS) names] in the session handshake, allowing a single server platform to serve large numbers of individual servers. The implications of the widespread use of NATs and the use of server sharing in service platforms has taken the pressure off the entire IPv4 address environment.

One of the best ways to illustrate the changing picture of address scarcity pressure in IPv4 is to look at the market price of address transfers over the past decade. Scarcity pressure is reflected in the market price. Figure 8 shows a time series of the price of traded IPv4 addresses.

The period of the COVID outbreak coincided with a rapid price escalation over 2021, but the price has since declined to between \$30 and \$40 per address, and this price, admittedly over a \$16 range from \$26 to \$42 per address, was stable across 2024. This price data indicates that IPv4 addresses were still in demand in 2024, but the level of demand appears to have equilibrated against available levels of supply, implying that there was no scarcity premium in evidence in the address market in 2024. This data points to the combination of the efficacy of NATs in extending the efficiency of IPv4 addresses by using the 16 bits of port address space plus the additional benefits of using shared address pools.

Figure 8: Market Prices of IPv4 Address Transfers. (Data from Hilco Streambank)



However, it's not just IPv4 that has alleviated the scarcity pressure for IPv4 addresses. Figure 1 indicates that over the past decade the level of IPv6 adoption has risen to encompass some 40% of the user base of the Internet. Most applications, including browsers, support *Happy Eyeballs*^[9], which is a shorthand notation for preferring to use IPv6 over IPv4 if both protocols are available for use in support of a service transaction. As network providers roll out IPv6 support, the pressure on their IPv4 address pools for NAT use is relieved because the applications prefer to use IPv6 where available.

How Much Longer?

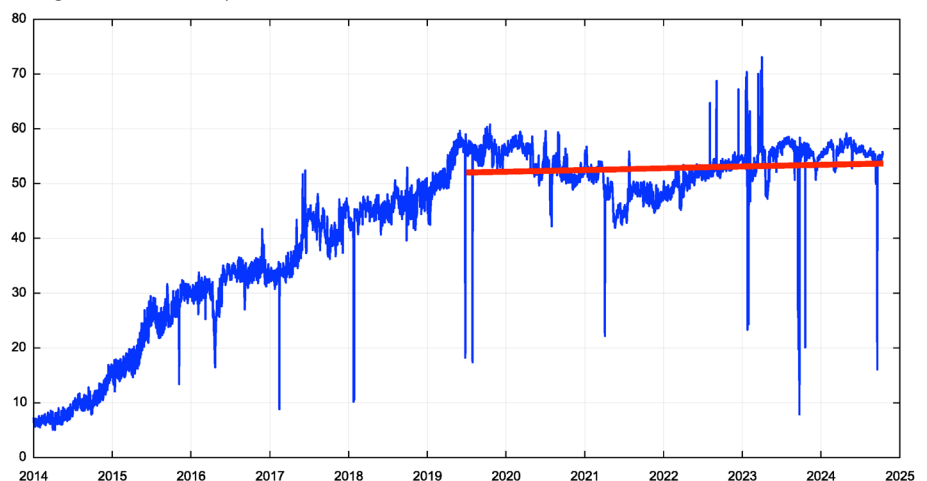
Now that we are somewhere in the middle of this transition, the question is: How much longer is this transition going to take?

This question seems simple, but it does need a little more elucidation. What is the “end point” when we can declare the transition to be over? When will this transition be “complete”? Is it the time when the Internet has no more IPv4-based traffic? Or is it the time when the Internet no longer requires IPv4 in public services? Or do we mean the point when IPv6-only services are viable? Or perhaps we should look at the market for IPv4 addresses and define the endpoint of this transition at the time when the price of IPv4 addresses completely collapses? Perhaps we can take a more pragmatic position here and rather than looking for completion as the point when the Internet is completely bereft of all use of IPv4 addresses and their use, we could define “completion” as the point when use of IPv4 is no longer necessary. The implication would be that when a service provider can operate a viable Internet service using only IPv6 and having no supported IPv4 access mechanisms at all, then we would have completed this transition.

What is the implication? Certainly, the ISP needs to provide IPv6. But all the connected edge networks and the hosts in these networks need to support IPv6 as well. After all, the ISP has no IPv4 services at this point of completion of the transition. It also implies that all the services the clients of this ISP use must be accessible over IPv6. Yes, this accessibility includes all the popular cloud services and cloud platforms, all the content streamers, and all the content-distribution platforms. It also includes specialised platforms such as Slack, Xero, Atlassian, and similar platforms. The data published at the Internet Society's *Pulse* page^[11] reports that only some 47% of the top 1000 web sites are reachable over IPv6, so clearly a lot of service platforms have work to do, and this work will take more time.

When we look at the IPv6 adoption data for the United States, another somewhat curious anomaly is evident (Figure 9).

Figure 9: IPv6 Adoption in the US - 2014 to 2024. (APNIC Labs Data)



The data shows that the level of IPv6 use in the US has remained constant since mid-2019. Why is there no further momentum to continue with the transition to IPv6 in this part of the Internet? I would offer the explanation that the root cause is a fundamental change in the architecture of the Internet.

Changes to the Internet Architecture

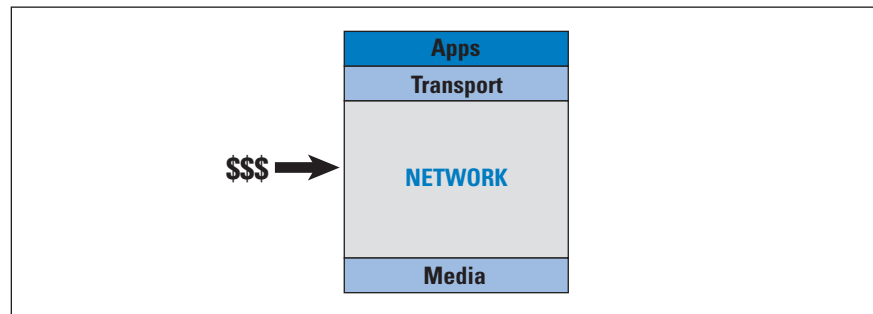
The major change to the Internet architecture is a shift away from a strict address-based architecture. Clients no longer need to use a persistent unique public IP address to communicate with servers and services. And servers no longer need to use a persistent unique public IP address to provide clients with access to the service or content. Address scarcity takes on an entirely different dimension when unique public addresses are not required to number every client and every distinct service.

Some of the clues that show the implications of this architectural shift are evident when you look at the changes in the internal economy of the Internet. The original model of IP was a network protocol that allowed attached devices to communicate with each other.

The network providers supplied the critical resource to allow clients to consume content and access services. At the time the costs of the network service dominated the entire cost of the operation of the Internet, and in the network domain distance was the dominant cost factor.

Network providers who offered distance services (so-called “transit providers”) were the dominant ones. Little wonder that we spent a lot of our time working through the issues of interconnection of network service providers, customer/provider relationships, and various forms of peering and exchanges. The ISPs were in effect brokers in the rationing of the scarce resource of distance capacity. This economy was a classic network economy (Figure 10).

Figure 10: The Classic Network Economy.



For many years the demand for communications services outstripped available inventory, and price was used as a distribution function to moderate demand against available capacity. However, everything changed because of the effects of *Moore's Law* consistently changing the cost of computing and communications.

The most obvious change has been in the count of transistors in a single integrated circuit. Figure 11 shows the transistor count over time since 1970.

The latest production chips at the end of 2024 were the Apple M3, a 3nm chip with up to 92 billion transistors. With perhaps the possible exception of powering AI infrastructure, these days processing capability is an abundant and cheap resource.

This continual refinement of integrated-circuit production techniques affects the size and unit cost of storage (Figure 12). While the speed of memory has been relatively constant for more than a decade, the unit cost of storage has been dropping exponentially for many decades. Storage is also an abundant resource.

These changes in the capabilities of processing have also profoundly affected communications costs and capacities. The constraining factor in fibre communications systems is the capabilities of the digital signal processors and the modulators. As silicon capabilities improve, it's possible to improve the signal-processing capabilities of transmitters and receivers, which allows for a greater capacity per wavelength on a fibre circuit (Figure 13).

Figure 12: Computer Memory and Storage Unit Costs over Time^[13].

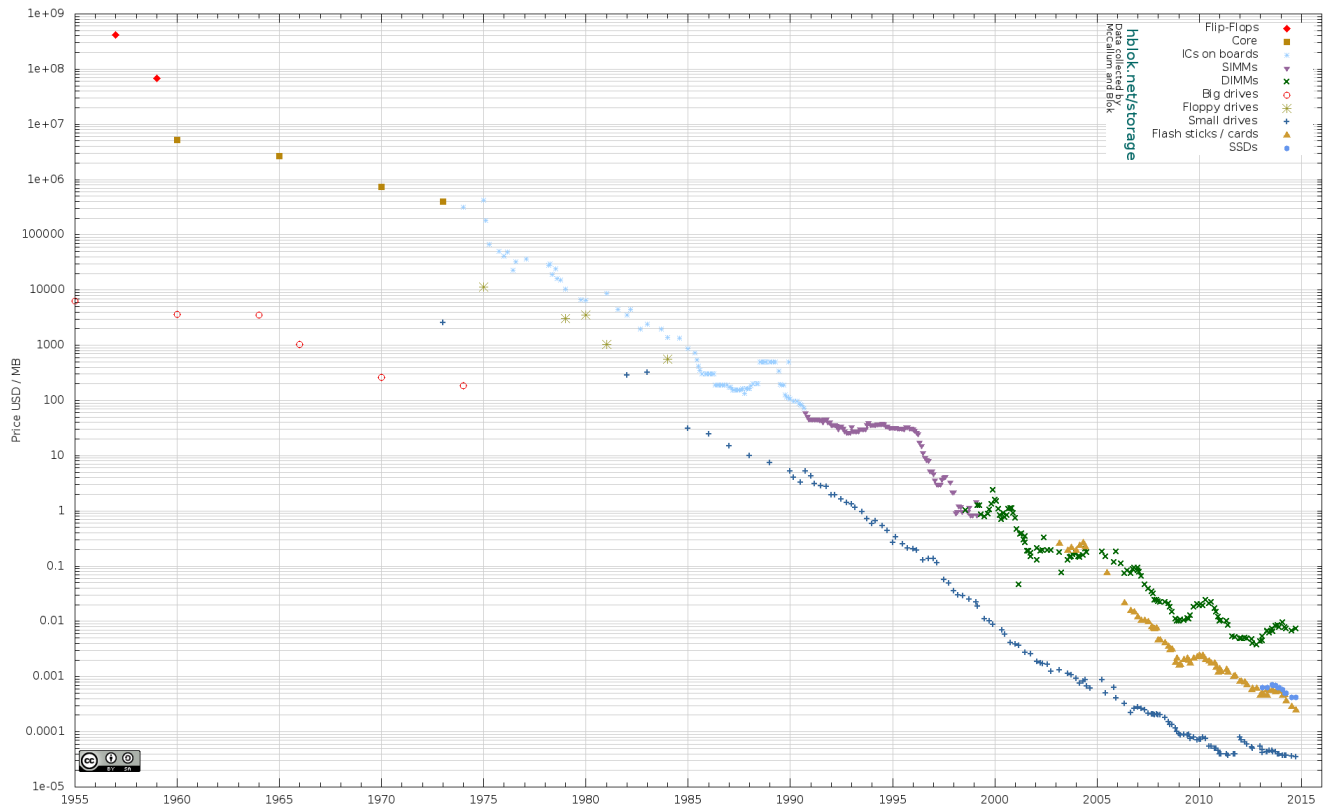
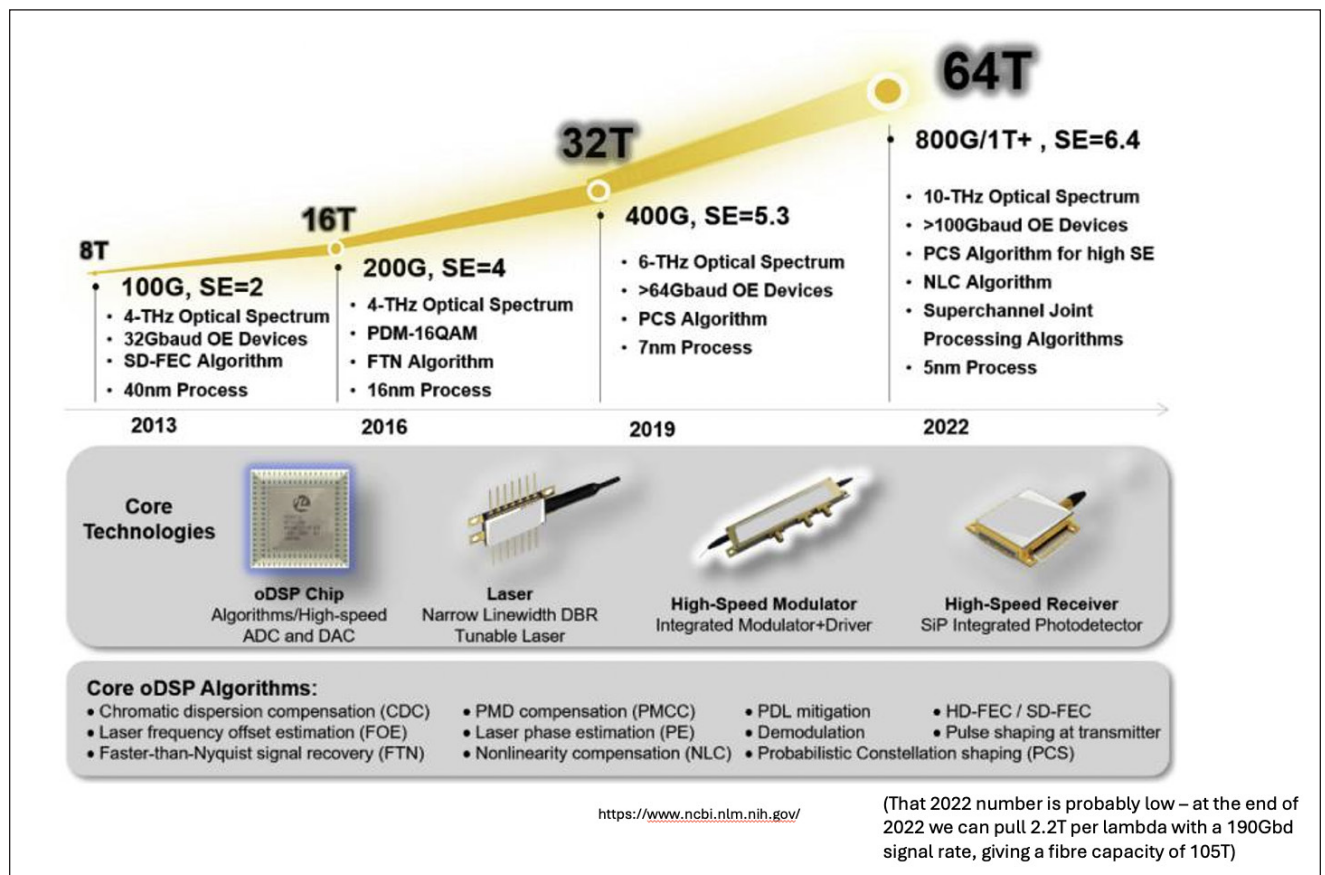


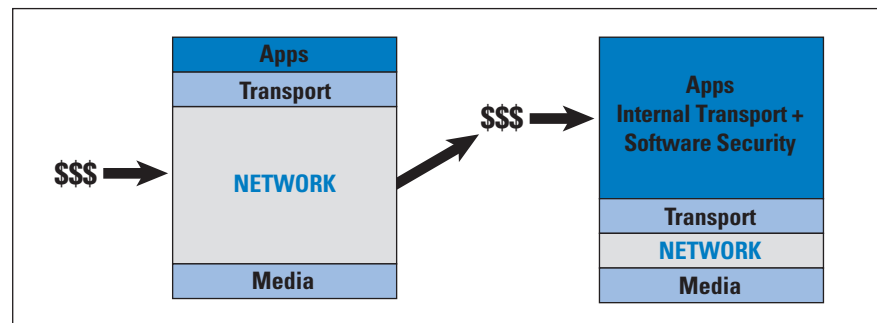
Figure 13: Fibre Capacity over Time^[14].



In addition to being bigger and faster, this environment of abundant communications, processing, and storage capacity is operating in an industry that enjoys significant economies of scale. And much of this environment is funded by capitalising a collective asset that is infeasible to capitalise individually, namely the advertisement market. The result of these changes is that a former luxury service accessible to just a few has been transformed into an affordable mass-market commodity service available to all.

However, it's more than just bigger, faster, and cheaper. This shift into abundance of basic inputs for the digital environment has changed the economics of the Internet as well. The role of the network as the arbiter of the scarce resource of communication capability has dissipated. In response, the economic focus of the Internet economy has shifted up the protocol stack to the level of applications and services (Figure 14).

Figure 14: The Transformation of the Network Economy.



Now let's return to the situation of the transition to IPv6. It is left to networks and network operators to make the investments to switch to a dual-stack platform initially (and then ultimately to remove support for IPv4). But this change is really not visible, or even crucial, to the content or service world. If IPv4 and NATs perform the carriage function adequately, then there is no motivation for the content and service operators to pay a network a premium to have a dual-stack platform.

It's domain names that operate as service identifiers, it's domain names that underpin the users' tests of authenticity of the online service, and it's the DNS that increasingly is used to steer users to the "best" service-delivery point for content or service. From this perspective addresses, IPv4 or IPv6, are not the critical resource for a service and its users. The "currency" of this form of CDN networking is *names*.

So where are we in 2025? Today's public Internet is largely a service-delivery network using CDNs to push content and service as close to the user as possible. The multiplexing of multiple services onto underlying service platforms is an application-level function tied largely to TLS and service selection using the SNI field of the TLS handshake. We use the DNS to perform "closest match" service platform selection. It's the objective of a CDN to directly attach to the access networks where its users are located, and the result is a BGP routing table inside the CDN with an average *Autonomous System (AS) Path Length* that is intended to converge to 1!

From this respect the DNS has supplanted the role of routing! We may not route “names” in today’s Internet, but it is certainly operating in a way that is largely isomorphic to such a named data network.

This architectural change has a few additional implications for the Internet. TLS, like it or not (and there is much to criticise about the robustness of TLS), is the sole underpinning of authenticity in the Internet. *Domain Name System Security Extensions* (DNSSEC) has not gathered much momentum to date. The protocol is too complex, too fragile, and just too slow to use for most services and their users. Some value its benefits highly enough that they are prepared to live with its shortcomings, but that’s not the case for most name holders and most users, and no amount of passionate exhortations about DNSSEC will change this situation! It supports the view that it’s not the mapping of a name to an IP address that’s critical. What is critical is that the named service can demonstrate that it operated by the owner of the name. Secondly, *Resource Public Key Infrastructure* (RPKI), the framework for securing information being passed in the BGP routing protocol, is really not all that useful in a service network where there is no routing!

The implication of these observations is that the transition to IPv6 is progressing very slowly not because this industry is chronically stupid or short-sighted—something else is going on here. IPv6 alone is not critical to a large set of end-user service-delivery environments. We’ve been able to take a 1980’s address-based architecture and scale it more than a billion-fold by altering the core reliance on distinguisher tokens from addresses to names. There was no real lasting benefit in trying to leap across to just another 1980’s address-based architecture (with only a few annoyingly stupid differences, apart from longer addresses!).

Where are we heading in the longer term? We are pushing everything, including value itself, out of the network and over to applications. Transmission infrastructure is becoming an abundant commodity. Network-sharing technology (multiplexing) is decreasingly relevant. We have so much network and computing resources available to us that we no longer have to take consumers to service-delivery points. Instead, we are taking services towards consumers and using the content frameworks to replicate servers and services. With so much computing and storage, the application is becoming the service, rather than just a window to a remotely operated service.

If that’s the case, then will networks matter anymore? The last couple of decades have seen us stripping out network-centric functionality and replacing it with an undistinguished commodity packet-transport medium. It’s fast and cheap, but it’s up to applications to overlay this common basic service with their own requirements. As we push these additional functions out to the edge and ultimately off the network altogether, we are left with simple dumb packet pipes!

You could argue that this situation is nothing new, and it's a continuation of the disruption that the Internet itself brought to bear on the predecessor telephone network infrastructure. The Internet architecture shifted functionality out of the core of the network and replaced synchronous real-time end-to-end virtual circuits with an extremely basic data packet-delivery service where networks were permitted to drop, duplicate, reorder, and re-time these packets in flight across the network.

It was left to the control functions that were embedded in the attached devices (such as the TCP protocol, for example) to create a functional, reliable, end-to-end communications service model. Internet hosts valued a network only to the level of a basic (and imperfect) packet-delivery service. Clients of a network were unwilling to pay a price premium for network-level services that were already being provided by the edge devices.

The result is a diminished network, dramatically reduced in both role and value. This diminished role impairs network operators to raise additional revenue through augmented services, whether it's through variable service responses through *Quality of Service* (QoS) responses or even as basic as IPv6 protocol support.

At this point it's useful to ask: What "defines" the Internet? Is the classic response, namely: "A common shared transmission fabric, a common suite of protocols, and a common protocol address pool" still relevant these days? Or is today's network more like: "A disparate collection of services that share common referential mechanisms using a common name space?"

When we think about what's important to the Internet these days, is the choice of endpoint protocol addressing really important? Is universal unique endpoint addressing a 1980's concept whose time has come and gone? If network transactions are localised, then what is the residual role of unique global endpoint addressing for clients or services? And if we cannot find a role for unique endpoint addressing, then why should we bother? Who decides when to drop this concept? Is this a market function, so that a network that uses local addressing can operate from an even lower cost base to gain a competitive market edge? Or are carriage services so cheap already that the relative benefits in discarding the last vestiges of unique global addresses are so small that it's just not worth bothering about?

And while we ponder such questions, what is the role of referential frameworks in networks? Without a common referential space, how do we usefully communicate? What do we mean by "common" when we think about referential frameworks? How can we join the "fuzzy" human language spaces with the tightly constrained deterministic computer-based symbol spaces?

Certainly, there is much to think about here!

And where does this situation leave the transition to IPv6?

I suspect that the dual-stack world we're in is a world we will be stuck in for quite some time. There seems to be no appetite to resolve this situation by completing the transition any time soon, and absolutely no desire to back out and revert to a IPv4-only network. We are here now, caught in a partial state of transition to IPv6 that is taking on an unfortunate air of permanence! And as the preponderance of value in this environment continues to move up the protocol stack into service, content, and today generative content in the guise of AI, there is little continued capacity to place collective attention on questions that have been left unresolved for decades.

It may well be that the question of when this IPv6 transition will end is a question that engenders decreasing levels of interest and attention in line with the larger picture of the decreasing relative economic value of the answer! Silicon abundance has enabled a few select content and service operators to privatise much of the former public communications platform, and in so doing they have managed to shrink the public Internet to a set of margins at the edges. That reality implies that the answer to the IPv6 transition question may soon be: "Who cares anyway?"

Disclaimer

The views in this article do not necessarily represent the views or positions of the Asia Pacific Network Information Centre.

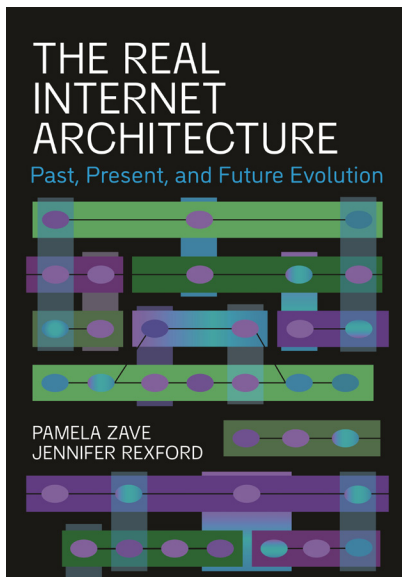
References and Further Reading

- [0] APNIC Labs Measurements and Data:
<https://labs.apnic.net/measurements/>
- [1] Steve Deering and Robert Hinden, "Internet Protocol, Version 6 (IPv6) Specification," RFC 2460, December 1998.
- [2] Phillip Gross and Philip Almquist, "IESG Deliberations on Routing and Addressing," RFC 1380, November 1992.
- [3] Paul Tsuchiya and Tony Eng, "Extending the IP Internet Through Address Reuse," ACM SIGCOMM *Computer Communications Review*, Volume 23, Issue 1, January 1993.
- [4] John Laugesen and Yufei Yuan, "What Factors Contributed to the Success of Apple's iPhone?," *Ninth International Conference on Mobile Business and 2010 Ninth Global Mobility Roundtable (ICMB-GMR)*, Athens, Greece, 2010, pp. 91–99, DOI: 10.1109/ICMB-GMR.2010.63.
- [5] Brian Carpenter and Keith Moore, "Connection of IPv6 Domains via IPv4 Clouds," RFC 3056, February 2001.
- [6] Christian Huitema, "Teredo: Tunneling IPv6 over UDP through Network Address Translations (NATs)," RFC 4380, February 2006.

- [7] Geoff Huston, “Stacking it Up,” presentation to the IPv6 Operations Working Group, IETF 80, March 2011,
<https://www.potaroo.net/presentations/2011-03-31-dualstack.pdf>
- [8] Tim Dierks and Christopher Allen, “The TLS Protocol Version 1.0,” RFC 2246, January 1999.
- [9] Dan Wing and Andrew Yourtchenko, “Happy Eyeballs: Success with Dual-Stack Hosts,” RFC 6555, April 2012.
- [10] Frank Solensky’s Presentation at IETF 18:
<https://www.ietf.org/proceedings/18.pdf>
- [11] Internet Society *Pulse*: <https://pulse.internetsociety.org/>
- [12] Transistor Count over Time:
<https://assets.ourworldindata.org/uploads/2020/11/Transistor-Count-over-time.png>
- [13] Computer Memory and Storage Unit Costs over Time:
http://aiimpacts.org/wp-content/uploads/2015/07/storage_memory_prices_large-_hblok.net_.png
- [14] Fibre Capacity over Time: <https://www.ncbi.nlm.nih.gov/>
- [15] Geoff Huston, “Anatomy: A Look inside Network Address Translators,” *The Internet Protocol Journal*, Volume 7, No. 3, September 2004.

GEOFF HUSTON AM, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

Book Review



The Real Internet Architecture: Past, Present, and Future Evolution, by Pamela Zave and Jennifer Rexford, Princeton University Press, ISBN 9780691255804, June 2024.

The goal of this book is to present a better way to describe the architecture of various networks, most notably the Internet. The initial portions of the book observe several deficiencies in how we both teach and talk about *networks* today. Most networking courses teach the bits and packets of the Internet, without a unifying framework. Networking papers suffer from “the lack of precise and consistent terminology.”

The authors present an alternative architectural model, in the hopes it will help with these issues. The core building block of the model is a network, where a network is defined as having members (hardware and software dedicated to participating in the network), names for members and groups of members, links, topology, a single administrator, and the capacity to route sessions of information. Networks offer users a simple service, namely the ability to send and receive information.

Networks can be connected three different ways:

- *Bridged Networks* are peers. Obviously IEEE 802 Ethernet networks are bridged. In this model, so too are IP networks run by different administrators. The Internet is a bridged (versus routed) network.
- *Layered Networks* put one network on top of another. A *Virtual Private Network* (VPN) is its own network and is placed on top of a set of bridged IP networks. IP networks are layered on bridged link-layer networks. As the examples illustrate, a network may be layered across multiple underlying (bridged) networks.
- Finally, there’s the case where a member in a VPN is engaged in a session with a member on a bridged IP network at the layer below the VPN. For that case, the authors repurpose the term *subduction* to describe interactions that cross layers.

This brief description summarizes a much richer conceptual framework in the book, but one can see that a large set of complex network interactions are simplified by putting a box around large chunks, declaring those chunks individual networks (a VPN, an Ethernet, etc.), and then using bridging, layering, and subduction to put them together.

In most cases, the simplification is a relief. The complex interactions among a 5G network, VPNs, and the larger Internet to serve a web-request on my phone become conceptually simpler. Similarly, tenant networks in data centers feel more tractable.

But sometimes the model falters. Treating a *Hypertext Transfer Protocol* (HTTP) session as its own network, as the book does, with a single point-to-point link as its topology, feels simplistic. It is also not at all clear who the single administrator of the HTTP session is (also an issue for some other networks described).

There are also some missed opportunities. I would suggest, in the authors' model, that there are subductive control protocols, of which the *Address Resolution Protocol* (ARP) is probably the most notable, and that these protocols present distinct challenges not encountered by protocols that stay within their network box.

The book is a thought-provoking read. I would be surprised if it managed to persuade many instructors to teach using its paradigms. I think a full textbook would be required to make that happen. At the same time, I think the notion of subduction as a way to reference the challenge of cross-boundary protocols may well catch on.

—Craig Partridge
craig@tereschau.net

Read Any Good Books Lately?

Then why not share your thoughts with the readers of IPJ? We accept reviews of new titles, as well as some of the “networking classics.” In some cases, we may be able to get a publisher to send you a book for review if you don’t have access to it. For more information, contact us at ipj@protocoljournal.org

Check your Subscription Details!

Make sure that both your postal and e-mail addresses are up-to-date since these are the only methods by which we can contact you. If you see the words “Invalid E-mail” on your printed copy, this means that we have been unable to contact you through the e-mail address on file. If this is the case, please contact us at ipj@protocoljournal.org with your new information. The subscription portal is located here: <https://www.ipjsubscription.org/>

JPNIC RPKI Guidelines Released

The *Japan Network Information Center* (JPNIC) has released a set of guidelines^[1] aimed at mitigating unauthorized routing incidents on the Internet using *Resource Public Key Infrastructure* (RPKI) *Route Origin Authorizations* (ROA). These guidelines provide technical and operational recommendations to enhance the security and reliability of Internet routing. The objective of these guidelines is to promote the adoption of RPKI-based security measures. Targeting a broad audience that includes both managerial and engineering professionals in the ISP and network operations sectors, the document offers a structured approach to implementing and maintaining RPKI.

Developed with inputs from the *Japanese Network Operators Group* (JANOG), research from the Ministry of Internal Affairs and Communications cybersecurity initiatives, and expert consultations, the guidelines offer practical insights based on real-world deployment experiences.

The guidelines cover both organizational and technical aspects of RPKI implementation. They explain the business risks associated with unauthorized routes and highlight the importance of adopting RPKI to mitigate these threats. By understanding these risks, decision-makers can justify investment in RPKI and align their security strategies with industry best practices. For network operators, the guidelines offer step-by-step instructions on creating ROAs and deploying *Route Origin Validation* (ROV). These measures ensure that only legitimate route announcements are propagated, reducing the risk of route hijacking and improving overall network security.

The guidelines also outline role-based measures for different types of network operators. IP holders are required to create ROAs and maintain consistency between their ROA records and routing information to prevent discrepancies. *Autonomous System* (AS) operators are encouraged to implement ROV to filter out invalid routes, strengthening the security of the global routing system. The guidelines include real-world configuration examples for routers and outline security measures for BGP beyond RPKI, ensuring that operators have practical resources to facilitate implementation.

Version 1 of the guidelines is available now in Japanese (translatable) in web and PDF formats and is supplemented with practical configuration examples for ROV deployment on routers. JPNIC plans to update the guidelines regularly in collaboration with experts to incorporate evolving best practices and emerging threats and has also developed an online tool, *rov-check*^[2], which allows network operators to verify whether their networks are effectively protected by ROV.

[1] <https://www.nic.ad.jp/ja/rpki/guideline/>

[2] <https://rov-check.nic.ad.jp/en>

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Kjetil Aas	Ilia Bromberg	Freek Dijkstra	Rodney Gehrke	Nils Johansson
Fabrizio Accatino	Lukasz Bromirski	Geert Van Dijk	Radu Cristian Gheorghiu	Brian Johnson
Michael Achola	Václav Brožík	David Dillow	Greg Giessow	Curtis Johnson
Martin Adkins	Christophe Brun	Richard Dodsworth	John Gilbert	Don Johnson
Melchior Aelmans	Gareth Bryan	Ernesto Doelling	Serge Van Ginderachter	Richard Johnson
Christopher Affleck	Ron Buchalski	Michael Dolan	Greg Goddard	Jim Johnston
Scott Aitken	Paul Buchanan	Eugene Doroniuk	Tiago Goncalves	Jose Enrique Diaz Jolly
Jacobus Akkerhuis	Stefan Buckmann	Michael Dragone	Ron Goodheart	Jonatan Jonasson
Antonio Cuñat Alario	Caner Budakoglu	Joshua Dreier	Octavio Alfageme	Daniel Jones
William Allaire	Darrell Budic	Lutz Drink	Gorostiaga	Gary Jones
Nicola Altan	BugWorks	Aaron Dudek	Barry Greene	Jerry Jones
Shane Amante	Scott Burleigh	Dmitriy Dudko	Jeffrey Greene	Michael Jones
Marcelo do Amaral	Chad Burnham	Andrew Dul	Richard Gregor	Amar Joshi
Matteo D'Ambrosio	Randy Bush	Joan Marc Riera	Martijn Groenleer	Javier Juan
Selva Anandavel	Colin Butcher	Duocastella	Geert Jan de Groot	David Jump
Jens Andersson	Jon Harald Bøvre	Pedro Duque	Ólafur Guðmundsson	Anders Marius Jørgensen
Danish Ansari	Olivier Cahagne	Holger Durer	Christopher Gumez	Merike Kaeo
Finn Arildsen	Antoine Camerlo	Karlheinz Dölger	Gulf Coast Shots	Andrew Kaiser
Tim Armstrong	Tracy Camp	Mark Eanes	Sheryll de Guzman	Vladislav Kalinovsky
Richard Artes	Brian Candler	Andrew Edwards	Rex Hale	Naoki Kambe
Michael Aschwanden	Fabio Caneparo	Peter Robert Egli	Jason Hall	Akbar Kara
David Atkins	Roberto Canonico	George Ehlers	James Hamilton	Christos Karayiannis
Jac Backus	David Cardwell	Peter Eisses	Darow Han	Daniel Karrenberg
Jaime Badua	Richard Carrara	Torbjörn Eklöv	Handy Networks LLC	David Kekar
Bent Bagger	John Cavanaugh	Jacobus Gerrit Elsenaar	Stephen Hanna	Stuart Kendrick
Eric Baker	Lj Cemerar	Y Ertur	Martin Hannigan	Robert Kent
Fred Baker	Dave Chapman	ERNW GmbH	John Hardin	Thomas Kernen
Santosh Balagopalan	Stefanos Charchalakakis	ESdatCo	David Harper	Jithin Kesavan
William Baltas	Molly Cheam	Steve Esquivel	Edward Hauser	Jubal Kessler
David Bandinelli	Christof Chen	Jay Etchings	David Hauweele	Shan Ali Khan
A C Barber	Pierluigi Checchi	Mikhail Evstiounin	Marilyn Hay	Nabeel Khatri
Benjamin Barkin-Wilkins	Greg Chisholm	Bill Fenner	Headcrafts SRLS	Dae Young Kim
Ryan Barnes	David Chosrova	Paul Ferguson	Hidde van der Heide	William W. H. Kimandu
Feras Batainah	Marcin Cieslak	Ricardo Ferreira	Johan Helsingius	John King
Michael Bazarewsky	Lauris Cikovskis	Kent Fichtner	Robert Hinden	Russell Kirk
David Belson	Brad Clark	Ulrich N Fierz	Michael Hippert	Gary Klesk
Richard Bennett	Narelle Clark	Armin Fisslthaler	Damien Holloway	Anthony Klopp
Matthew Best	Horst Clausen	Michael Fiumano	Alain Van Hoof	Henry Kluge
Hidde Beumer	James Cliver	The Flirble Organisation	Edward Hotard	Michael Kluk
Pier Paolo Biagi	Guido Coenders	Jean-Pierre Forcioli	Bill Huber	Andrew Koch
Arturo Bianchi	Robert Collet	Gary Ford	Hagen Hultzs	Ia Kochiashvili
John Bigrow	Joseph Connolly	Susan Forney	Kauto Huopio	Carsten Koempe
Orvar Ari Bjarnason	Steve Corbató	Christopher Forsyth	Asbjørn Højmark	Richard Koene
Tyson Blanchard	Brian Courtney	Andrew Fox	Kevin Iddles	Alexader Kogan
Axel Boeger	Beth and Steve Crocker	Craig Fox	Mika Ilvesmaki	Matthijs Koot
Keith Bogart	Dave Crocker	Fausto Franceschini	Karsten Iwen	Antonin Kral
Mirko Bonadei	Kevin Croes	Erik Fredriksson	Joseph Jackson	Robert Krejčí
Roberto Bonalumi	John Curran	Valerie Fronczak	David Jaffe	John Kristoff
Lolke Boonstra	André Danthine	Tomislav Futivic	Ashford Jaggernaut	Terje Krogdahl
Cente Cornelis Boot	Morgan Davis	Laurence Gagliani	Thomas Jalkanen	Bobby Krupczak
Julie Bottorff Photography	Jeff Day	Edward Gallagher	Jozef Janitor	Murray Kuchera
Gerry Boudreaux	Nicholas Dean	Andrew Gallo	Martijn Jansen	Warren Kumari
Leen de Braal	Fernando Saldana	Chris Gamboni	John Jarvis	George Kuo
Stephen Bradley	Del Castillo	Xosé Bravo Garcia	Dennis Jennings	Dirk Kurfuerst
Kevin Breit	Rodolfo Delgado-Bueno	Oswaldo Gazzaniga	Edward Jennings	Mathias Körber
Thomas Bridge	Julien Dhallenne	Kevin Gee	Aart Jochem	Darrell Lack

Andrew Lamb	Carsten Melberg	David Phelan	Peter Schoo	Peter Tomsu Fine Art
Richard Lamb	Kevin Menezes	Harald Pilz	Dan Schrenk	Photography
Yan Landriault	Bart Jan Menkveld	Derrell Piper	Richard Schultz	Joseph Toste
Edwin Lang	Sean Mentzer	Rob Pirnie	Timothy Schwab	Rey Tucker
Sig Lange	Eduard Metz	Jorge Ivan Pincay	Roger Schwartz	Sandro Tumini
Markus Langenmair	William Mills	Ponce	SeenThere	Angelo Turetta
Fred Langham	David Millsom	Marc Vives Piza	Scott Seifel	Brian William Turnbow
Tracy LaQuey Parker	Desiree Miloshevic	Victoria Poncini	Paul Selkirk	Michael Turzanski
Christian de Larrinaga	Joost van der Minnen	Blahoslav Popela	Andre Serralheiro	Phil Tweedie
Alex Latzko	Thomas Mino	Andrew Potter	Yury Shefer	Steve Ulrich
Jose Antonio Lazaro	Rob Minshall	Ian Potts	Yaron Sheffer	Unitek Engineering AG
Lazaro	Wijnand Modderman-	Eduard Llull Pou	Doron Shikmoni	John Urbanek
Antonio Leding	Lenstra	Tim Pozar	Tj Shumway	Martin Urwaleck
Rick van Leeuwen	Mohammad Moghaddas	David Preston	Jeffrey Sicuranza	Bart Vanautgaerden
Simon Leinen	Charles Monson	David Raistrick	Thorsten Sideboard	Betsy Vanderpool
Anton van der Leun	Andrea Montefusco	Priyan R Rajeevan	Greipur Sigurdsson	Surendran Vangadasalam
Robert Lewis	Fernando Montenegro	Balaji Rajendran	Fillipe Cajaiba da Silva	Ramnath Vasudha
Christian Liberale	Roberto Montoya	Paul Rathbone	Andrew Simmons	Randy Veasley
Martin Lillepuu	Joel Moore	William Rawlings	Pradeep Singh	Philip Venables
Roger Lindholm	Joseph Moran	Mujtiba Raza Rizvi	Henry Sinnreich	Buddy Venne
Link Light Networks	John More	Bill Reid	Geoff Sisson	Alejandro Vennera
Art de Llanos	Maurizio Moroni	Petr Rejhon	John Sisson	Luca Ventura
Mike Lochocki	Brian Mort	Robert Remenyi	Helge Skrivervik	Scott Vermillion
Chris and Janet Lonvick	Soenke Mumm	Rodrigo Ribeiro	Terry Slattery	Tom Vest
Mario Lopez	Tariq Mustafa	Glenn Ricart	Darren Sleeth	Peter Villemoes
Sergio Loreti	Stuart Nadin	Justin Richards	Richard Smit	Vista Global Coaching &
Eric Louie	Michel Nakhla	Rafael Riera	Bob Smith	Consulting
Adam Loveless	Mazdak Rajabi Nasab	Mark Risinger	Courtney Smith	Dario Vitali
Josh Lowe	Krishna Natarajan	Fernando Robayo	Eric Smith	Marc Vives
Guillermo a Loyola	Naveen Nathan	Michael Roberts	Mark Smith	Rüdiger Volk
Hannes Lubich	Darryl Newman	Gregory Robinson	Tim Sneddon	Jeffrey Wagner
Dan Lynch	Mai Nguyen	Ron Rockrohr	Craig Snell	Don Wahl
David MacDuffie	Thomas Nikolajsen	Graziano G Rodegari	Job Snijders	Michael L Wahrman
Sanya Madan	Paul Nikolich	Carlos Rodrigues	Ronald Solano	Lakhinder Walia
Miroslav Madić	Travis Northrup	Magnus Romedahl	Asit Som	Laurence Walker
Alexis Madriz	Marijana Novakovic	Lex Van Roon	Ignacio Soto Campos	Randy Watts
Carl Malamud	David Oates	Marshall Rose	Evandro Sousa	Andrew Webster
Jonathan Maldonado	Ovidiu Obersterescu	Alessandra Rosi	Peter Spekrijse	Jd Wegner
Michael Malik	Jim Oplotnik	David Ross	Thayumanavan Sridhar	Tim Weil
Tarmo Marners	Tim O'Brien	William Ross	Paul Stancik	Westmoreland
Yogesh Mangar	Mike O'Connor	Boudhayan	Ralf Stempfner	Engineering Inc.
John Mann	Mike O'Dell	Roychowdhury	Matthew Stenberg	Rick Wesson
Bill Manning	John O'Neill	Carlos Rubio	Martin Štěpánek	Peter Whimp
Diego Mansilla	Carl Örne	Rainer Rudigier	Adrian Stevens	Russ White
Harold March	Packet Consulting Limited	Timo Ruiter	Clinton Stevens	Jurrien Wijlhuizen
Vincent Marchand	Carlos Astor Araujo	RustedMusic	John Streck	Joseph Williams
Normando Marcolongo	Palmeira	Babak Saberi	Martin Streule	Derick Winkworth
Gabriel Marroquin	Gordon Palmer	George Sadowsky	David Strom	Pindar Wong
David Martin	Alexis Panagopoulos	Scott Sandefur	Colin Strutt	Brian Woods
Jim Martin	Gaurav Panwar	Sachin Sapkal	Viktor Sudakov	Makarand Yerawadekar
Ruben Tripiana Martin	Chris Parker	Arturas Satkovskis	Edward-W. Suor	Phillip Yialeloglou
Timothy Martin	Alex Parkinson	PS Saunders	Vincent Surillo	Janko Zavernik
Carles Mateu	Craig Partridge	Richard Savoy	Terence Charles Sweetser	Bernd Zeimet
Juan Jose Marin Martinez	Manuel Uruena Pascual	John Sayer	T2Group	Muhammad Ziad
Ioan Maxim	Ricardo Patara	Phil Scarr	Roman Tarasov	Ziauddin
David Mazel	Dipesh Patel	Gianpaolo Scassellati	David Theese	Tom Zingale
Miles McCredie	Dan Paynter	Elizabeth Scheid	Rabbi Rob and	Matteo Zovi
Gavin McCullagh	Leif-Eric Pedersen	Jeroen Van Ingen	Lauren Thomas	Jose Zumalave
Brian McCullough	Rui Sao Pedro	Schenau	Douglas Thompson	Romeo Zwart
Joe McEachern	Juan Pena	Carsten Scherb	Kerry Thompson	廖明沂.
Alexander McKenzie	Luis Javier Perez	Ernest Schirmer	Lorin J Thompson	
Jay McMaster	Chris Perkins	Benson Schliesser	Jerome Tissieres	
Mark Mc Nicholas	Michael Petry	Philip Schneck	Fabrizio Tivano	
Olaf Mehlberg	Alexander Peuchert	James Schneider		

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Follow us on X and Facebook



@protocoljournal



<https://www.facebook.com/newipj>

Supporters and Sponsors

Supporters



Internet
Society



Diamond Sponsors

Your logo here!

Ruby Sponsors



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

John Crain, Senior Vice President and Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Merike Kaeo, Founder and vCISO
Double Shot Security

Olaf Kolkman, Principal – Internet Technology, Policy, and Advocacy
The Internet Society

Dr. Jun Murai, Founder, WIDE Project
Distinguished Professor, Keio University
Co-Director, Keio University Cyber Civilization Research Center, Japan

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.



The Internet Protocol Journal

August 2025

Volume 28, Number 2

A Quarterly Technical Publication for
Internet and Intranet Professionals

FROM THE EDITOR

In This Issue

From the Editor	1
ShowNet 2024 Highlights	2
The Root of the DNS.....	14
Letters to The Editor	31
In Memoriam: Dave Täht.....	34
In Memoriam: Fred Baker.....	35
Fragments.....	37
Thank You!	40
Call for Papers.....	42
Supporters and Sponsors.....	43

You can download IPJ
back issues and find
subscription information at:
www.protocoljournal.org

ISSN 1944-1134

The *TCP/IP Interoperability Conference*—later renamed *Interop*—began as a small workshop in August 1986. It quickly grew in scope to incorporate tutorials, and by 1988 an exhibition network connected 51 exhibitors to each other and to the global Internet. This network was designed and deployed by a group of volunteers, and it became the proving ground for many emerging technologies. In 1994, Interop added Tokyo to its international venues, where 30 years later the conference and exhibition attracts more than 120,000 visitors annually. Following an article in our October 2024 issue describing the history and evolution of the Interop show network, and a second article detailing the Tokyo *ShowNet* in our previous issue, we now bring you the final installment in this series with an article that highlights some of the technology demonstrations performed during the 2024 Interop Tokyo event. The article is by Ryo Nakamura, Haruki Nakamura, Kazuya Okada, and Ryosuke Kato.

The *Domain Name System* (DNS) is one of the core components of the Internet. We have covered many aspects of the DNS over the years, but we have not discussed the *root server system* since an article in Volume 20, No. 2, June 2017. In this issue, Geoff Huston returns to the topic with a detailed tutorial and analysis of today's DNS root server operations.

We always welcome feedback and suggestions on any aspect of this journal. Included in this issue are two Letters to the Editor in response to the IPv6 Transition article in our May 2025 edition. If you'd like to get in touch, send your comments to: ipj@protocoljournal.org.

In late June, I attended the *Internet Governance Forum* (IGF) meeting in Lillestrøm, Norway. Lillestrøm happens to be the place where I attended high school. During my summer breaks, I worked at the nearby *Norwegian Defence Research Establishment* (NDRE), which had one of the first connections to the ARPANET starting in 1973. The IGF exhibition area had a series of posters highlighting the evolution of the Internet in Norway. I was pleased to see that the first poster featured Internet Hall of Fame Inductees Pål Spilling and Yngvar Lundh, my former managers at NDRE. See page 39.

—Ole J. Jacobsen, Editor and Publisher
ole@protocoljournal.org

Technology Highlights of ShowNet 2024

by Ryo Nakamura, Haruki Nakamura, Kazuya Okada, and Ryosuke Kato

Interop Tokyo 2024 was held from June 12 to 14 in the *Makuhari Messe* exhibition halls. With 542 organizations exhibiting and 124,482 visitors attending the exhibition, Interop Tokyo is one of the largest IT shows in Japan. *ShowNet*^[0], the large demonstration network for Interop Tokyo, was also built at the venue. In 2024, the ShowNet comprised approximately 2,300 products and services in more than 20 full-height racks built and operated by 650 engineers including 31 *Network Operations Center* (NOC) team members, 38 volunteer members, and 581 engineers from vendors who contributed their products to ShowNet. These engineers gathered at Makuhari Messe on May 31 and built the network in two weeks. Figure 1 is a picture of the second day of the ShowNet construction in 2024.

Figure 1: A snapshot of ShowNet under construction at Makuhari Messe on June 1, 2024.



The fundamental role of ShowNet is to provide network connectivity to Interop exhibitors and visitors. Furthermore, ShowNet conducts various experiments and demonstrations of new protocols, technologies, and products while serving user traffic.

In 2024, ShowNet featured the following technical topics in each field:

- *Facility*: A high-density *Main Distribution Frame* (MDF) with SN connector-based patch panels^[1].
- *Optical Transport*: Multi-vendor optical transport network with emerging optics such as 400GBASE-ZR+ and XR Optics.
- *Backbone Network*: An SRv6 uSID-based backbone network and *Ethernet VPN* (EVPN) and *Virtual eXtensible Local Area Network* (VXLAN) for access.
- *Data Center and Cloud*: Distributed container clusters and testing lossless networks for *Remote Direct Memory Access* (RDMA) over *Converged Ethernet* (RoCE) traffic.
- *Wireless Network*: Multi-band Wi-Fi access with Wi-Fi6E- and Wi-Fi7-capable Access Points, and multi-vendor *OpenRoaming*^[2].
- *Monitoring*: Integrated monitoring systems with various sensors and user interfaces, and experimentation of how to exploit AI for future monitoring.
- *Security*: Incorporating multiple aspects of protection and hardening such as SASE, ZTNA, EASM, and NGFW.
- *Tester*: Testing upper layers with protocol emulation for routing and penetration tests for security, and demonstrating automating test processes.
- *5G*: Multiple private 5G systems of RAN and cores with multiple vendors, and demonstration of live streaming over the 5G networks.
- *Media-over-IP*: Professional audio and media are now migrating from SDI to IP: demonstrating real-time broadcasting over IP networks.

In this article, we describe four of these topics, namely: *The Backbone Network*, *Optical Transport*, *5G*, and *Media-over-IP*.

The Backbone Network

The backbone network of ShowNet is the core of all the experiments and demonstrations. In 2024, the backbone network was composed of ten routers of nine products listed in Table 1. In addition, two containerized routers, XRd from Cisco Systems and cRPD from Juniper Networks, performed route reflectors for *Border Gateway Protocol* (BGP). With those routers, we built the backbone network based on Segment Routing while conducting SRv6 uSID interoperability tests.

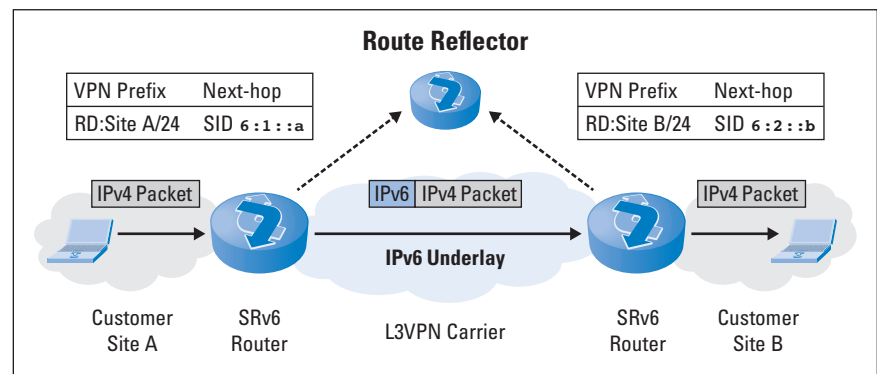
Table 1: Routers composing the backbone network of ShowNet in 2024.

Vendor	Product
Cisco Systems	Cisco 8201-32FH, Cisco 8608, NCS-57B1
Furukawa Electric	FX2
Huawei Technologies	NE8000-M4
Juniper Networks	ACX7348, MX204, MX304, PTX10002-36QDD

Segment Routing (SR)^[3] is a recent routing and forwarding paradigm that enables source routing. In SR, topological entities are represented by segments; for example, nodes, links, and adjacency. SR nodes control where packets should flow and how packets are processed by embedding a series of segments into a packet. SR has two concrete data-plane implementations: SR-MPLS leveraging *Multi-Protocol Label Switching (MPLS)* labels as *Segment Identifiers (SIDs)* and SRv6 leveraging IPv6 addresses as SIDs. An MPLS label stack encapsulating a packet indicates a SID list in SR-MPLS and IPv6 addresses in a *Segment Routing Header*^[4]—which is a new IPv6 extension header—it also indicates a SID list in the SRv6 data plane.

A major use case of SR is *Layer-3 VPN (L3VPN)*. Figure 2 illustrates a simple example of SRv6-based L3VPN. Two SRv6 routers perform Provider Edge functions for two customer sites, and exchange VPN prefixes via *Multi-Protocol BGP (MP-BGP)*. Note that the next-hops for those VPN prefixes are SRv6 SIDs: the ingress SRv6 router encapsulates packets from the customer site A to site B with IPv6 headers whose destination address is the SID (`6:2::b`) of the egress SRv6 router.

Figure 2: A simple example of SRv6-based L3VPN.

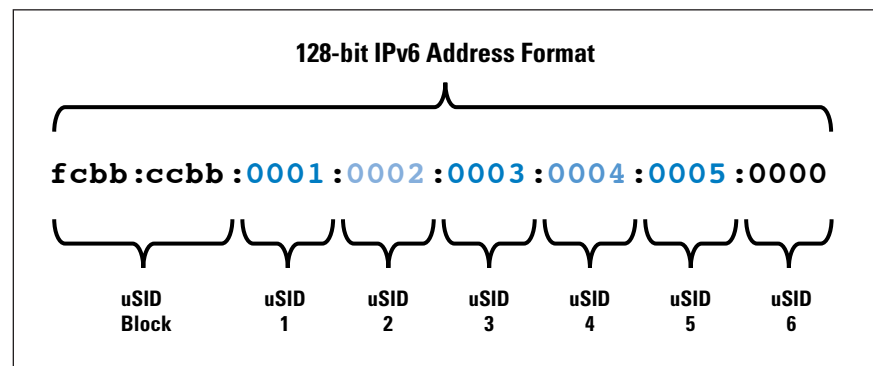


For ShowNet at Interop Tokyo, we have worked on Segment Routing continuously since 2018. In 2018 we conducted a simple and small interoperability test of the SR-MPLS and SRv6 data planes, and in 2019 we demonstrated service chaining over SRv6 with multiple vendors' products. Since 2021, we have deployed SR on the ShowNet backbone networks. The backbone network of ShowNet 2021 was composed of SR-MPLS, and we further conducted a measurement experiment on Internet latency using SR-MPLS-based *Egress Peer Engineering*, which enables steering specific egress traffic to given *External BGP (eBGP)* peers. The results of the experiment were published in a paper^[5] and in an APNIC blog post^[6]. In 2022 and 2023, the ShowNet backbone was fully SRv6-enabled, and IPv4 addresses were eliminated—interfaces of backbone links had no IP addresses configured thanks to IPv6 link-local addresses. Our chronicle with SR was summarized in a presentation at the *Asia Pacific Regional Internet Conference on Operational Technologies (APRICOT) 2024*^[7].

In 2024, a main topic in the backbone network was SRv6 micro SIDs (uSID). uSID, also known as the NEXT-C-SID flavor in^[8], is a mechanism for compressing SID lists in SRv6. A SID in the original SRv6 is a 128-bit IPv6 address, encapsulating packets with multiple SIDs. For example, traffic engineering, involves significant overhead on MTU sizes. uSID encodes multiple SIDs into a 128-bit IPv6 address format to avoid the overhead. Figure 3 illustrates a uSID structure with F3216 format^[9], which implementations must support at present.

The first 32-bit is a uSID block that all routers in an SRv6 domain share. The 16-bit blocks shown in Figure 3 are uSIDs. When an SRv6 node processes the first uSID (**fcbb:ccbb:0001:...**), the node shifts the 80 bits from the second to the last uSID 16 bits to the left and overwrites the first uSID. In other words, the new destination address of the packet is **fcbb:ccbb:0002:0003:0004:0005::**, and the packet is forwarded to the next SRv6 node that has the uSID **0002**. This new packet forwarding mechanism is currently being implemented in router products of multiple vendors, and we confirmed that uSID interoperability between the devices listed in Table 1 was successfully achieved in ShowNet 2024.

Figure 3: A uSID structure with the F3216 format.



The second topic is a demonstration for campus and enterprise networks. The “customers” of ShowNet are exhibitors connecting equipment in their booths to the network. This means that the last hop to the booths consists of several hundred UTP cables spread over the exhibition halls. Accommodating those access circuits becomes a technical demonstration of campus and enterprise networks. This year we built those access networks as L2 and L3VPN with *Ethernet VPN* (EVPN) and *Virtual eXtensible Local Area Network* (VXLAN)^[10] with campus switches from multiple vendors.

VXLAN is an Ethernet-over-IP tunneling protocol, and EVPN is a BGP-based control plane that can construct overlay fabrics using VXLAN as its data plane^[11]. EVPN-VXLAN was originally designed and introduced for data-center use; therefore, switches and routers that were intended primarily for use in data centers supported these protocols in the early days.

Over the years, recent switches for campus and enterprise networks, which are different product lines from those for data centers, have begun to support EVPN-VXLAN for campus use. Adopting Ethernet overlays for campus networks will eliminate (often fragile) spanning-tree protocols and provide scalability and resiliency by using underlying dynamic routing protocols.

The access network in ShowNet 2024 was composed of three routers and eight switches of seven models listed in Table 2. All devices exchanged EVPN routes via route reflectors, constructed a VXLAN fabric, and forwarded user traffic over the fabric. User VLANs could be extended between the switches over the IP underlay. In addition, EVPN can construct L3VPNs using EVPN Type-5 routes^[12]. We also confirmed that the EVPN Type-5 route interoperability works well with these devices.

Table 2: Routers and switches composing the access network with EVPN-VXLAN.

Vendor	Product
Cisco Systems	Catalyst 9300, Nexus 93108TC-FX
Huawei Technologies	CloudEngine S5732, NE8000 M4
Juniper Networks	EX4400, MX304, SRX4600

While SRv6 uSID and EVPN-VXLAN for access were major topics, the demonstrations were not limited to just these two. Other demonstrations and technical challenges were also conducted at the ShowNet backbone network; for example, an experiment of SRv6 over a satellite for disaster recovery, testing *Path Computation Element Protocol* (PCEP), and a total of 2 Tbps external circuits including a capacity of 1.8 Tbps provided by Open APN.

Optical Transport

The optical transport network in ShowNet multiplexes waves on fibers to optimize fiber use while showcasing products in this area. Furthermore, the optical transport network in 2024 faced challenges, including interoperability tests, and tests with other layers above Layer 2. The topics in 2024 were as follows:

- Using multi-band connections of C-band and L-band.
- Interoperability between 400GBASE-ZR+ transceivers based on OpenZR+.
- 1:N point-to-multipoint connections as defined by the *Open XR Optics Forum*.

The optical transport network in ShowNet 2024 consisted of multiple *Wavelength Division Multiplexing* (WDM) networks. One of the WDM networks used a *Reconfigurable Optical Add-Drop Multiplexer* (ROADM) with C-band and L-band wavelengths, connecting transponders and muxponders with capacities ranging from 400 to 800 Gbps.

This WDM network also provided connections of 100GBASE-LR4 and 400GBASE-FR4 to the backbone routers. In addition, we conducted an interoperability test of 400GBASE-ZR+ transceivers at ShowNet. 400GBASE-ZR+^[13] employs coherent optics that enable configuring and transmitting multiple wavelengths so that they can remove transponders. Different manufacturers provide coherent optics equipped with *Digital Signal Processing* (DSP), and we confirmed that they operated correctly in various combinations. Using this infrastructure, we also tried to transfer wavelengths directly from a carrier through the optical transport network built at ShowNet in collaboration with the carrier.

Another WDM network conducted a test of coherent 100GBASE-ZR in the QSFP28 form factor, which was developed after 400GBASE-ZR+ emerged, with ROADMs using C-band wavelengths, in addition to the interoperability of 400GBASE-ZR+ transceivers. Further, we deployed XR Optics^[14], which enables point-to-multipoint optical connections. Deploying Open XR Optics with a ROADM was the first challenge, and it was successfully completed by strong cooperation with each vendor of the transceiver, transponder, ROADM, and *Erbium-Doped Fiber Amplifier* (EDFA) at ShowNet.

5G

Private 5G networks are wholly owned and operated 5G networks that enable individual companies to possess some radio spectrum for their purposes. In Japan, private 5G networks are recognized as Local 5G. This type of private 5G and local 5G is defined as a *Standalone Non-Public Network* (SNPN) in the 3GPP standards. We have been conducting private 5G experiments in a part of ShowNet with 5G-related vendors and integrators since 2022. This year, we deployed three different private 5G networks with multiple vendors and conducted two demonstrations: live streaming in *Network Operations Center* (NOC) guided tours in the exhibition using the 5G networks to improve participants' experience and provided Internet connectivity to several exhibition booths. In addition, we designed a stable and redundant *Precision Time Protocol* (PTP)^[15] network for the 5G networks. In this demonstration, we constructed three private 5G networks that use licensed n79 spectrums in Japan. The demonstration highlighted the advantages of private 5G networks over mobile carriers' 5G services, including low latency and guaranteed access in licensed areas.

NOC guided tours in the exhibition adopted real-time video streaming with the private 5G systems for this year. On the tours, called the *ShowNet Walking Tours*, a NOC team member gives a talk about design concepts and underlying technologies for every rack. However, the areas around the racks were crowded and noisy during the exhibition, so it was difficult for tour participants to see the equipment that NOC members were describing. Furthermore, technologies and devices introduced during the tour were extensive; therefore, conveying this information clearly to the tour participants through only verbal explanations was also challenging.

To address the uncomfortable situation in the tours, we streamed the voice and live movie of the tour guide describing the racks to attendees' 5G-capable tablets and smartphones. Encoded movies and audio were transported to a decode server located in a ShowNet rack via a 5G system. Then, the decoded movie was mixed with supplemental slides and was presented on the attendees' tablets at the right time. An on-premises streaming server delivered the edited movie and audio to attendees' 5G tablets and smartphones via two different 5G networks. Figures 4 and 5 show a camera recording a tour guide describing a ShowNet rack, and the video is mixed with slides. Tour attendees watch the mixed stream, as shown in Figure 6.

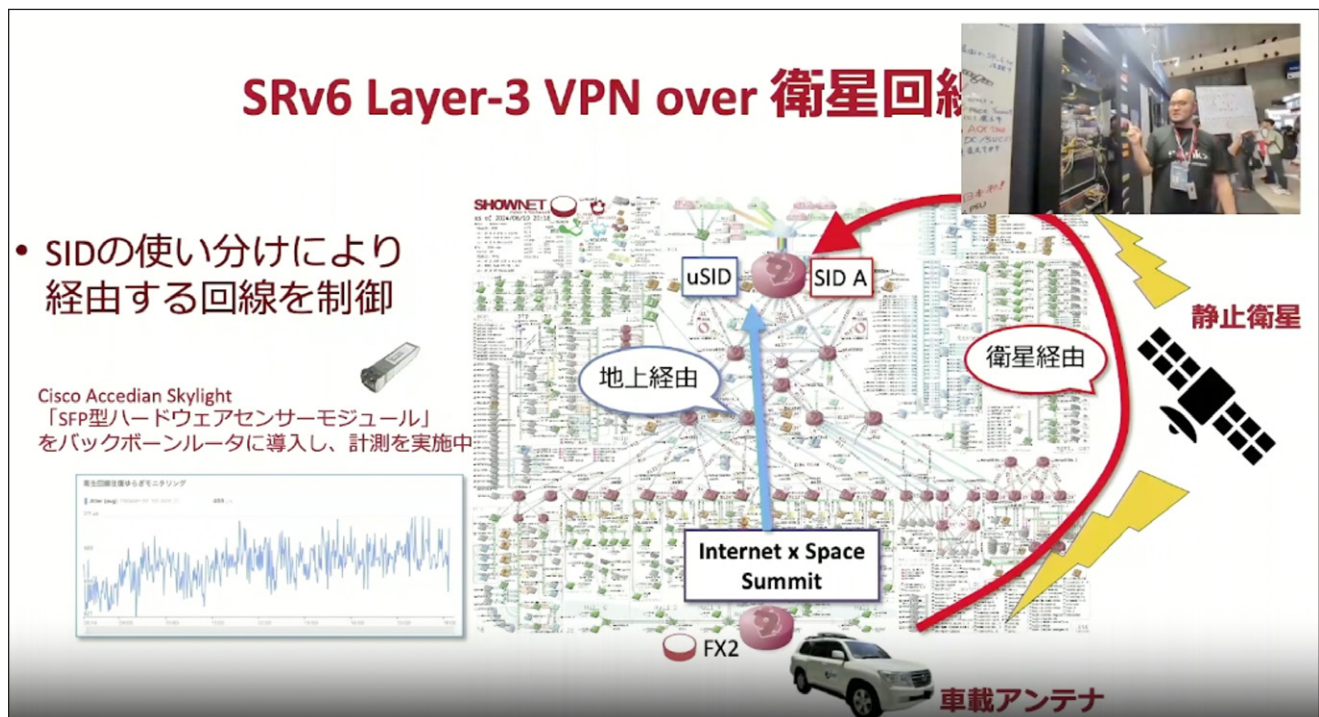
Figure 4: Live broadcasting with private 5G-enabled smartphone cameras.



Figure 5: Mixing the received video image with a slide related to what the NOC member is describing.



Figure 6: Video images delivered to tour attendees' tablets and smartphones.



This demonstration using the 5G systems provided very stable live streaming in the exhibition halls, in contrast to using Wi-Fi. Wi-Fi access was also available, but the Wi-Fi public bands of 2.4 and 5 GHz were already experiencing congestion due to massive numbers of visitors' mobile Wi-Fi devices. Therefore, latency and available network bandwidth were unstable, and it was not easy to provide guaranteed streaming. Wi-Fi 6E, which uses 6-GHz channels, still has not been congested because of the small number of capable devices. However, it is anticipated that this situation will change next year.

Media-over-IP

Professional audio and media are now migrating from *Serial Digital Interface* (SDI) cables, which have low transfer rates and high costs, to Ethernet/IP-based systems for higher transfer rates and lower costs because of the availability of commodity equipment. ShowNet has featured these media-over-IP solutions as one of the main topics since 2022. In 2024, we collaborated with broadcasters pursuing the transition to IP in broadcasting to explore the possibilities of media-over-IP networks and services for broadcasting industries; we attempted to connect and exchange media between the ShowNet booth and three geographically distributed broadcast stations over IP networks.

In the ShowNet booth, we built the *Media Operation Center* (MOC), a broadcast control room for media production and remote operation with IP-based systems. Using this MOC (Figure 7), we demonstrated real-time recording, editing, and broadcasting of a stage (Figure 8) where many sessions were held during the exhibition. This facility also supported live mixing and streaming on the tours with the 5G demonstration described previously.

Media-over-IP technologies are standardized by the *Society of Motion Picture and Television Engineers* (SMPTE)^[16], and its standards are prefixed with SMPTE. For example, the SMPTE ST 2110 series^[17] defines protocols and parameters for professional video, audio, and data-over-IP transport.

From the network viewpoint, that media traffic is *Real-time Transport Protocol* (RTP) streams over IP multicast, and media endpoints speaking the protocols require *Precision Time Protocol* (PTP) to synchronize clocks. Thus, in ShowNet 2024, we built a Layer-3 multicast network with *Open Shortest Path First Version 2* (OSPFv2) and *Protocol Independent Multicast – Sparse Mode* (PIM-SM) for the control room by using Cisco Nexus series and Huawei Cloud Engine switches. These switches are capable of PTP for broadcast profiles (SMPTE ST 2059-2). Furthermore, we configured Layer-2 VPN and Layer-3 VPN connections using VPN devices for media transmission and control between two broadcast stations in Tokyo (30 km away from the venue) and a station in Sapporo (830 km away from the venue) over the Internet. These connections established a remote production environment between the broadcast stations and the MOC booth at ShowNet.

Figure 7: The Media Operation Center at the ShowNet booth.



Figure 8: A stage presentation broadcasted by the media-over-IP systems deployed on ShowNet.



We demonstrated media production with the remote broadcast stations over IP networks during the three-day exhibition. Traffic transferred through the networks included bidirectional uncompressed video streams (SMPTE ST 2110-20, 1080i with 59.94 Hz, up to 1.3 Gbps per stream) and compressed video streams by JPEG-XS (SMPTE ST 2110-22, 1080i with 59.94 Hz, up to 200 Mbps per stream). Additionally, sensors embedded in *Small Form-factor Pluggable* (SFP) modules from Accedian were placed at a ShowNet rack and the broadcast stations to enable active monitoring by *Two-Way Active Measurement Protocol* (TWAMP) measurements. This setup allowed us to observe real-time network performance impacts on media traffic.

During the event, we collaborated with broadcasting industry members to conduct live broadcast and video production of sessions at the exhibition. Eventually, all media transport and equipment operations between the broadcast stations and the Media Operation Center at the ShowNet booth were conducted entirely over IP.

Conclusion

In this article, we introduced technology highlights from ShowNet in 2024. ShowNet covers broader aspects of networking technologies and conducts demonstrations from Layer 1 to Layer 7. Unfortunately, explanations of all the topics discussed in this article are not possible because of the amount of material it would necessitate. So, in this article we covered only four topics: the backbone network, optical transport, 5G, and media-over-IP, and briefly described these technical overviews.

ShowNet is a show in the Interop exhibition; different from ordinary networks, it is an ephemeral network built and operated for just three days. However, we do not let the show network end as just a show. Through conducting various experiments and demonstrations, as described in this article, we aim to encourage network communities in Japan, foster relationships between engineers, and contribute the knowledge and insights gained at ShowNet to society.

Acknowledgments

The design, construction, and demonstration of the ShowNet network were made possible through the collaboration of NOC team members, contributing vendors and their teams, and ShowNet Team members. We want to thank all the people involved in the ShowNet at Interop Tokyo 2024.

References and Further Reading

- [0] Takashi Tomine, Ryo Nakamura, and Ryota Motobayashi, "ShowNet at Interop Tokyo: A Continuously Evolving Demonstration Network," *The Internet Protocol Journal*, Volume 28, No. 1, May 2025.
- [1] SENKO Advance Co., Ltd. SN 1.6mm Standard Connector:
[https://www.senko.com/product/
sn-1-6mm-standard-connector/](https://www.senko.com/product/sn-1-6mm-standard-connector/)

- [2] Wireless Broadband Alliance, OpenRoaming: <https://wballiance.com/openroaming/>
- [3] Clarence Filsfils, Stefano Previdi, Les Ginsberg, Bruno Decraene, Stephane Litkowski, and Rob Shakir, “Segment Routing Architecture,” RFC 8402, July 2018.
- [4] Clarence Filsfils, Darren Dukes, Stefano Previdi, John Leddy, Satoru Matsushima, and Daniel Voyer, “IPv6 Segment Routing Header (SRH),” RFC 8754, March 2020.
- [5] Ryo Nakamura, Kazuki Shimizu, Teppei Kamata, and Cristel Pelsser, “A first measurement with BGP Egress Peer Engineering,” in *Proceedings of 23rd International Conference on Passive and Active Measurement*, PAM 2022, pages 199–215, Springer International Publishing.
- [6] Ryo Nakamura, “Measuring the potential benefit of egress traffic engineering with Segment Routing, *APNIC Blog*, March 10, 2022.
- [7] Teppei Kamata, “Segment Routing Deployments and Demonstrations at Interop Tokyo ShowNet,” Asia Pacific Regional Internet Conference on Operational Technologies (APRICOT) 2024, February 2024.
- [8] Weiqiang Cheng, Clarence Filsfils, Zhenbin Li, Bruno Decraene, and Francois Clad, “Compressed SRv6 Segment List Encoding,” RFC 9800, June 2025.
- [9] Bell Canada, “uSID Address Allocation, How to Assign SRv6 Locators to Network Nodes,” Presentation during IETF srv6 Working Group Meeting at IETF 119, March 2024.
- [10] Mallik Mahalingam, Dinesh Dutt, Kenneth Duda, Puneet Agarwal, Larry Kreeger, T. Sridhar, Mike Bursell, and Chris Wright, “Virtual eXtensible Local Area Network (VXLAN): A Framework for Overlaying Virtualized Layer 2 Networks over Layer 3 Networks,” RFC 7348, August 2014.
- [11] Ali Sajassi, John Drake, Nabil Bitar, Ravi Shekhar, Jim Uttaro, and Wim Henderickx, “A Network Virtualization Overlay Solution Using Ethernet VPN (EVPN),” RFC 8365, March 2018.
- [12] Jorge Rabadan, Wim Henderickx, John Drake, Wen Lin, and Ali Sajassi, “IP Prefix Advertisement in Ethernet VPN (EVPN),” RFC 9136, October 2021.
- [13] OpenZR+: <https://www.openzrplus.org/>
- [14] Open XR Optics Forum: <https://openxropticsforum.org/>
- [15] IEEE, “1588-2008 IEEE Standard for a Precision Clock Synchronization Protocol for Networked Measurement and Control Systems,” pages 1–300, 2008.
- [16] SMPTE The home of Media Professionals, Technologists, and Engineers: <https://www.smpte.org/>

- [17] SMPTE, “ST 2110 Suite of Standards”:
<https://www.smpte.org/standards/st2110>
- [18] David Strom, “The Interop ShowNet,” *The Internet Protocol Journal*, Volume 27, No. 3, October 2024.
- [19] Interop 2024 ShowNet concept:
<https://www.interop.jp/2024/shownet/concept/>
- [20] Interop 2024 ShowNet Brochure:
<https://www.interop.jp/2024/assets/file/arukikata.pdf>
- [21] Interop 2024 ShowNet map:
<https://www.interop.jp/2024/assets/file/e-web.pdf>
- [22] ShowNet map icons:
<https://github.com/interop-tokyo-shownet/shownet-icons>
- [23] “Behind the Scenes - Interop Tokyo 2019 ShowNet,” Interop Tokyo YouTube video:
<https://www.youtube.com/watch?v=X-JhPs1T7sc>

RYO NAKAMURA received his Ph.D. degree in Information Science and Technology from the University of Tokyo, Tokyo, Japan, in 2017. He is currently an Associate Professor at the Information Technology Center, the University of Tokyo, where he operates the university’s campus network. His research interests include networking in operating systems, network virtualization, and network operations. Since 2009, he has been involved in Interop Tokyo ShowNet as a ShowNet team member until 2011, and as a member of the NOC team from 2012 to the present. He has been primarily responsible for the backbone network of ShowNet, and he led demonstrations of SDN-related technologies from 2013 to 2017. He can be reached at:
ryo@interop-tokyo.net

HARUKI NAKAMURA received a Master’s degree from Keio University Graduate School of Media Design in 2019. He started his career as a Solutions Engineer at Cisco Systems G.K., focusing on data-center networking and computing technologies. Since 2022, he has expanded his responsibilities to include IP Media Networking and joined Interop Tokyo ShowNet as a contributor. In 2024, he served as a ShowNet NOC member in the Media-over-IP Working Group, where he led initiatives to collaborate with contributors and broadcasting companies on proof-of-concept projects for the next-generation Media-over-IP system. He can be reached at:
hanakamu@interop-tokyo.net

KAZUYA OKADA received a PhD degree in computer science from the Nara Institute of Science and Technology (NAIST), Japan, in 2014. He is currently a Principal Researcher at the InfoTech (Research Division of Information Technology), Toyota Motor Corporation, Japan. His research interests include cyber resilience and cybersecurity for connected vehicles. He has been a NOC team member of ShowNet since 2011. He can be reached at: **okada@interop-tokyo.net**

RYOSUKE KATO has been employed by BroadBand Tower, a Japanese data-center company, since 2013. His role involves the investigation of essential network technologies for data centers, with a focus on IP closed network services between data centers and optical transmission technologies. Additionally, he has been an active member of the NOC team for ShowNet since 2017. He contributed to the demonstration of data-center network interconnection technologies from 2017 to 2019 and was involved in optical transmission technology from 2021 to 2024. He can be reached at: **kato@interop-tokyo.net**

The Root of the DNS

by Geoff Huston, APNIC

The *Domain Name System* (DNS) of the Internet is a remarkably simple system. You send queries into this system via a call to the name resolution library of your local host, and you get answers back. If you peek into the DNS system you'll see exactly the same simplicity: The DNS resolver that receives your query may not know the answer, so it, in turn, will send queries deeper into the system and collect the answers. This query/response process is the same, applied recursively. Simple.

However, the DNS is simple in the same way that Chess or Go are simple. They are all constrained environments governed by a small set of rigid rules, but they all generate surprising complexity in their operation.

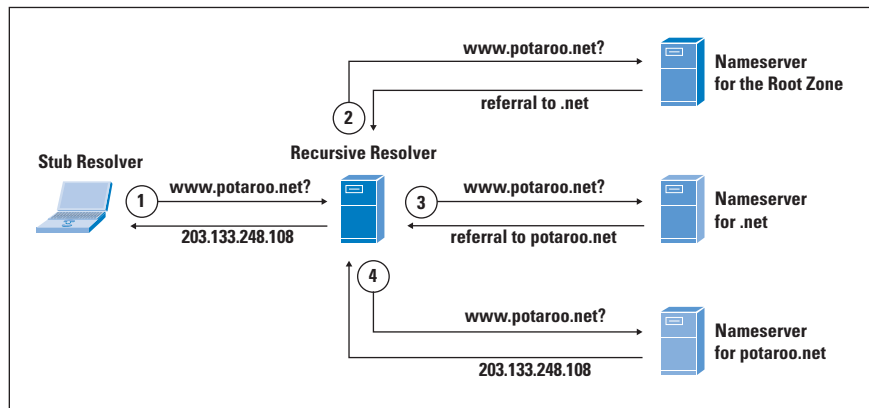
The Root Zone

The DNS is not a dictionary of any natural language, although these days when we use DNS names in our written and spoken communications we might be excused from getting the two concepts confused! The DNS is a hierarchical namespace. Individual domain names are constructed using an ordered sequence of labels. This ordered sequence of labels serves numerous functions, but perhaps most usefully it can be used as an implicit procedure to translate a domain name into an associated attribute value through the DNS name resolution protocol.

For example, I operate a web server that is accessed using the DNS name **www.potaroo.net**. If you direct your browser to load the contents of this DNS name, your system first needs to resolve this DNS name to an IP address, so that your browser knows where to send the IP packets to perform a transaction with my server. But how does the system know which nameserver is authoritative for the zone that includes the name **www.potaroo.net**?

This point is where the structure of the namespace is used to discover the nameserver. In this case, the DNS resolver will query a *root server* to resolve the name. As this name is not defined within the *Root Zone* (the zone that is served by the root servers), the response from any root server to such a query will be a *referral* response. In this example, this response is a redirection that lists the set of nameservers that are authoritative for the **.net** zone. Ask any of these **.net** nameservers for this same DNS name and again you will get back a redirection response, consisting of the list of nameservers that are authoritative for the **potaroo.net** zone. Ask any of these **potaroo.net** nameservers for the same name, **www.potaroo.net**, and you will receive the IP address you are looking for (Figure 1).

Figure 1: Name Resolution in the DNS.



Every DNS name is resolved in the same way. The name itself defines the order of name resolution processing, and it defines the path to be followed through the distributed database that leads to the answer you seek.

In this entire process, there is one starting point for every DNS resolution operation: the *Root Zone*.

Some criticize any exceptional consideration given to the root zone of the DNS; they think it is just another DNS zone, like any other. It is a set of authoritative servers that receive queries and answer them, like any other zone. There is no magic in the root zone, and all this attention on the root zone as *special* in some way is entirely unwarranted.

However, I think this view understates the criticality of the root zone in the DNS. The DNS is a massive, distributed database. Indeed, it is so massive that there is no single static map that identifies every authoritative source of information and the collection of data points about which it is authoritative. Instead, we use a process of dynamic discovery, where the resolution of a DNS name is first directed to locating the authoritative server that has the data relating to the name we want resolved, and then querying this server for the data. The beauty of this system is that these discovery queries and the ultimate query are precisely the same query in every case.

But everyone has to start somewhere. A DNS recursive resolver does not know all the DNS authoritative servers in advance, and it never will. But it does know one thing: It knows the IP address of at least one of the root servers in its provided configuration. From this starting point everything can be constructed in real time. The resolver can ask a root server for the names and IP addresses of all other root servers (the so-called *priming query*), and it can store that answer in a local cache. When the resolver is given a name to resolve, it can then start with a query to a root server to find the next point in the name delegation hierarchy and go on from there in a recursive manner.

If this description illustrates how the DNS actually works, then it is pretty obvious that the entire DNS system would have melted down years ago. What makes this approach viable is *local caching*. A DNS resolver stores the answers in a local cache and uses this locally held information to answer subsequent queries for the life of the cached entry. So perhaps a more refined statement of the role of the root servers is that every DNS resolution operation starts with a query to the cached state of the root zone. If the local cache cannot answer the query, then a root server must be queried.

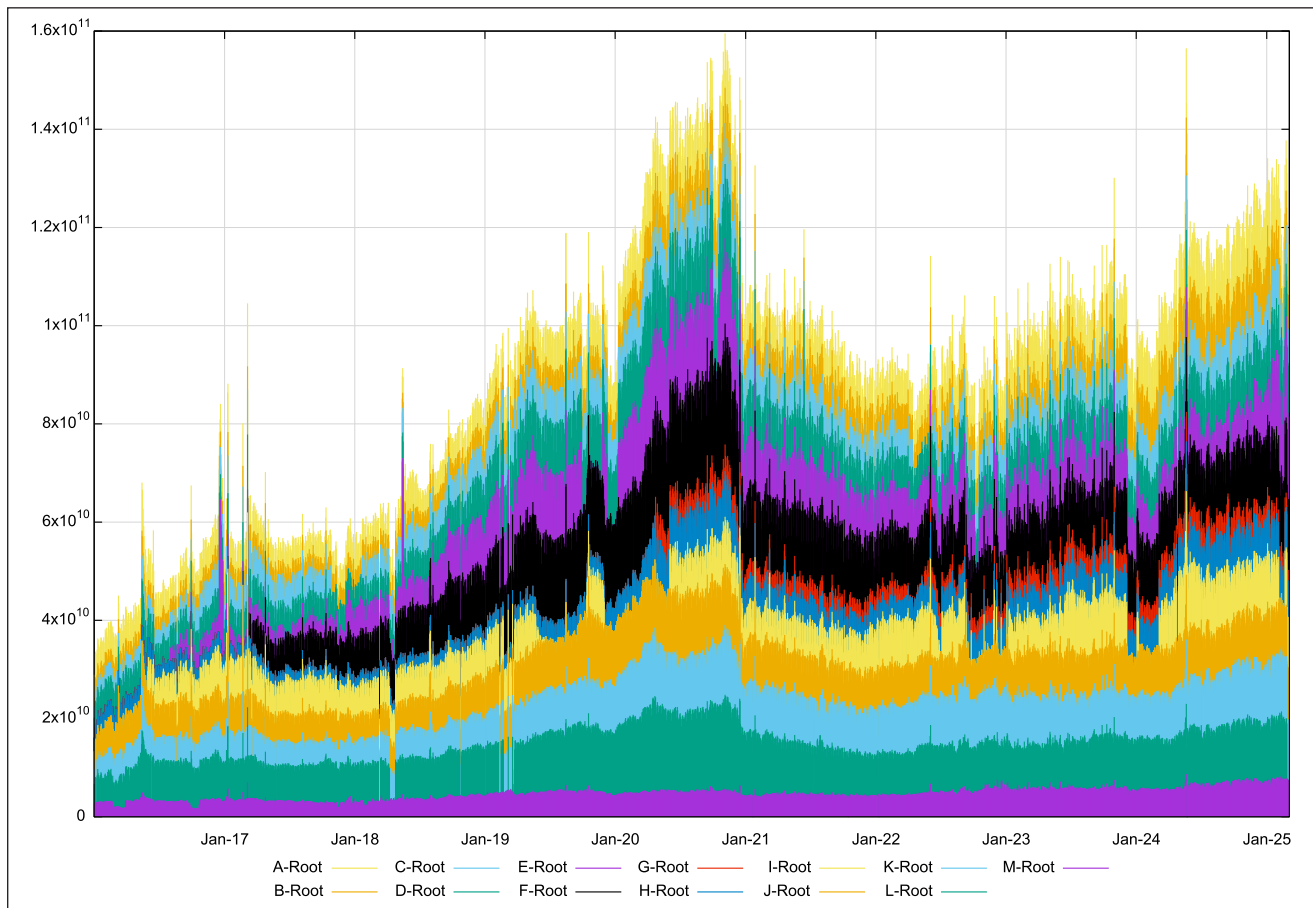
However, behind this statement lurks an uncomfortable observation: If all of the root servers are inaccessible, then the entire DNS ceases to function. This is perhaps a dramatic overstatement in some respects, as there would be no sudden collapse of the DNS and the Internet along with it. In the hypothetical situation where all the instances of the root servers were inaccessible, then DNS resolvers would continue to work using locally cached information. However, as these cached entries time out, they would be discarded from these local resolvers (as they could not be refreshed by re-querying the root servers). The light of the DNS would slowly fade to black bit by bit as these cached entries time out and are removed. The DNS root zone is the master lookup for every other zone. That's why it deserves particular attention. For that reason, the DNS root zone is uniquely different from every other zone.

Root zone servers are not used for every DNS lookup because of local caching. The theory is that the root servers will only see queries as a result of cache misses in resolvers. With a relatively small root zone and a relatively small set of DNS recursive resolvers, the root zone query load should be small. Even as the Internet expands its user base the query load at the root servers does not necessarily rise in direct proportion. It is the number of DNS resolvers that supposedly determines root server query load if we believe in this model of the function of the root in the DNS.

However, the model does not appear to hold up under operational experience. Figure 2 shows the total volume of queries per day recorded by the root servers since January 2016.

Over the period from 2016 to 2020, the volume of queries seen by the collection of root servers tripled. The query volume decreased in 2021 and stabilised over 2022. It is likely that changes to the behaviour of the Chrome browser may explain this abrupt change. Chrome used to probe the local DNS environment by making a sequence of queries to non-existent names (so-called *Chromeoids*) upon startup, and because the query names referred to undelegated top-level domains, these queries were a significant component of the queries seen at the root servers. Changing this behaviour in Chrome at the end of 2020 appears to have resulted in a dramatic change to the DNS query profile as seen by the root servers. However, over 2023 and 2024 the aggregate volume of queries seen by the root servers resumed its upward trend, rising by 40% from some 90 billion queries per day at the start of 2023 to more than 130 billion queries per day at the start of 2025.

Figure 2: Root Service Queries per Day – from ^[1].



What are we doing in response to this trend in the growth of queries to the root zone? How are we ensuring that the root zone service can continue to grow in capacity in response to this resumption in the growth of query rates?

Digression – The Economics of the DNS

In conventional markets, when a good is consumed, the consumer pays the producer a fee for the consumption of that good. As long as the fee covers the cost of production of the good, increasing consumption generates increasing revenue that can cover the costs associated with expanding the means of production of the good. Obviously, that's a very simplistic view of the operation of markets, but the key assumption is that greater consumption generates more revenue for producers, which, in turn, allows producers to produce greater volumes of the good. The essential assumption is that there is an underlying market-based discipline associated with the production and consumption of the good.

This assumption breaks down in the DNS, and in the root zone servers in particular. DNS queries are essentially unfunded. Like many Internet users, I have an *Internet Service Provider* (ISP), and I pay an access fee for its service.

Typically, an ISP operates a DNS recursive resolver for its clients, and my access fee contributes to my ISP's costs in running this resolver service. However, it's a fixed access fee, not a metered fee, so I contribute the same sum to the running of this shared resolver whether I submit one DNS query per day or one million!

As well as the costs in operating this resolver service, does the ISP incur any other cost in operating a DNS service to resolve my queries? No! All of the authoritative nameservers that are queried by my ISP's resolver are not funded by my ISP. More generally, all DNS queries in the public Internet are not directly funded by the querier!

Obviously, there are costs associated with operation of authoritative nameservers, and, for the most part, these costs are met by the "owners" of the zones that are served by these nameservers. There are various funding models for authoritative nameservers, ranging from metered costs per answered query, flat-rate costs, and even free services in some circumstances. But the essential aspect of this service is that authoritative nameservers do not derive revenue from the entities that query them. If there is a revenue stream, it comes from the DNS zone administrators who are paying for the nameservers to serve their zone.

I did note that this fact holds "for the most part," and there is one very notable exception here, namely the root zone. The twelve entities who provide the nameservers for the root zone do so as a collection of independent, largely autonomous volunteers who meet their own costs.

This situation is in many ways a curious relic of an earlier Internet that had a spirit of cooperative enterprise in many of its endeavours, but at the scale where each *Root Service Operator* is operating a service platform capable of responding to an average query load of some 10 billion queries per day, then it is no slight donation of effort and resources to a common-good outcome. Such a core of altruism in the centre of a market-driven frenzy of activity that operates today's digital world is unusual to see.

Given the criticality of the role that these operators collectively undertake, and the observation that directly or indirectly we are all beholden to the outcomes of these efforts to maintain a functional namespace for the Internet, then perhaps, odd as it may be, this situation is better than many of the alternatives.

In a market economy, a monopoly supplier of a critical resource is able to extract a monopoly rental from all others, while customers cannot seek relief through competitive offerings because of the very nature of the monopoly. Today's world looks to market regulators and the associated public regulatory frameworks to protect markets from such forms of abuse. But in the Root Service function we find a service that is both universal across the entire collection of individual public regimes and a collective monopoly.

A self-imposition by these operators of a freely offered service is perhaps not the only possible response to counter such risks of potential abuse of role. So far, however, the ethos of these twelve root service operators has proved to be an adequate and sufficient measure.

But perhaps it's now time to consider the outstanding question, namely "How are we ensuring that the root-zone service can continue to grow in capacity in response to this resumption in the growth of query rates?", and now factor in the apparent need to escalate the level of resources that are in effect donated to this service by this small collection of operators.

Root Zone Scaling

The original model of authoritative servers in the DNS was based on the concept of *unicast* routing. A server name had a single IP address, and this single server was located at a single point in the network. Augmenting server capacity entailed using a larger server and adding network capacity. However, such a model does not address the issues of a single point of vulnerability, nor does it provide an optimal service for distant clients.

More Servers

The DNS approach to this problem is to use multiple nameserver records. A DNS resolver was expected to retry its query with a different server if its original query did not elicit a response. That way, a collection of servers could provide a framework of mutual backup. To address the concept of optimal choice, DNS resolvers were expected to maintain a record of the query/response delay for each of the root servers and prefer to direct the majority of their queries to the fastest server.

Why not use multiple address records for a single common server name? The two approaches (multiple server names and multiple address records for a name) look similar. Once a resolver has assembled a collection of IP addresses that represent the nameservers for a domain, then it seems to me that a resolver could be justified for treating the list of IP addresses consistently, irrespective of whether the list was assembled from multiple IP addresses associated with a single name, or from multiple names. The use of multiple names allows for the use of multiple paths through the DNS to resolve these names of the nameservers that can remove a potential single point of failure, although I wonder as to the true benefit of using a set of nameserver names within a common single DNS zone as compared to using a single name with multiple IPv4 and IPv6 *Resource Records*, particularly when the bulk of DNS zones are provisioned with 2 or 4 nameservers, so there are typically 2 or 4 IPv4 and IPv6 addresses. I suspect that the use of multiple names is a policy compliance outcome rather than a true effort to provision nameservers with resilience through diversity.

If we want to increase the capacity of the root zone, then why not just add more nameserver names to the root zone?

What's so special about this zone's use of 13 named nameservers and a total of 26 IP addresses? For the root zone, the scaling issue with multiple nameservers is the question of completeness and the size of the nameserver response to the *priming query*. The question here is: If a resolver asks for the nameservers of the root zone, should the resolver necessarily be informed of all such servers in the response? The size of the response will increase with the number of servers, and the size of the response may exceed the default maximal DNS over a *User Datagram Protocol* (UDP) payload size of 512 bytes.

The choice of the number of server names for the root zone, 13, was based on the calculation that this was the largest list of a server list that could fit into a DNS response that was under 512 bytes in size. This choice assumed that only the IPv4 address records were being used in the response. With the addition of the IPv6 AAAA records, the response size has expanded. The size of the priming response for the root zone with 13 dual-stack authoritative servers is 823 bytes, or 1,097 bytes if the *Domain Name System Security Extensions* (DNSSEC) signature is included, and slightly larger if DNS cookies are added.

In today's DNS environment, if the query does not include an *Extension Mechanisms for DNS* (EDNS)(0)^[2] indication that they can accept a DNS response over UDP larger than 512 bytes, then the root servers will provide a partial response in any case, usually listing all 13 names, but truncating the list of addresses of these services in the *Additional Section* of the response to fit with a 512-byte payload.

Past experiments have been conducted with more than 13 nameservers at the apex of a DNS-like name system (such as the *Yeti*^[3] project, of some 5–8 years ago), and while it is technically feasible to do so, some vexing questions remain, such as how to select new root service operators, what is a safe ceiling of the number of such services, and how would it impact the stability and coherence of the name system.

Until we have much broader levels of adoption of query name minimisation than we appear to have today, root servers are privy to the myriad of domain names that users are querying. Such data is effectively a real-time view into the activity in the Internet through this meta-data query stream. If we opened up the root service to a broader set of operators, would a temptation to monetise this unique and highly valuable data stream prove overwhelming? In this space is it even possible to enforce constraints that would preclude any such activity?

So far, we appear to have avoided such difficult questions by leaving the number of root nameservers constant and scaling the root service in other ways.

If we can't, or don't want, to just keep on adding more root servers to the nameserver set in the root zone, then what are the other scaling options for serving the root zone?

More Service Platforms

The first set of responses to these scaling issues was in building root servers that have greater network capacity and greater processing throughput. But with just 13 servers to work with, this capacity was never going to scale at the pace of the Internet. We needed something more.

The next scaling step was the conversion from unicast to *anycast*^[32–37] services. There may be 26 unique IP addresses for root servers (13 in IPv4 and 13 in IPv6), but each of these service operators now uses anycast to replicate the root service in different locations. The current number of root server sites is described at root-servers.org (Table 1). Now the routing system is used to optimise the choice of the “closest” location for each root server.

Table 1: Anycast Site Counts for Root Servers, March 2025^[4].

Root	A	B	C	D	E	F	G	H	I	J	K	L	M	Total
Sites	59	6	13	220	328	359	6	12	85	148	131	123	23	1,513

The root server system has embraced anycast, some parts more enthusiastically than others. Currently a total of 1,513 sites have one or more instances of root servers. Some 24 months earlier, in January 2023, the root server site count was 1,396, so that’s an 8% increase in the number of sites in a little over two years.

The number of authoritative server instances is larger than the number of sites, as it is common these days to use multiple server engines within a site and use some form of query distribution at the front end to distribute the incoming query load across multiple back-end engines at each site. Today, the total of root server system instances is 1,907.

Even this form of expanding the distributed service may not be enough in the longer term. We are seeing the resumption of the growth profile last seen in 2016–2020. With a 25% compound annual query growth rate, in four years we may need double the root service capacity from the current levels, and in a further four years we’ll need to double it again. Exponential growth is a very harsh master.

Can this anycast model of replicated root servers expand indefinitely? Or should we look elsewhere for scaling solutions?

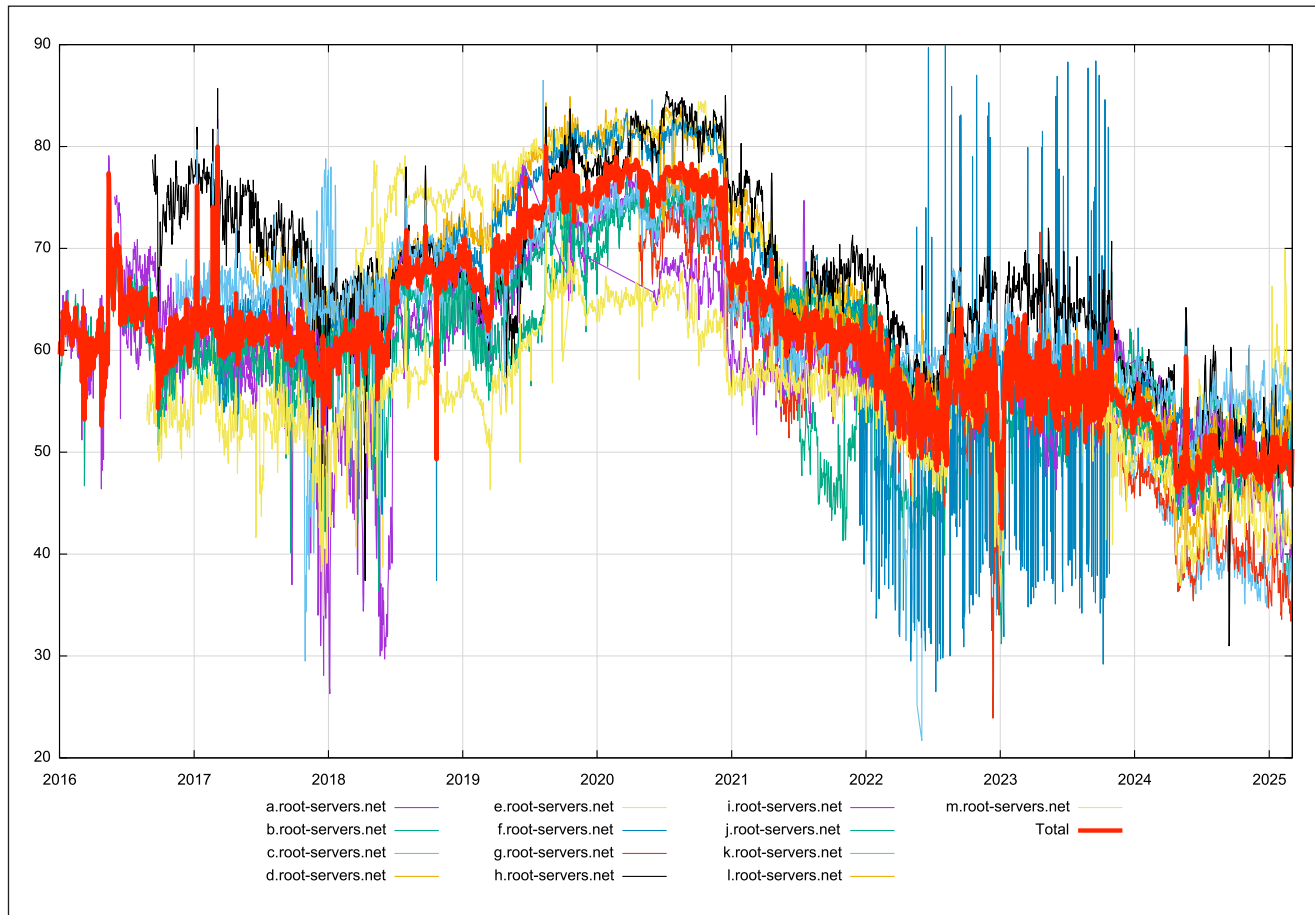
Query Deflection for Negative Responses

There have been many studies of the root service and the behaviour of the DNS over the past few decades. If the root servers were meant simply to respond to the cache misses of DNS resolvers, then whatever is happening at the root is not entirely consistent with such a model of behaviour. Indeed, it’s not clear what is going on at the root!

It has been reported that the majority of queries to the root servers result in NXDOMAIN (“non-existent domain”) error responses.

In looking at the published response code data, it appears that some 50% of root zone queries result in NXDOMAIN responses (Figure 3). The NXDOMAIN response rate was as high as 75% in 2020, and dropped presumably when the default behaviour of the Chrome browser in using Chromeoids changed. In theory these queries are all cache misses at the recursive resolver level, so the problem is that the DNS is not all that effective in handling cases where the name itself does not exist.

Figure 3: Proportion of Root Zone NXDOMAIN Responses per Day ^[1].



If we want to reduce the query pressure on the root servers, one possible approach is to alter the way DNS resolvers handle queries for non-existent names, and in particular names where the top-level label in the queried name is not delegated in the root zone. How else can we deflect these queries away from the root server system?

One such approach is described in RFC 8198^[5], “Aggressive NSEC Caching.” When a top-level label does not exist in a DNSSEC-signed zone and the query has the EDNS(0) DNSSEC “OK” flag enabled, the NXDOMAIN response from a root server includes a signed NSEC record that gives the two labels that exist in the root zone that “surrounds” the non-existent label.

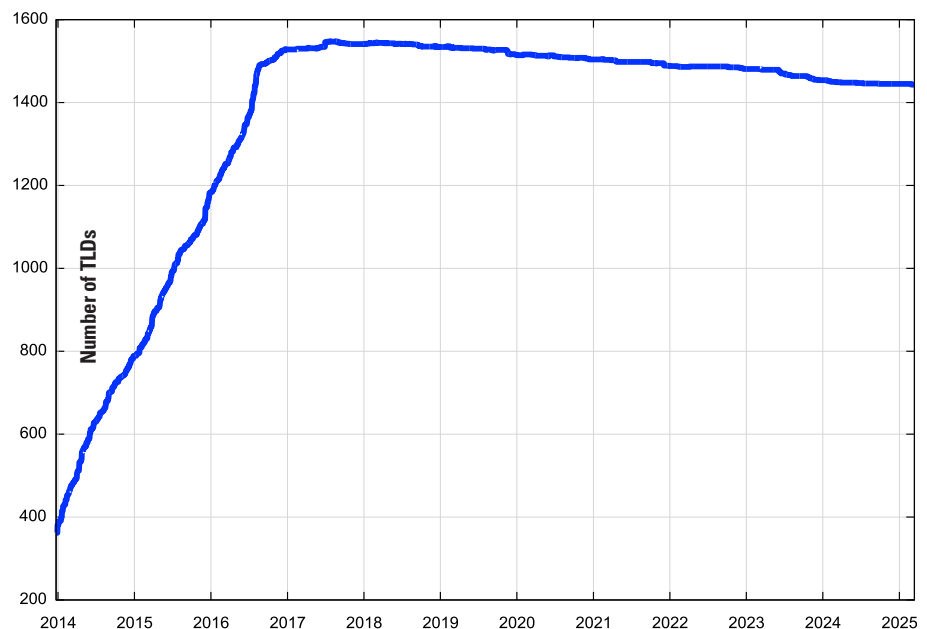
NSEC records say more than “this label is not in this zone.” It says that no label that is lexicographically between these two labels exists. If the recursive resolver caches this NSEC record, then it can use this same cached record to respond to all subsequent queries for names in this label range, in the same way that it conventionally uses “positive” cached records.

If a recursive resolver cached both the 1,443 top-level delegated labels and the 1,444 NSEC records in the root zone, then the resolver would not need to pass any queries to a root server for the lifetime of the cached entries. If all recursive resolvers performed this form of NSEC caching of the root zone, then the query volumes seen at the root from recursive resolvers would fall significantly for non-existent labels.

How Many TLDs Are in the Root Zone?

There were 1,443 *Top-Level Domains* (TLDs) in the root zone of the DNS in March 2025. It has not always been this size. The root zone started with a small set of generic labels, and in the late 1980’s expanded to include the set of two-letter country codes. There were some tentative steps to augment the number of generic top-level domain names, and then in the 2010s *The Internet Corporation for Assigned Names and Numbers* (ICANN) embarked on a larger program of generic TLD expansion. Figure 4 shows the daily count of TLDs in the root zone since 2014.

Figure 4: Daily Count of Root Zone TLDs.



What was surprising to me was that TLDs are not necessarily permanent. The largest TLD count occurred in August 2017, with 1,547 TLDs, and since then the number of TLDs has been declining.

Aggressive use of NSEC caching in recursive resolvers appears to play a contributory role in helping us scale the root zone. *Bind* supports this function as of release 9.12, *Unbound* supports it as of release 1.7.0, and *Knot* resolver supports it as of version 2.0.0. But the queries at the root zone keep growing despite the declining proportion of queries, resulting in an NXDOMAIN response. While this measure may have dampened the relative growth of queries for non-existent names seen at the root servers, to some extent it has not significantly affected the overall problem of the growth of queries directed to the root servers; other factors appear to be causing it.

I'd characterise the situation as aggressive NSEC caching representing a tactical response to root zone scaling concerns, as distinct from a strategic response. The technique is still dependent on the root server infrastructure, and it uses a query-based method of promulgating the contents of the root zone. Nothing really changes in the root service model. What NSEC caching does is allow the resolver to make full use of the information in the NSEC response.

Root Zone Mirroring

Another scaling option is to jump completely out of the query/response model where recursive resolvers incrementally learn the contents of the root zone query-by-query and simply load the entire root zone into their local cache and refresh this local copy with a period of several hours or even a day or so. The idea here is that if a recursive resolver is loaded with a copy of the root zone, then it can operate autonomously with respect to the root servers for the period of validity of the local copy of the root zone contents. It will send no further queries to the root servers.

The procedures to follow to load a local root zone are well documented in RFC 8806^[6], and I should also note here the existence of the *LocalRoot*^[7] service that apparently offers DNS NOTIFY messages when the root zone changes. The root zone is not a big data set. A signed, uncompressed plaintext copy of the root zone as of March 14, 2025, is 2.2 MB in size.

However, this approach has its potential drawbacks. How do you know that the zone you might have received via some form of zone transfer or other is the current genuine root zone? Yes, the zone is signed, but not every element in the zone is signed (NS records for delegated zones are unsigned). The client is left with the task of performing a validation of every digital signature in the zone, and at present there are some 1,444 *Resource Record Digital Signature* (RRSIG) records in the root zone. Even then the client cannot confirm that its local copy of the root zone is complete and authentic because of the unsigned NS delegation records in the root zone.

The IETF published RFC 8976^[8], the specification of a message-digest record for DNS zones, in February 2021. This RFC defines the *Message Digest for DNS Zones* (ZONEMD) record.

What's a Message Digest?

A *Message Digest* is a form of a condensed digital signature of a digital artefact. If the digital artefact has changed in any way, the digest will necessarily change in value as well. If a receiver of this artefact is given the data object and its digest value, then the receiver can be assured, to some extent, that the contents of the file have been unaltered since the digest was generated.

These digital signatures are typically generated using a *Cryptographic Hash Function*. These functions have several useful properties. They are normally a fixed-length output function, so that the resulting value is a fixed size, irrespective of the size of the data for which the hash has been generated.

They constitute a *unidirectional* function, in that knowledge of the hash function value will not provide any assistance in trying to recreate the original data. They are *deterministic*, in that the same hash function applied to the same data will always produce the same hash value. Any form of change to the data should generate a different hash value. Hash functions do not necessarily produce a unique value for each possible data collection, but it should be exhaustively challenging (unfeasible) to synthesise or discover a data set that produces a given hash value (*preimage resistance*), and equally challenging to find or generate two different data sets that have the same hash function value (*collision resistance*).

In other words, an adversary, malicious or otherwise, cannot replace or modify the data set without changing its digest value. Thus, if two data sets have the same digest, one can be relatively confident that they are identical. Second pre-image resistance prevents an attacker from crafting a data set with the same hash as a document the attacker cannot control. Collision resistance prevents an attacker from creating two distinct documents with the same hash.

The root zone includes a ZONEMD record, signed with the *Zone Signing Key* of the root zone. When a client receives the root zone it should look for this record, validate the *Resource Record Digital Signature* (RRSIG) of the ZONEMD record in the same way that it DNSSEC-validates any other RRSIG entry in the root zone, and then compare the value of this record with a locally calculated message digest value of the local copy of the root zone. If the digest values match, then the client has a high level of assurance that this copy of the root zone is authentic and has not been altered in any way.

The dates in the DNSSEC signatures can indicate some level of currency of the data, but further assurance at a finer level of granularity than the built-in key validity dates that the local copy of the root zone data is indeed the current value of the root zone is a little more challenging in this context. DNSSEC does not provide any explicit concept of revocation of prior versions of data, so all “snapshots” of the root zone within the DNSSEC key validity times are equally valid for a client.

The root zone uses a two-week signature validity period (Figure 5).

Figure 5: Root Zone Start of Authority (SOA) Signature.

```
. 86400 IN SOA a.root-servers.net. nstld.verisign-grs.com.
2025031303 1800 900 604800 86400

. 86400 IN RRSIG SOA 8 0 86400 20250326200000 20250313190000
26470 . nYhmvV[...]Ng==
```

This approach of a whole-of-zone signature has some real utility in terms of the distribution of the root zone to DNS resolvers and thereby reduces the dependency on the continuous availability and responsiveness of the root zone servers. The use of the ZONEMD record allows any client to use a local copy of the root zone irrespective of the way in which the zone file was obtained. Within the limits of the authenticated currency of the zone file, as already noted, any party can redistribute a copy of the root zone, and clients of such a redistributed zone can answer queries using this data with some level of confidence that the responses so generated are authentic. It would be useful to augment the existing in-band root zone retrieval using *Authoritative Transfer* (AXFR) with a simple memorable web-retrieval object, such as https://1.2.3.4/root_zone.txt, for example, to allow the zone distribution function to be undertaken by *Content Distribution Networks* (CDNs) as well as by DNS servers.

Resolvers that elect to use a locally managed copy of the root zone can use the ZONEMD record to verify the authenticity of a received root zone. Resolver implementations that perform this verification using ZONEMD include *Unbound* (from v1.13.23) and *PowerDNS Recursor* (from v4.7.04) and *Bind* (v9.17.13).

Notification mechanisms that could prompt a resolver to work from a new copy of the root zone are not addressed in this ZONEMD framework. To me that's the last piece of the framework that could promote every recursive resolver into a peer root server. We've tried numerous approaches to scalable distribution mechanisms over the years. There is the structured *push* mechanism, where clients sign up to a distributor and the distributor pushes updated copies of the data to them. Routing protocols use this mechanism. There also is the *pull* approach, where the client probes its feed point to see if the data has changed and pulls a new copy if it has changed. This mechanism has some scaling issues in that aggressive probing by clients may overwhelm the distributor. We've also seen hybrid approaches where a change indication signal is pushed to the client, and it is up to the client to determine when to pull the new data.

This model of local root zone distribution has the potential to change the nature of the DNS root service, unlike NSEC caching. If there is one thing that we've learned to do astonishingly well in recent times it is distribution of content.

Indeed, we've concentrated on this activity to such an extent that it appears that the entire Internet is nothing more than a small set of CDNs. If the root zone is signed in its entirety with zone signatures that allow a recursive resolver to confirm its validity and currency and is submitted into these distribution systems as just another digital object, then the CDN infrastructure is perfectly capable of feeding this zone to the entire collection of recursive resolvers with ease. Perhaps if we changed the management regime of the root zone to generate a new zone file every 24 hours according to a strict schedule, we could eliminate the entire notification superstructure. Each iteration of the root zone contents would be published 2 hours in advance and it would be valid for a period of precisely 48 hours, for example. At that point the root zone could be served by the existing millions of recursive resolvers rather than the twelve operators and some 2,000 server instances we use today. That's a thousand-fold increase in the capacity of the root system, and at the same time it eliminates the general reliance on a narrow neck of incremental queries being directed to the 12 root server operators that underpin today's DNS.

Futures

We operate the root service in its current framework because it represents a set of compromises that have been functionally adequate so far. That is to say the predominate query-based approach to root zone distribution hasn't visibly collapsed in a screaming heap of broken DNS yet! And it will probably continue to operate in a robust manner for many years to come.

But we don't have to continue relying on this query-based approach just because it hasn't broken so far. Our need to further scale this function is ongoing, and it makes a lot of sense to take a broader view of available options and the just-in-time delivery process used by the DNS incremental query name-resolution algorithm.

We have some choices as to how the root service can evolve and scale.

With Aggressive NSEC Caching we can have recursive resolvers make better use of signed NSEC records and we appear to have staved off some of the more pressing immediate issues about further scaling of the root system. But that's probably not enough.

We can either wait for the DNS system to collapse and then try to salvage the DNS from the broken mess, or perhaps we could explore some alternatives now. For example, we could look at how we can break out of a query-based incremental root content promulgation model and view the root zone as just another content "blob" in the larger ecosystem of content distribution. If we can cost-efficiently load every recursive resolver with a current copy of the root zone, and these days that's not even a remotely challenging target, then perhaps we can put aside the issues of how to scale the root server system to serve ever greater volumes of queries to ever more demanding clients, and perhaps also provide an alternate answer to the continual questions about the politics and finances relating to root servers and their operation.

The reason why content distribution networks have revolutionised the Internet in recent years is that pre-provisioning at the edge makes for a faster, cheaper, and more scalable network in the current context of abundant computing and storage capabilities. If we are prepared to allow this same thinking to intrude into the way we provision the DNS, I suspect we could realise similar benefits for the DNS as well.

Disclaimer

The views shared herein do not necessarily represent the views or positions of the Asia Pacific Network Information Centre.

References and Further Reading

- [0] ICANN Root Server System Advisory Committee (RSSAC), “RSSAC Advisory on Measurements of the Root Server System,” RSSAC002, June 1, 2016.
- [1] A collection of collected RSSAC002 data:
<https://github.com/rssac-caucus/RSSAC002-data>
- [2] Joao Damas, Michael Graff, and Paul Vixie, “Extension Mechanisms for DNS (EDNS(0)),” RFC 6891, April 2013.
- [3] The Yeti Project: <https://yeti-dns.org/>
- [4] <https://root-servers.org/>
- [5] Kazunori Fujiwara, Akira Kato, and Warren Kumari, “Aggressive Use of DNSSEC-Validated Cache,” RFC 8198, July 2017.
- [6] Warren Kumari and Paul Hoffman, “Running a Root Server Local to a Resolver,” RFC 8806, June 2020.
- [7] *LocalRoot*—Serve Yourself the Root:
<https://localroot.isi.edu/>
- [8] Duane Wessels, Piet Barber, Matt Weinberg, Warren Kumari, and Wes Hardaker, “Message Digest for DNS Zones,” RFC 8976, February 2021.
- [9] Jon Postel, “Internet Name Server,” IEN 61, October 1978.
- [10] Paul Mockapetris, “Domain names: Concepts and facilities,” RFC 882, November 1983.
- [11] Paul Mockapetris, “Domain names: Implementation specification,” RFC 883, November 1983.
- [12] Zi Hu, Liang Zhu, John Heidemann, Allison Mankin, Duane Wessels, and Paul Hoffman, “Specification for DNS over Transport Layer Security (TLS),” RFC 7858, May 2016.
- [13] Christian Huitema, Sara Dickinson, and Allison Mankin, “DNS over Dedicated QUIC Connections,” RFC 9250, May 2022.
- [14] Paul Hoffman and Patrick McManus, “DNS Queries over HTTPS (DoH),” RFC 8484, October 2018.
- [15] Eric Kinnear, Patrick McManus, Tommy Pauly, Tanya Verma, and Christopher A. Wood, “Oblivious DNS over HTTPS,” RFC 9230, June 2022.

- [16] Stephane Bortzmeyer, “DNS Query Name Minimisation to Improve Privacy,” RFC 7816, March 2016.
- [17] Roy Arends, Rob Austein, Matt Larson, Dan Massey, and Scott Rose, “DNS Security Introduction and Requirements,” RFC 4033, March 2005.
- [18] Ben Schwartz, Mike Bishop, and Erik Nygren, “Service Binding and Parameter Specification via the DNS (SVCB and HTTPS Resource Records),” RFC 9460, November 2023.
- [19] Ben Schwartz, “Service Binding Mapping for DNS Servers,” RFC 9461, November 2023.
- [20] Carlo Contavalli, Wilmer van der Gaast, David C. Lawrence, and Warren Kumari, “Client Subnet in DNS Queries,” RFC 7871, May 2016.
- [21] Miek Gieben, “DNSSEC: The Protocol, Deployment, and a Bit of Development,” *The Internet Protocol Journal*, Volume 7, No. 2, June 2004.
- [22] Richard Barnes, “Let the Names Speak for Themselves: Improving Domain Name Authentication with DNSSEC and DANE,” *The Internet Protocol Journal*, Volume 15, No.1, March 2012.
- [23] M. Stuart Lynn, “A Unique Root,” *The Internet Protocol Journal*, Volume 4, No. 3, September 2001.
- [24] Geoff Huston, “A Question of DNS Protocols,” *The Internet Protocol Journal*, Volume 17, No. 1, September 2014.
- [25] Geoff Huston, “Scaling the Root,” *The Internet Protocol Journal*, Volume 18, No. 1, March 2015.
- [26] Geoff Huston, “What’s in a DNS Name?” *The Internet Protocol Journal*, Volume 19, No. 1, March 2016.
- [27] Geoff Huston, “The Root of the DNS,” *The Internet Protocol Journal*, Volume 20, No. 2, June 2017.
- [28] Geoff Huston, “DNS Privacy and the IETF,” *The Internet Protocol Journal*, Volume 22, No. 2, July 2019.
- [29] Geoff Huston, “DNS Trends,” *The Internet Protocol Journal*, Volume 24, No. 1, March 2021.
- [30] Geoff Huston, “DNS Evolution,” *The Internet Protocol Journal*, Volume 27, No. 2, July 2024.
- [31] Burton Kaliski Jr., “Minimized DNS Resolution: Into the Penumbra,” *The Internet Protocol Journal*, Volume 25, No. 3, December 2022.
- [32] Craig Partridge, Trevor Mendez, and Walter Milliken, “Host Anycasting Service,” RFC 1546, November 1993.
- [33] David B. Johnson and Steve Deering, “Reserved IPv6 Subnet Anycast Addresses,” RFC 2526, March 1999.
- [34] Dino Farinacci and Yiqun Cai, “Anycast-RP Using Protocol Independent Multicast (PIM),” RFC 4610, August 2006.

- [35] Joe Abley and Kurt Eric Lindqvist, “Operation of Anycast Services,” RFC 4786, December 2006.
- [36] Danny McPherson, Dave Oran, Dave Thaler, and Eric Osterweil, “Architectural Considerations of IP Anycast,” RFC 7094, January 2014.
- [37] Sebastian Kiesel and Reinaldo Penno, “Port Control Protocol (PCP) Anycast Addresses,” RFC 7723, January 2016.

GEOFF HUSTON AM, B.Sc., M.Sc., is the Chief Scientist at APNIC, the Regional Internet Registry serving the Asia Pacific region. He has been closely involved with the development of the Internet for many years, particularly within Australia, where he was responsible for building the Internet within the Australian academic and research sector in the early 1990s. He is author of numerous Internet-related books, and was a member of the Internet Architecture Board from 1999 until 2005. He served on the Board of Trustees of the Internet Society from 1992 until 2001. At various times Geoff has worked as an Internet researcher, an ISP systems architect, and a network operator. E-mail: gih@apnic.net

Our Privacy Policy

The *General Data Protection Regulation* (GDPR) is a regulation for data protection and privacy for all individual citizens of the *European Union* (EU) and the *European Economic Area* (EEA). Its implementation in May 2018 led many organizations worldwide to post or update privacy statements regarding how they handle information collected in the course of business. Such statements tend to be long and include carefully crafted legal language. We realize that we may need to provide similar language on our website and in the printed edition, but until such a statement has been developed here is an explanation of how we use any information you have supplied relating to your subscription:

- The mailing list for *The Internet Protocol Journal* (IPJ) is entirely “opt in.” We never have and never will use mailing lists from other organizations for any purpose.
- You may unsubscribe at any time using our online subscription system or by contacting us via e-mail. We will honor any request to remove your name and contact information from our database.
- We will use your contact information only to communicate with you about your subscription; for example, to inform you that a new issue is available, that your subscription needs to be renewed, or that your printed copy has been returned to us as undeliverable by the postal authorities.
- We will never use your contact information for any other purpose or provide the subscription list to any third party other than for the purpose of distributing IPJ by post or by electronic means.
- If you make a donation in support of the journal, your name will be listed on our website and in print unless you tell us otherwise.

Letters to the Editor

As a long-time subscriber to the *Internet Protocol Journal*, I have always found the articles to be timely, extensively researched, and presented with clarity. Because of these qualities I've often used IPJ articles as reading assignments to my students in several of the courses that I teach. They give students a depth that is not too technical but yet contains enough of the necessary technical details that help them better understand the technology while also grasping the impact of the technology in the context of how it is being used today while looking ahead to the future.

I especially enjoy reading and sometimes sharing articles by Geoff Huston. Geoff has a real knack for presenting technical details within the larger scope of “how we got here and where we are going” that always makes his articles a pleasure to read.

However, Geoff's article on “The IPv6 Transition” in Volume 28, No. 1, May 2025 was especially beneficial. Students often ask me why, after all these years, IPv6 still plays “second fiddle” to IPv4 and is not more widely adopted. Geoff's article does a superb job of explaining the numerous reasons behind the slow transition to IPv6. And his analysis of how today's Internet is moving away from a strict address-based architecture offers an excellent assessment as to what lies ahead in the future for IPv6.

And the timeliness of Geoff's article could not have been better: after arriving in my email inbox that afternoon I had enough time to add some of his observations to my class discussion the very next morning!

Thanks, Geoff, and keep up the good work!

Regards,

—*Dr. Mark Ciampa*
Professor, Western Kentucky University
Bowling Green, KY
mark.ciampa@wku.edu

The Author responds:

Thank you Mark for your kind words. The Internet has not followed a path driven as much by market pressures as it is by technical evolution, and the outcomes are often surprising. This topic was the main theme of my article. One thing is sure, however, that the pressures to innovate will continue, and tomorrow will be as surprising as today!

Kind regards,

—*Geoff Huston*
gih@apnic.net

Hi Ole,

As always, I enjoy reading Geoff's articles in IPJ and I appreciate Geoff's continuing to write and your willingness to publish his material. What I so appreciate, Ole and Geoff, about IPJ is the *context* which you provide that gives me the larger frame/larger picture into which to place much of what I do. And even when I don't specifically need the context—I don't do anything with cellular networks for example—I find the articles fun reading regardless.

My “two bits” on Geoff's “IPv6 Transition” article: I logged into my first network-attached computer in 1981, as a student. I configured my first IP router in 1991, and I have spent my career to date supporting IT infrastructure (compute/network/storage) for academic or other non-profit research institutes, mostly in the life sciences. My perspective over the years, as I have sat in seminars at Interop or read news or attended internal meetings about IPv4 address exhaustion and the need for IPv6:

- Adding IPv6 support to our network would take money, staff hours, and training time, not only from network engineers but also from desktop, server, and storage system engineers; smells like a lot of effort to me.
- I find it difficult to prioritize distant risks over immediate priorities.
- Our user base has yet to ask for access to an IPv6-only resource (such a request would affect how we prioritize, but I have not seen even one).
- I suspect that any institution that wants to make a resource broadly available will invest significant effort into making it available via IPv4 because there are so many IPv4 users out there.

As a result, I have yet to configure any device to support IPv6. I am not opposed to IPv6 ... but since I still get up at 2am when my [pager](#) phone buzzes, in response to our network malfunctioning, I have—reasonably I propose—allocated my time to other priorities.

—Stuart Kendrick

Allen Institute, Seattle, WA USA
stuartk@alleninstitute.org

The Author responds:

Hi Stuart. We share a similar vintage, as I first logged into a computer as a student in 1976 (A Sperry Rand Univac mainframe), although building a network connection to the Internet for all Australian Universities would take a further 13 years, when the project that I was leading, AARNet, had managed to complete its initial mission.

In the early 1990's when the IETF was debating the approach to be used for the "next generation" IP protocol, there were many general approaches. One approach, originally called "SIP" was intended to change the IPv4 design as little as possible. It lengthened the address fields to 128 bits, but not much else changed. Other approaches described a more radical set of design changes. SIP won the day, and IPv6 is, to all intents and purposes, just IPv4 with bigger address fields.

In other words, IPv6 was not intrinsically "better" than IPv4 for any particular use case. It wasn't intrinsically faster, nor more secure, nor more agile. It just had bigger address fields. The result was that deploying IPv6 did not provide a network operator with a compelling competitive product. If a network operator already had secured ample pools of IPv4 addresses, then it was not necessarily impacted by IPv4 address scarcity, and the case for incurring the cost of deploying IPv6 in a dual-stack network scenario was extremely challenging to make. A direct result of this situation is the protracted transition to IPv6 for many parts of the Internet, an issue I explored in this article.

Frankly, it does not really make much sense to comment on the IPv6 design as "right" or "wrong." It represented the common wisdom of the IETF at that time. However, we did fail to predict just how long the dual-stack transition was going to take, or even if this transition would ever come to an end. Some 30 years after we started down this path I suspect we are no closer to answering these two very fundamental questions.

Regards,

—*Geoff Huston*
gih@apnic.net

Check your Subscription Details!

Make sure that both your postal and e-mail addresses are up-to-date since these are the only methods by which we can contact you. If you see the words "Invalid E-mail" on your printed copy, this means that we have been unable to contact you through the e-mail address on file. If this is the case, please contact us at **ipj@protocoljournal.org** with your new information. The subscription portal is located here:
<https://www.ipjsubscription.org/>

In Memoriam



Dave Täht

Dave Täht, formerly known as Michael David Täht (August 11, 1965 – April 1, 2025) was our friend, colleague, and mentor at LibreQoS. To the rest of the world, Dave was an American network engineer, musician, lecturer, asteroid exploration advocate, Internet activist, and much more.^[1]

The fruits of Dave’s work are everywhere. Most people will never notice—a testament to his engineering. Dave’s work on creating algorithms like *Flow Queueing with Controlled Delay*^[2] and *Common Applications Kept Enhanced* (CAKE)^[3] was instrumental, and now it’s part of Linux, OpenWrt, and Starlink; mainstream networking equipment vendors like MikroTik use it as well.

Dave and Jim Gettys spearheaded the networking industry’s effort to eliminate bufferbloat, latency, and jitter on today’s interactive Internet, where bandwidth matters less.^[4] Around 2010, Dave was semi-retired in Nicaragua—and Jim in the USA. They independently came to the realization that Voice over IP and videoconferencing were suffering from the same issues: *lag* and *jitter*, caused by the proverbial time difference between Dave speaking on one continent and Jim hearing his voice on another. Jim coined the term “bufferbloat” to describe the culprit: the extensive and ever-increasing size of buffering on network devices. They started **Bufferbloat.net** and began solving the problem.

Besides his work on solving bufferbloat, Dave also spent years in Nicaragua trying to find ways to bring the Internet (and power, lighting, food, medicine, and books) as an outgrowth of Nicholas Negroponte’s *One Laptop Per Child Project*^[5].

Dave was also known for his little ditties, songs he liked to play while presenting at conferences or podcasts. He wrote “One First Landing,” for example, to cheer up people at SpaceX when they were not doing well with landing their rocket, and he was forced to rewrite it when they started to land their rockets successfully.^[6] And he was happy to do it!

For the last couple of years Dave lived on a boat in Half Moon Bay, a small city on the California coast, south of San Francisco; Dave always liked to sail.

Ad astra per aspera, Dave, you are an astronaut now!

—Robert, Herbert, and Frank - LibreQoS

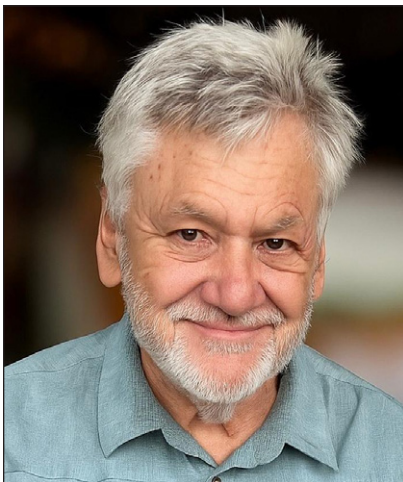
LibreQoS is an Open Source project founded by Robert Chacón (and quickly joined by Dave Täht). LibreQoS provides a drop-in middlebox for *Internet Service Providers* (ISPs), applying CAKE to all of the ISP’s users, as well as an array of network monitoring tools.

References

- [1] Wikipedia article on Dave Täht, May 2025:
https://en.wikipedia.org/wiki/Dave_T%C3%A4ht
- [2] Toke Hoeiland-Joergensen, Paul McKenney, Dave Täht, Jim Gettys, and Eric Dumazet, “The Flow Queue CoDel Packet Scheduler and Active Queue Management Algorithm,” RFC 8290, January 2018.
- [3] Toke Høiland-Jørgensen, Dave Täht, and Jonathan Morton, “Piece of CAKE: A Compre-hensive Queue Management Solution for Home Gateways.” <https://arxiv.org/abs/1804.07617>
- [4] Jim Gettys, “The Blind Men and the Elephant,” Jim Gettys’ ramblings on random topics, and occasional rants, February 11, 2018. <https://gettys.wordpress.com/2018/02/11/the-blind-men-and-the-elephant/>
- [5] Wikipedia article on One Laptop per Child:
https://en.wikipedia.org/wiki/One_Laptop_per_Child
- [6] YouTube video, “One First Landing (Thank you SpaceX for a Wonderful Year!).”
<https://www.youtube.com/watch?v=wjurORG-v-I>

“In my own mind, I like to think of him as the person who added the most effective capacity to the Internet.”

—Karl Auerbach



Frederick Juergens Baker

The Internet has lost a generous long-time contributor. Frederick Juergens Baker (February 28, 1952 – June 18, 2025)^[1] was one of the original members of the *Internet Systems Corporation* (ISC) Board of Directors, appointed at ISC’s incorporation in 1994.

Fred had a long career in the communications industry, working for Control Data Corporation, Vitalink Communications, Advanced Computer Communications, and for 22 years, at Cisco Systems.

After retiring from Cisco, Fred worked as a contractor, notably for the Internet Society and ISC. In addition to serving on the ISC BOD, in 2017 he joined the *Root Server System Advisory Committee* (RSSAC) of the *Internet Corporation for Assigned Names and Numbers* (ICANN), representing ISC. He served as co-chair of RSSAC from October 2018 to December 2019, and as chair from January 2020 through December 2022.

Fred volunteered a lot of his time to working with the *Internet Engineering Task Force* (IETF), the body that develops standards for the Internet. He chaired numerous IETF working groups, including several that specified the *Management Information Bases* (MIB) used to manage network bridges and popular telecommunications links, and the IPv6 Operations working group.

He served as IETF chair from 1996 to 2001, and he served on the Internet Architecture Board from 1996 through 2002. Fred co-authored or edited at least 60 *Request for Comments* (RFC) documents ^[2,3] on Internet protocols, and contributed to others. The subjects covered include network management, *Open Shortest Path First* (OSPF) and *Routing Information Protocol* (RIPv2) routing, *Quality of Service* (using both the *Integrated Services* and *Differentiated Services* models), *Lawful Interception*, precedence-based services on the Internet, and others.

In addition, Fred served as a member of the Board of Trustees of the Internet Society from 2002 through 2008, and as its chair from 2002 through 2006. He was a member of the *Technical Advisory Council* of the US *Federal Communications Commission* in 2004. He worked as a liaison to other standards organizations such as the ITU-T. In 2009–2010, he served as chair of the *RFC Series Oversight Committee*.

He represented IETF on the *National Institute of Standards and Technology Smart Grid Interoperability Panel* and the *Architecture Committee* from 2008 through 2013, and was Cisco's representative to the *Broadband Internet Technical Advisory Group* (BITAG). He also holds several patents.

Fred was committed to the collaborative, consensus-driven process of creating open standards for the Internet, and he demonstrated his commitment throughout his long career with years of active volunteering. Besides his leadership roles, he also welcomed and mentored new participants in the IETF.

Fred was a wonderful guy, an Internet luminary, and a great friend to ISC over the course of decades of board membership, as well as representing ISC at RSSAC as chair and in many other roles in the IETF, ICANN, and ISOC worlds. We all will dearly miss him.

We extend our deepest condolences to Fred's family.

—Jeff Osborn, ISC
jsosborn@isc.org

References

- [1] Wikipedia article on Fred Baker:
[https://en.wikipedia.org/wiki/Fred_Baker_\(engineer\)](https://en.wikipedia.org/wiki/Fred_Baker_(engineer))
- [2] RFCs authored or co-authored by Fred Baker:
<https://datatracker.ietf.org/doc/search?name=&sort=&rfcs=on&activedrafts=on&by=author&author=fred+baker>
- [3] Datatracker Profile for Fred Baker:
<https://datatracker.ietf.org/person/fredbaker.ietf@gmail.com>
- [4] In Memory of Fred Baker, *Ever Loved*:
<https://everloved.com/life-of/frederick-baker/>

IETF-developed MLS set to be used on 100s of Millions of Mobile Devices

Less than two years after *Messaging Layer Security* (MLS) was published as an RFC^[0], it is poised to be deployed on Android phones and Apple iPhones and other devices, thanks to newly updated RCS specifications, enabling interoperable encryption between different platform providers for the first time.

The GSM Association^[1] announced that the latest *Rich Communications Services* (RCS) standard includes end-to-end encryption based on the MLS protocol. RCS enhances traditional SMS messaging by offering a suite of service capabilities, including group chat, file transfers, typing notifications, and more. Key stakeholders for RCS implementation include device manufacturers, telecommunications operators, and business service providers.

MLS, developed by the IETF *Messaging Layer Security Working Group*^[2], provides unsurpassed security and privacy for users of group communications applications. Using MLS, participants always know which other members of a group will receive the messages they send, and the validity of new participants joining a group is verified by all the other participants. During its development^[3] in the IETF, MLS underwent formal security analysis and industry review. It currently supports multiple cipher suites, and makes it straightforward to add quantum attack resistant cipher suites in the future^[4].

The open processes and “running code” that are hallmarks of the IETF, mean that MLS is already proven to be efficient at Internet scale, working efficiently with groups that have thousands of participants. MLS is already available from, and implemented and deployed by a wide range of companies and organizations^[5]. This includes real-time platforms such as Webex, Wire, and Discord, as well as in devices such as drones.

MLS is also extensible, meaning it can be easily updated in a number of ways. Work is continuing in the MLS Working Group in a number of areas and the IETF *More Instant Messaging Interoperability* (mimi)^[6] working group is looking to build on MLS as they aim to specify the minimal set of mechanisms required to make Internet messaging services interoperable.

[0] Richard Barnes, Benjamin Beurdouche, Raphael Robert, Jon Millican, Emad Omara, and Katriel Cohn-Gordon, “The Messaging Layer Security (MLS) Protocol,” RFC 9420, July 2023.

[1] Tom Van Pelt, GSMA, “RCS Encryption: A Leap Towards Secure and Interoperable Messaging,”
<https://www.gsma.com/newsroom/article/rcs-encryption-a-leap-towards-secure-and-interoperable-messaging/>

[2] IETF Messaging Layer Security Working Group:
<https://datatracker.ietf.org/wg/mls/about/>

- [3] Nick Sullivan and Sean Turner, “Messaging Layer Security: Secure and Usable End-to-End Encryption,” *IETF Blog*, March 29, 2023.
- [4] Rohan Mahy and Richard Barnes, “ML-KEM and Hybrid Cipher Suites for Messaging Layer Security,” Internet-Draft, Work in Progress, **draft-mahy-mls-pq-00**, March 2025.
- [5] “Support for MLS,” *IETF Blog*, July 18, 2023.
<https://www.ietf.org/blog/support-for-mls-2023/>
- [6] mimi Working Group:
<https://datatracker.ietf.org/group/mimi/about/>

ICANN and ISOC Joint Report on 20 Years of IGF

The *Internet Corporation for Assigned Names and Numbers* (ICANN) and the *Internet Society* (ISOC) recently released a report that offers a substantive look at the global impact of the *Internet Governance Forum* (IGF). It demonstrates how coordination—rather than control—has driven tangible progress in the Internet’s resilience, reach, and trust. Structured not as a retrospective but as a practical record of outcomes, the report draws from two decades of work across infrastructure, access, security, and policy. It offers grounded evidence of what coordination has made possible and what could be lost if support for multistakeholder cooperation erodes. The key insights of the report are as follows:

- *Infrastructure and Access*: In Africa, *Internet Exchange Points* (IXPs) more than doubled in over a decade. In countries like Kenya and Nigeria, this growth helped localize traffic, cutting the delay in data travel (latency) from around 200–600 milliseconds to 2–10 milliseconds, and saving millions annually in international connectivity costs. The IGF enabled the sharing of best practices that directly contributed to this expansion.
- *Multilingual Access*: Nearly 4.4 million domain names are now registered in non-Latin scripts. Through IGF-hosted sessions and stakeholder coalitions, *Internationalized Domain Names* (IDNs) and *Universal Acceptance* (UA) have gained critical momentum. In 2025, more than 50 global events marked *UA Day*, promoting linguistic access across the Internet ecosystem.
- *Security and Resilience*: Today, 93% of top-level domains are secured using *Domain Name System Security Extensions* (DNSSEC), which protect against forged DNS responses. In parallel, over 1,000 networks have adopted the *Mutually Agreed Norms for Routing Security* (MANRS), a global initiative to reduce routing attacks. The IGF has catalyzed awareness, collaboration, and implementation of these safeguards.
- *Local Engagement and Policy Influence*: More than 180 *National and Regional IGFs* (NRI) now form a decentralized backbone of year-round Internet governance dialogue. Initiatives like *Youth IGFs* and the *IGF Parliamentary Track* are shaping national and international policy—including formal declarations on digital trust, user rights, and multistakeholder governance.

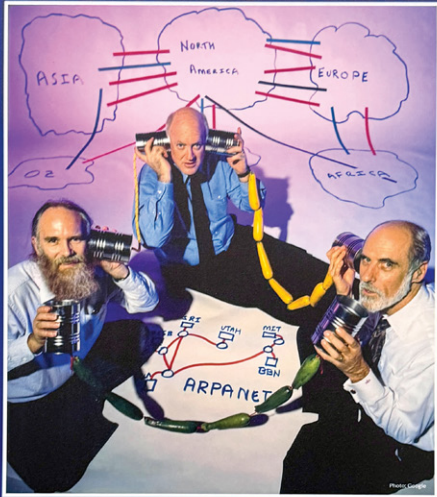
- *Community-Centric Innovation*: From the Arctic to the Andes, community networks have grown through IGF platforms and the *Dynamic Coalition on Community Connectivity* (DC3). These grass-roots efforts now inform regulatory change, including *International Telecommunication Union* (ITU) resolutions and national endorsements, and have helped close connectivity gaps in underserved regions.
- *Global Coordination Platform*: The IGF has evolved from an annual convening into a living ecosystem. It bridges technical and policy domains, connects local with global perspectives, and enables distributed but aligned Internet governance. That structure is now both a model and a necessity.



The report launches ahead of the 20-year review of the *World Summit on the Information Society* (WSIS+20)—a pivotal moment that will shape the next phase of global digital cooperation. It serves as both a record of achievement and a warning: coordination works but is not self-sustaining. The Internet's openness, security, and interoperability depend on it. If that cooperation falters, the conditions that have made the Internet thrive may not hold. The full report is available here:

https://www.internetsociety.org/wp-content/uploads/2025/06/20-Years-of-IGF_EN.pdf

Poster featuring Internet Pioneers Pål Spilling and Yngvar Lundh displayed at IGF 2025 in Lillestrøm, Norway.

1973



Pål Spilling, Yngvar Lundh and Dag Belsnes took part in the development of the TCP/IP protocol and established with Vint Cerf and his team the first connection to the Internet (Arpanet) at Kjeller – visit us on Tuesday and Thursday.

Thank You!

Publication of IPJ is made possible by organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol. The following individuals have provided support to IPJ. You can join them by visiting <http://tinyurl.com/IPJ-donate>

Kjetil Aas	Ilia Bromberg	Freek Dijkstra	Radu Cristian Gheorghiu	Nils Johansson
Fabrizio Accatino	Lukasz Bromirski	Geert Van Dijk	Greg Giessow	Brian Johnson
Michael Achola	Václav Brožík	David Dillow	John Gilbert	Curtis Johnson
Martin Adkins	Christophe Brun	Richard Dodsworth	Serge Van Ginderachter	Don Johnson
Melchior Aelmans	Gareth Bryan	Ernesto Doelling	Greg Goddard	Richard Johnson
Christopher Affleck	Ron Buchalski	Michael Dolan	Tiago Goncalves	Jim Johnston
Scott Aitken	Paul Buchanan	Eugene Doroniuk	Ron Goodheart	Jose Enrique Diaz Jolly
Jacobus Akkerhuis	Stefan Buckmann	Michael Dragone	Octavio Alfageme	Jonatan Jonasson
Antonio Cuñat Alario	Caner Budakoglu	Joshua Dreier	Gorostiaga	Daniel Jones
William Allaire	Darrell Budic	Lutz Drink	Barry Greene	Gary Jones
Nicola Altan	BugWorks	Aaron Dudek	Jeffrey Greene	Jerry Jones
Shane Amante	Scott Burleigh	Dmitriy Dudko	Richard Gregor	Michael Jones
Marcelo do Amaral	Chad Burnham	Andrew Dul	Martijn Groenleer	Amar Joshi
Matteo D'Ambrosio	Randy Bush	Joan Marc Riera	Geert Jan de Groot	Javier Juan
Selva Anandavel	Colin Butcher	Duocastella	Ólafur Guðmundsson	David Jump
Jens Andersson	Jon Harald Bøvre	Pedro Duque	Christopher Guemez	Anders Marius Jørgensen
Danish Ansari	Olivier Cahagne	Holger Durer	Rafael Leon Guerrero	Merike Kao
Finn Arildsen	Antoine Camerllo	Karlheinz Dölger	Gulf Coast Shots	Andrew Kaiser
Tim Armstrong	Tracy Camp	Mark Eanes	Galen Guyer	Vladislav Kalinovsky
Richard Artes	Brian Candler	Andrew Edwards	Sheryll de Guzman	Naoki Kambe
Michael Aschwanden	Fabio Caneparo	Peter Robert Egli	Rex Hale	Akbar Kara
David Atkins	Roberto Canonico	George Ehlers	Jason Hall	Christos Karayiannis
Jac Backus	David Cardwell	Peter Eisses	James Hamilton	Daniel Karrenberg
Jaime Badua	Richard Carrara	Torbjörn Eklöv	Darow Han	David Kekar
Bent Bagger	John Cavanaugh	Jacobus Gerrit Elsenaar	Handy Networks LLC	Stuart Kendrick
Eric Baker	Lj Cemerias	Y Ertur	Stephen Hanna	Robert Kent
Fred Baker†	Dave Chapman	ERNW GmbH	Martin Hannigan	Thomas Kernen
Santosh Balagopalan	Stefanos Charchalak	ESdatCo	John Hardin	Jithin Kesavan
William Baltas	Molly Cheam	Steve Esquivel	David Harper	Jubal Kessler
David Bandinelli	Christof Chen	Jay Etchings	Edward Hauser	Shan Ali Khan
A C Barber	Pierluigi Checchi	Mikhail Evstiounin	David Hauweele	Nabeel Khatri
Benjamin Barkin-Wilkins	Greg Chisholm	Bill Fenner	Marilyn Hay	Dae Young Kim
Ryan Barnes	David Chosrova	Paul Ferguson	Headcrafts SRLS	William W. H. Kimandu
Feras Batainah	Marcin Cieslak	Ricardo Ferreira	Hidde van der Heide	John King
Michael Bazarewsky	Lauris Cikovskis	Kent Fichtner	Johan Helsingius	Russell Kirk
Robert Beckett	Brad Clark	Ulrich N Fierz	Robert Hinden	Gary Klesk
David Belson	Narelle Clark	Armin Fisslthaler	Michael Hippert	Anthony Klopp
Richard Bennett	Horst Clausen	Michael Fiumano	Damien Holloway	Henry Kluge
Matthew Best	James Cliver	The Flirble Organisation	Alain Van Hoof	Michael Kluk
Hidde Beumer	Guido Coenders	Jean-Pierre Forcioli	Edward Hotard	Andrew Koch
Pier Paolo Biagi	Robert Collet	Gary Ford	Bill Huber	Ia Kochiashvili
Arturo Bianchi	Joseph Connolly	Susan Forney†	Hagen Hultzs	Carsten Koempe
John Bigrow	Steve Corbató	Christopher Forsyth	Kauto Huopio	Richard Koene
Orvar Ari Bjarnason	Brian Courtney	Andrew Fox	Asbjørn Højmark	Alexader Kogan
Tyson Blanchard	Beth and Steve Crocker	Craig Fox	Kevin Iddles	Matthijs Koot
Axel Boeger	Dave Crocker	Fausto Franceschini	Mika Ilvesmaki	Antonin Kral
Keith Bogart	Kevin Croes	Erik Fredriksson	Karsten Iwen	Robert Krejčí
Mirko Bonadei	John Curran	Valerie Fronczak	Joseph Jackson	John Kristoff
Roberto Bonalumi	Sergio Danelli	Tomislav Futivic	David Jaffe	Terje Krogdahl
Lolke Boonstra	André Danthine†	Laurence Gagliani	Ashford Jaggernauth	Bobby Krupczak
Cente Cornelis Boot	Morgan Davis	Edward Gallagher	Thomas Jalkanen	Murray Kucherawy
Julie Bottorff Photography	Jeff Day	Andrew Gallo	Jozef Janitor	Warren Kumari
Gerry Boudreaux	Nicholas Dean	Chris Gamboni	Martijn Jansen	George Kuo
Leen de Braal	Fernando Saldana	Xosé Bravo Garcia	John Jarvis	Dirk Kurfuerst
Stephen Bradley	Del Castillo	Oswaldo Gazzaniga	Dennis Jennings	Mathias Körber
Kevin Breit	Rodolfo Delgado-Bueno	Kevin Gee	Edward Jennings	Darrell Lack
Thomas Bridge	Julien Dhallenne	Rodney Gehrke	Aart Jochem	Andrew Lamb

Richard Lamb	Bart Jan Menkveld	Derrell Piper	Timothy Schwab	Sandro Tumini
Yan Landriault	Sean Mentzer	Rob Pirnie	Roger Schwartz	Angelo Turetta
Edwin Lang	Eduard Metz	Jorge Ivan Pincay	SeenThere	Brian William Turnbow
Sig Lange	William Mills	Ponce	Scott Seifel	Michael Turzanski
Markus Langenmair	David Millsom	Marc Vives Piza	Paul Selkirk	Phil Tweedie
Fred Langham	Desiree Miloshevic	Victoria Poncini	Andre Serralheiro	Steve Ulrich
Tracy LaQuey Parker	Joost van der Minnen	Blahoslav Popela	Yury Shefer	Unitek Engineering AG
Christian de Larrinaga	Thomas Mino	Andrew Potter	Yaron Sheffer	John Urbanek
Alex Latzko	Rob Minshall	Ian Potts	Doron Shikmoni	Martin Urwaleck
Jose Antonio Lazaro	Wijnand Modderman-	Eduard Llull Pou	Tj Shumway	Bart Vanautgaerden
Lazaro	Lenstra	Tim Pozar	Jeffrey Sicuranza	Betsy Vanderpool
Antonio Leding	Mohammad Moghaddas	David Preston	Thorsten Sideboard	Surendran Vangadasalam
Rick van Leeuwen	Charles Monson	David Raistrick	Greipur Sigurdsson	Ramnath Vasudha
Simon Leinen	Andrea Montefusco	Priyan R Rajeevan	Fillipe Cajaiba da Silva	Randy Veasley
Anton van der Leun	Fernando Montenegro	Balaji Rajendran	Andrew Simmons	Philip Venables
Robert Lewis	Roberto Montoya	Paul Rathbone	Pradeep Singh	Buddy Venne
Christian Liberale	Joel Moore	William Rawlings	Henry Sinnreich	Alejandro Vennera
Martin Lillepuu	Joseph Moran	Mujtiba Raza Rizvi	Geoff Sisson	Luca Ventura
Roger Lindholm	John More	Bill Reid	John Sisson	Scott Vermillion
Link Light Networks	Maurizio Moroni	Petr Rejhon	Helge Skrivervik	Tom Vest
Art de Llanos	Brian Mort	Robert Remenyi	Terry Slattery	Peter Villemoes
Mike Lochocki	Soenke Mumm	Rodrigo Ribeiro	Darren Sleeth	Vista Global Coaching &
Chris and Janet Lonvick	Tariq Mustafa	Glenn Ricart	Richard Smit	Consulting
Mario Lopez	Stuart Nadin	Justin Richards	Bob Smith	Dario Vitali
Sergio Loreti	Michel Nakhla	Rafael Riera	Courtney Smith	Marc Vives
Eric Louie	Mazdak Rajabi Nasab	Mark Risinger	Eric Smith	Rüdiger Volk
Adam Loveless	Krishna Natarajan	Fernando Robayo	Mark Smith	Jeffrey Wagner
Josh Lowe	Naveen Nathan	Michael Roberts	Tim Sneddon	Don Wahl
Guillermo a Loyola	Ryan Nelson	Gregory Robinson	Craig Snell	Michael L Wahrman
Hannes Lubich	Darryl Newman	Ron Rockrohr	Job Snijders	Lakhinder Walia
Dan Lynch [†]	Mai Nguyen	Graziano G Rodegari	Ronald Solano	Laurence Walker
David MacDuffie	Thomas Nikolajsen	Carlos Rodrigues	Asit Som	Randy Watts
Sanya Madan	Paul Nikolich	Magnus Romedahl	Ignacio Soto Campos	Andrew Webster
Miroslav Madić	Travis Northrup	Lex Van Roon	Evandro Sousa	Jd Wegner
Alexis Madriz	Marijana Novakovic	Marshall Rose	Peter Spekrijse	Tim Weil
Carl Malamud	David Oates	Alessandra Rosi	Thayumanavan Sridhar	Westmoreland
Jonathan Maldonado	Ovidiu Obersterescu	David Ross	Paul Stancik	Engineering Inc.
Michael Malik	Jim Oplotnik	William Ross	Ralf Stempfer	Rick Wesson
Tarmo Mamers	Tim O'Brien	Boudhayan	Matthew Stenberg	Peter Whimp
Yogesh Mangar	Mike O'Connor	Roychowdhury	Martin Štěpánek	Russ White
John Mann	Mike O'Dell	Carlos Rubio	Adrian Stevens	Jurrien Wijlhuizen
Bill Manning [†]	John O'Neill	Rainer Rudigier	Clinton Stevens	Joseph Williams
Diego Mansilla	Carl Ötne	Timo Ruiters	John Streck	Derick Winkworth
Harold March	Packet Consulting Limited	RustedMusic	Martin Streule	Pindar Wong
Vincent Marchand	Carlos Astor Araujo	Babak Saberi	David Strom	Brian Woods
Normando Marcolongo	Palmeira	George Sadowsky	Colin Strutt	Makarand Yerawadekar
Gabriel Marroquin	Gordon Palmer	Scott Sandefur	Viktor Sudakov	Phillip Yialeloglou
David Martin	Alexis Panagopoulos	Sachin Sapkal	Kathleen Summers	Janko Zavernik
Jim Martin	Gaurav Panwar	Arturas Satkovskis	Edward-W. Suor	Bernd Zeimetz
Ruben Tripiana Martin	Chris Parker	PS Saunders	Vincent Surillo	Muhammad Ziad
Timothy Martin	Alex Parkinson	Richard Savoy	Terence Charles Sweetser	Ziayuddin
Carles Mateu	Craig Partridge	John Sayer	T2Group	Tom Zingale
Juan Jose Marin Martinez	Manuel Uruena Pascual	Phil Scarr	Roman Tarasov	Matteo Zovi
Ioan Maxim	Ricardo Patara	Gianpaolo Scassellati	David Theese	Jose Zumalave
David Mazel	Dipesh Patel	Elizabeth Scheid	Rabbi Rob and	Romeo Zwart
Miles McCredie	Dan Paynter	Jeroen Van Ingen	Lauren Thomas	廖明沂.
Gavin McCullagh	Leif-Eric Pedersen	Schenau	Douglas Thompson	
Brian McCullough	Rui Sao Pedro	Carsten Scherb	Kerry Thompson	
Joe McEachern	Juan Pena	Ernest Schirmer	Lorin J Thompson	
Alexander McKenzie	Luis Javier Perez	Benson Schliesser	Jerome Tissieres	
Jay McMaster	Chris Perkins	Philip Schneek	Fabrizio Tivano	
Mark Mc Nicholas	Michael Petry	James Schneider	Peter Tomsu Fine Art	
Olaf Mehlberg	Alexander Peuchert	Peter Schoo	Photography	
Carsten Melberg	David Phelan	Dan Schrenk	Joseph Toste	
Kevin Menezes	Harald Pilz	Richard Schultz	Rey Tucker	

Call for Papers

The *Internet Protocol Journal* (IPJ) is a quarterly technical publication containing tutorial articles (“What is...?”) as well as implementation/operation articles (“How to...”). The journal provides articles about all aspects of Internet technology. IPJ is not intended to promote any specific products or services, but rather is intended to serve as an informational and educational resource for engineering professionals involved in the design, development, and operation of public and private internets and intranets. In addition to feature-length articles, IPJ contains technical updates, book reviews, announcements, opinion columns, and letters to the Editor. Topics include but are not limited to:

- Access and infrastructure technologies such as: Wi-Fi, Gigabit Ethernet, SONET, xDSL, cable, fiber optics, satellite, and mobile wireless.
- Transport and interconnection functions such as: switching, routing, tunneling, protocol transition, multicast, and performance.
- Network management, administration, and security issues, including: authentication, privacy, encryption, monitoring, firewalls, troubleshooting, and mapping.
- Value-added systems and services such as: Virtual Private Networks, resource location, caching, client/server systems, distributed systems, cloud computing, and quality of service.
- Application and end-user issues such as: E-mail, Web authoring, server technologies and systems, electronic commerce, and application management.
- Legal, policy, regulatory and governance topics such as: copyright, content control, content liability, settlement charges, resource allocation, and trademark disputes in the context of internetworking.

IPJ will pay a stipend of US\$1000 for published, feature-length articles. For further information regarding article submissions, please contact Ole J. Jacobsen, Editor and Publisher. Ole can be reached at ole@protocoljournal.org or olejacobsen@me.com

The Internet Protocol Journal is published under the “CC BY-NC-ND” Creative Commons Licence. Quotation with attribution encouraged.

This publication is distributed on an “as-is” basis, without warranty of any kind either express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, or non-infringement. This publication could contain technical inaccuracies or typographical errors. Later issues may modify or update information provided in this issue. Neither the publisher nor any contributor shall have any liability to any person for any loss or damage caused directly or indirectly by the information contained herein.

Follow us on X and Facebook



@protocoljournal



<https://www.facebook.com/newipj>

Supporters and Sponsors

Supporters



Internet
Society



Diamond Sponsors

Your logo here!

Ruby Sponsors



Sapphire Sponsors



Emerald Sponsors



Corporate Subscriptions



For more information about sponsorship, please contact sponsor@protocoljournal.org

The Internet Protocol Journal
Link Fulfillment
7650 Marathon Dr., Suite E
Livermore, CA 94550

CHANGE SERVICE REQUESTED

The Internet Protocol Journal

Ole J. Jacobsen, Editor and Publisher

Editorial Advisory Board

Dr. Vint Cerf, VP and Chief Internet Evangelist
Google Inc, USA

John Crain, Senior Vice President and Chief Technology Officer
Internet Corporation for Assigned Names and Numbers

Dr. Steve Crocker, CEO and Co-Founder
Shinkuro, Inc.

Dr. Jon Crowcroft, Marconi Professor of Communications Systems
University of Cambridge, England

Geoff Huston, Chief Scientist
Asia Pacific Network Information Centre, Australia

Dr. Cullen Jennings, Cisco Fellow
Cisco Systems, Inc.

Merike Kaeo, Founder and vCISO
Double Shot Security

Olaf Kolkman, Principal – Internet Technology, Policy, and Advocacy
The Internet Society

Dr. Jun Murai, Founder, WIDE Project
Distinguished Professor, Keio University
Co-Director, Keio University Cyber Civilization Research Center, Japan

The Internet Protocol Journal is published quarterly and supported by the Internet Society and other organizations and individuals around the world dedicated to the design, growth, evolution, and operation of the global Internet and private networks built on the Internet Protocol.

Email: ipj@protocoljournal.org
Web: www.protocoljournal.org

The title "The Internet Protocol Journal" is a trademark of Cisco Systems, Inc. and/or its affiliates ("Cisco"), used under license. All other trademarks mentioned in this document or website are the property of their respective owners.

Printed in the USA on recycled paper.

