

# **Application Delivery – An Enhanced Internet Based Solution**



**October 2007**

## Executive Summary

Up until a few years ago, enterprise wide area networks were deployed primarily in a hub and spoke design and utilized technologies such as Frame Relay and ATM<sup>1</sup>. Now a number of factors are causing IT organizations to move away from a hub and spoke network design and to adopt a meshed network infrastructure. In order to efficiently support a meshed network most IT organizations have begun to utilize either MPLS or the Internet. In key areas such as cost, complexity, lead-time and support for nomadic workers, Internet based services are superior to MPLS. One area where MPLS services have traditionally been superior to Internet-based services is that MPLS promises low predictable delay. Now, however, with the advent of new techniques that optimize routing and transport over the Internet, that advantage is diminishing.

One of the reasons why the Internet has not traditionally been able to provide low predictable delay is because the Internet uses the Border Gateway Protocol (BGP) to determine the routes from one subtending network to another. When choosing a route, BGP strives to minimize the number of hops between the origin and the destination. In particular, BGP does not strive to choose a route with the optimal performance characteristics (i.e., lowest delay, lowest packet loss). Given the dynamic nature of the Internet, a network or a particular peering point router can go through periods where they exhibit severe delay and/or packet loss. As a result, the route with the fewest hops is not necessarily the route with the best performance.

In the last few years, a number of techniques have been developed to accelerate the delivery of applications delivered over the Internet. For example, many IT organizations have deployed an appliance to mitigate the concerns about the performance of the servers in the data center. These devices are sometimes referred to as an Application Front End (AFE) or as an Application Delivery Controller (ADC). The use of an AFE alleviates some of the factors that limit the performance of the servers in the data center. An AFE, however, does not alleviate all of the data center performance issues, nor does it alleviate the concerns about the performance of applications that transit the Internet.

There are, however, some additional techniques that can alleviate the concerns about the performance of Internet based applications such as HTTPS, Citrix ICA, UDP, SSL based VPNs, and IPsec based VPNs. For example, optimization techniques such as compression, caching, pre-fetching and differencing can often improve the performance of these applications. Other techniques that improve the performance of these applications include route optimization, which determines the optimal route through the Internet, and transport optimization that eliminates some of the features of TCP that can cause poor application performance; i.e., the TCP slow start algorithm and TCP's retransmission timeout.

There is a strong synergy between route optimization and transport optimization. For example, because route optimization chooses the optimum path through the Internet, it is more likely than is BGP to choose a path that has minimum congestion and hence not experience the performance problems associated with TCP. In addition, because the path is optimized, it is possible to get

---

<sup>1</sup> For readability purposes, a few key acronyms will be defined in the main text of the white paper. The rest will be defined in Appendix A.

more aggressive with both the TCP slow start algorithm and TCP's retransmission timeout without incurring additional congestion.

Managed service providers (MSPs) have begun to deploy services that are applicable in two fundamentally different application delivery environments. One of these environments is the delivery of applications to branch offices using a WAN technology other than the Internet; i.e., Frame Relay, ATM, or MPLS. These services involve deploying a piece of equipment often referred to as a WAN Optimization Controller (WOC) in selected branch offices as well as in the data center.

The other environment being addressed by MSPs is the delivery of applications over the Internet. When choosing a managed service provider for delivering applications over the Internet there are some specific technological capabilities that IT organizations should look for. For example, it is critical that the MSP can perform route, transport and application specific optimization in an optimal fashion. In order to maximize the benefits of these optimization techniques an intelligent distributed infrastructure is required. For example, the determination of the best route through the Internet requires dynamic information on the performance characteristics of the path, not just between the origin and the destination, but also between numerous intermediary points that are between the origin and termination. An intelligent distributed infrastructure is not required in order to implement TCP optimization techniques. However, the beneficial impact of these techniques is magnified if the delay and packet loss of the Internet is minimized. Minimizing this delay and packet loss requires an intelligent distributed infrastructure.

In addition, the most effective way to implement application optimization techniques such as caching is to implement them as close to the user as possible so that the information can be delivered to the user with minimum delay. This requires an intelligent distributed infrastructure so that there are servers that are close to the end user. It also requires the ability to establish a connection between the user and the optimal server based on factors such as the real-time Internet conditions and the server load.

## **Introduction**

In his recent book entitled "The World is Flat", Thomas Friedman discussed how the traditional geographical barriers to commerce are falling in large part due to the deployment of technologies that enable people around the world to communicate easily and effectively. Friedman listed the PC, the Internet and Web browsers as examples of this type of technology. As Friedman pointed out, these technologies allow for the development of large, complex global supply chains.

In line with Friedman's comments, many enterprises are beginning to use the Internet as an alternative to traditional network services such as Frame Relay and ATM. As a result, Internet based applications have become quite popular. Enterprise IT organizations now commonly develop IP-based applications themselves or acquire them from myriad software vendors including Microsoft, IBM and SAP. There are many compelling factors driving the popularity of Internet based applications. First, by implementing Internet based applications companies can achieve Friedman's vision of extending their reach to virtually anybody in the world. Second, since Internet based applications often rely on the ubiquitous browser and not specialized client software, the IT resources required to deploy and manage the desktop are significantly reduced.

Third, the majority of the IT resources required to support Internet based applications (i.e., servers, databases, operating systems) are located in a centralized site. The centralization of IT resources reduces the burden of deploying, managing and securing these resources.

Until recently, it was difficult to find a network organization with responsibility for ensuring acceptable application performance. Now that concern is a top of mind issue for virtually all network organizations. The primary reason for the dramatic shift in the importance that network organizations place on ensuring acceptable application performance is that companies increasingly run their business operations over the network. As a result, if there is a performance problem of some kind, key business operations such as sales and customer service are impacted. However, the complicating factor is that ensuring acceptable application performance is extremely difficult.

One of the goals of this white paper is to identify why the Internet is beginning to be viewed as an acceptable service for transporting business-critical, delay-sensitive traffic. Another goal is to analyze the delivery of Internet based applications from an architectural perspective. This white paper will first identify the primary functional components of an Internet based application delivery system, describe the bottlenecks associated with each component and suggest the functionality that should be present in each component to alleviate these bottlenecks.

A final goal of this white paper is to discuss the viability of using a managed service provider to deliver applications over the Internet. As will be demonstrated, one of the primary reasons why using a managed service provider makes sense for so many IT organizations is that managed service providers are well positioned to deploy the functionality that can alleviate the bottlenecks in each of the components of an Internet based application delivery system.

## **The Changing Wide Area Network Environment**

Historically data networks have been built using some form of a hub-and-spoke design and were based on technologies such as Frame Relay and ATM. Hub-and-spoke designs have been common in large part because that design reflected what had been the natural flow of traffic between a branch office and a headquarters site. However, many factors are changing the traffic flow in data networks and hence causing IT organizations to move away from simple hub-and-spoke network designs. One of these factors is that the natural flow of data traffic is changing. For example, today branch offices typically need access to multiple data centers, either for disaster recovery or for access to applications that are only hosted in one of the company's multiple data centers. This type of traffic is often characterized as being one-to-many. Another factor is that the vast majority of companies have deployed VoIP and voice traffic does not tend to follow a hub-and-spoke pattern. Voice traffic tends to follow an any-to-any traffic pattern.

It is certainly possible to support one-to-many and any-to-any traffic using Frame Relay and ATM. This approach, however, means that as the traffic flows from the origin to the destination, it will transit through an intermediate point. If a large percentage of the traffic has to transit through an intermediate point this tends to force the organization to deploy a significant amount of additional WAN bandwidth, which adds cost. This transit traffic also consumes significant additional resources on the organization's routers, which also adds cost. The majority of IT organizations have also come to recognize that Frame Relay and ATM are legacy technologies,

and that few service providers or equipment vendors are investing in these technologies. This lack of investment means that the price to performance ratio of these technologies will not improve, and that it is highly unlikely that any new significant features or management functionality will be added to these technologies.

The shift away from hub and spoke network designs combined with the growing awareness that Frame Relay and ATM are legacy technologies has caused many IT organizations to evaluate alternative WAN transmission services. One alternative that many IT organizations have expressed an interest in is MPLS. MPLS has garnered a lot of attention over the last few years in part because the vast majority of the major carriers have implemented MPLS within their backbone networks.

Like any WAN service, MPLS has advantages and disadvantages. One of the advantages of MPLS is that it is widely deployed. One of the disadvantages of MPLS is that similar to Frame Relay and ATM, MPLS services tend to be expensive. Another disadvantage is that there tends to be a long lead-time associated with deploying new MPLS services.

In addition to being widely deployed, another one of the advantages of MPLS is that it efficiently supports meshed traffic whether the traffic is one-to-many or many-to-many. One of the other characteristics of MPLS that gets a lot of attention is MPLS' ability to support differing service classes. Appendix B contains a table that illustrates a representative set of MPLS service classes.

One observation that can be made from the table in appendix B is that MPLS is complex. That fact was reflected in recent market research on MPLS services<sup>2,3</sup>. The complexity of MPLS has precluded virtually all IT organizations from deploying private MPLS networks. In addition, the complexity of MPLS has meant that when an IT organization acquires MPLS services from a carrier they usually acquire a managed service.

In addition to the complexity associated with MPLS, the SLAs associated with MPLS services tend to be weak which mitigates to some degree the promise of supporting disparate service classes. For example, the SLAs are primarily reactive in focus. The computation of an outage begins only when the customer opens a trouble ticket. Not only is the computation of the SLA metrics done in a way that is not favorable to the customer, the level of compensation for violation of service level agreements remains quite modest. Further weakening the value of these SLAs is the fact that the SLA metrics are primarily calculated as network-wide averages rather than for a specific customer's traffic or for a particular site. As a result, it would be possible for a company's data center to receive notably poor service in spite of the fact that the network-wide SLA metrics remain within agreed-upon bounds. In addition, many of the service providers have some unique quirks in their service level agreements. For example, one service provider excludes from their availability target any network outage of less than a minute in duration.

---

<sup>2</sup> Innovation in MPLS Based Services <http://www.webtorials.com/main/resource/papers/kubernan/paper2/MPLS-Innovation.pdf>

<sup>3</sup> The Cost and Management Challenges of MPLS

[http://www.webtorials.com/main/resource/papers/netscout/briefs/05-06/0506\\_Cost\\_MPLS.pdf](http://www.webtorials.com/main/resource/papers/netscout/briefs/05-06/0506_Cost_MPLS.pdf)

The general availability of MPLS services in advanced countries is extremely high. If an organization, however, needs to provide access in emerging countries, MPLS access may not be available. Gaining access to MPLS services in emerging countries is made more difficult by the fact that the majority of service providers rarely implement a Network-to-Network Interface (a.k.a., NNI) between themselves and other service providers. As a result, if the service provider of choice does not provide MPLS service in a particular country, but there is a local provider in the country that does, those local services will typically not be usable.

An area where MPLS is notably weak is the support of home and nomadic workers. It is not possible to use MPLS to support nomadic workers who need connectivity from virtually anywhere, such as a hotel room, a coffee shop, or an airport. And, while it is possible to use MPLS to support home workers, it is typically prohibitively expensive.

Since the Internet also supports meshed traffic flows, it also is being used by many enterprises as an alternative to frame relay and ATM. Table 1 compares MPLS and Internet services using a wide range of criteria. As seen in Table 1, with one exception the Internet provides service that is equal or better than MPLS. That exception is the ability to provide low, predictable delay. As will be demonstrated in the next section of the white paper, there are techniques that can be applied that greatly improve the ability of the Internet to provide low, predictable delay.

<b>Criteria</b>	<b>MPLS</b>	<b>Internet</b>
<b>Complexity</b>	High	Low
<b>Low predictable delay</b>	Part of the service description, but weak SLAs are the norm	Not usually possible without enhancements
<b>Availability</b>	Extremely high in advanced countries. Can be spotty in emerging countries.	Ubiquitous
<b>Lead time for new service</b>	Can be lengthy	Typically short
<b>Support for meshed traffic</b>	High	High
<b>Cost</b>	High	Low
<b>Reliability</b>	High	High
<b>Security</b>	High	Can be made very high by using protocols such as HTTPS and SSL
<b>Support of home and nomadic workers</b>	Very low	Very high

**Table 1: Comparison of MPLS and Internet Services**

## **An Internet Based Application Delivery System**

The two primary physical components of an Internet based application delivery system are the data center that hosts the applications and the end-to-end network that is used to support the communications between the end user and the application. These two components will be referred to in this white paper as the Host Data Center and the Network. There is a third component of an Internet based application delivery system that is more logical than physical. The third component is the type of information that has to flow between the end user and the

application servers. As will be discussed, in some instances a number of large objects need to be sent from the application servers to the users. Alternatively, in some other instances, the information flow between the user and the application server consists of hundreds of short transmissions. Each type of information flow presents a unique performance challenge.

One of the primary bottlenecks associated with the Host Data Center concerns the utilization of the servers in the data center. One example of this involves situations in which there are multiple servers for a given application. Based on how the users are assigned to a server, some servers could sit idle while other servers are over-utilized and are hence providing poor response. Another example involves the use of SSL. There is a lot of overhead associated with processing SSL traffic. If there is only a small amount of SSL traffic, there is no problem processing it on the server. However, if there is a lot of SSL traffic, all of the server's resources could be consumed with overhead processing leaving no processing resources to support the application.

Another bottleneck associated with the Host Data Center concerns the provisioning of servers. In particular, most application usage is peaked. Application usage is notably higher during a few hours of the day than it is during the rest of the day. The peaked nature of application usage has traditionally presented IT organizations with two alternative approaches to provisioning servers, neither one of which is highly desirable. One alternative is to provision enough servers to support the maximum application demand. The advantage of this approach is that it ensures that the lack of servers is not a cause of poor application performance. The disadvantage of this approach is that it tends to be very costly as the servers sit underutilized most of the time and the additional servers drive the need for additional software licenses, real estate, power, and maintenance. The second alternative is to provision fewer servers than are necessary to support the maximum demand. The advantage of this approach is that it is less costly than the alternative approach. The disadvantage of this approach is that it guarantees that applications will perform poorly during times of peak demand.

The Network actually comprises three different components:

- The First Mile

This refers to the connection between the end user and the Middle Mile (see below) and is typically some form of private line. While this may be more than a mile in length, it is typically a short connection.

- The Middle Mile

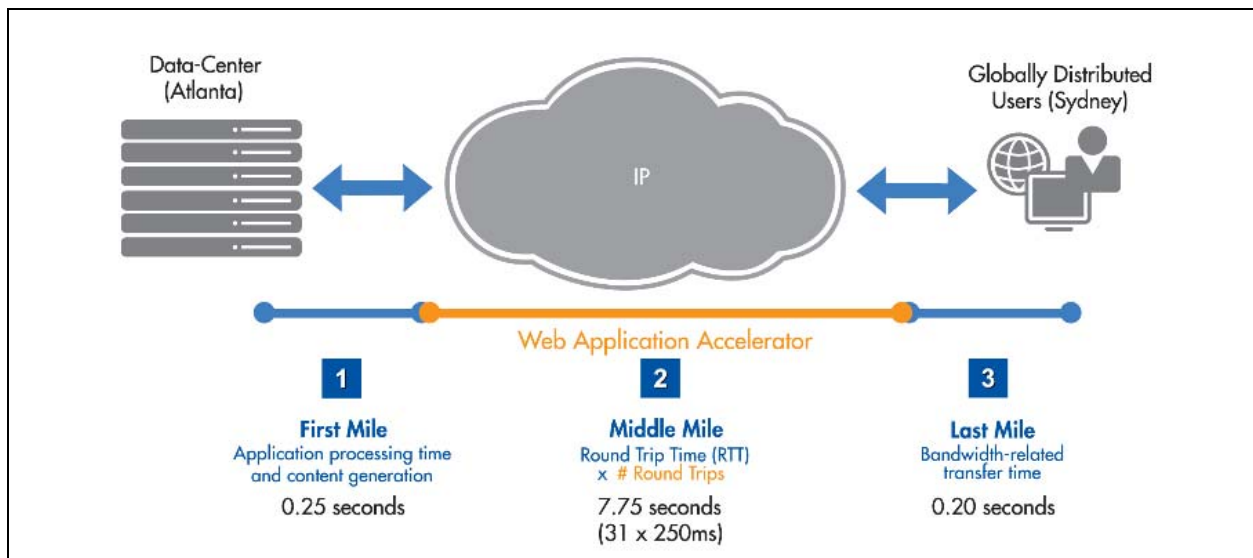
This refers to the connectivity through the Internet. This is typically hundreds or even thousands of miles in length.

- The Last Mile

This refers to the connectivity from the Middle Mile to the Host Data Center Internet and is typically some form of private line. While this may be more than a mile in length, it is typically a short connection.

Given the deployment of high-speed transmission services, the First Mile and the Last Mile seldom add to the problem of ensuring acceptable application performance. The Middle Mile, however, is often the source of application degradation. One of the primary bottlenecks associated with the Middle Mile is the variability associated with the Internet. In particular, the Internet uses BGP to determine the routes from one subtending network to another. When choosing a route, BGP strives to minimize the number of hops between the origin and the destination. In particular, BGP does not strive to choose a route with the optimal performance characteristics; i.e., lowest delay, lowest packet loss. Given the dynamic nature of the Internet, a given network or a particular peering point router can go through periods where it exhibits severe delay and/or packet loss. As a result, the route that has the fewest hops is not necessarily the route that has the best performance.

Figure 1 demonstrates the performance issues associated with the Middle Mile. In particular, Figure 1 depicts the situation of a user in Sydney, Australia accessing a Web page in Atlanta, GA. The Web page contains 25 objects and a total of 70 kilobytes of data. It takes 31 round trips and 8.4 seconds to download the Web page. The Middle Mile consumed ninety-two percent of the 8.4 seconds.



**Figure 1: Impact of the Middle Mile**

The performance issues associated with the Middle Mile are not new. They are, however, having more of an impact currently because of some trends in business and IT. One of these trends, the globalization of business, was referred to in the introduction. Because of the cost, time to deploy, and reach associated with the Internet, many IT organizations use the Internet to support globalization. The Middle Mile for a global network tends to be lengthy and complex. Both of these characteristics tend to cause degraded network performance that typically results in degraded application performance.

Another of these trends is that many companies either already have, or are in the process of consolidating servers out of branch offices and into centralized data centers. This consolidation typically reduces cost and enables IT organizations to have better control over the company's

data. It does mean, however, that employees now must use a WAN to access applications that they used to access over a LAN.

In addition to consolidating servers out of branch offices and into centralized data centers, many companies are also reducing the number of data centers they support worldwide. Many of these companies are also adopting a single-hosting model whereby users from all over the globe transit the WAN to access an application that the company hosts in just one of its data centers. Both data center consolidation and single hosting increase the distance and complexity of the Middle Mile between remote users and the applications they need to access. As previously stated, both of these characteristics tend to cause degraded network and application performance.

Another trend that contributes to poor application performance is the deployment of chatty protocols such as HTTP (Hypertext Transfer Protocol). A chatty protocol requires hundreds or even thousands of round trips or applications turns to complete a single transaction. Assume that the round-trip delay through the middle mile is 150 ms and that for a particular transaction HTTP requires one hundred round trips. As such, the Middle Mile adds fifteen seconds to the overall application delay.

## **Optimization Techniques**

This section will discuss the techniques that can be applied to each component of an Internet based application delivery system to ensure better application performance.

### **The Host Data Center**

Many IT organizations have deployed an appliance to mitigate the concerns about the utilization of the servers in the Host Data Center. These devices are sometimes referred to as an Application Front End or as an Application Delivery Controller.

The genesis of this category of solution dates back to the IBM mainframe-computing model of the late 1960s and early 1970s. Part of that computing model was to have a Front End Processor (FEP) reside in front of the IBM mainframe. The primary role of the FEP was to free up processing power on the general purpose mainframe computer by performing communications processing tasks in a device that was designed just for these tasks.

Analogous to a FEP, an AFE is intended to maximize the performance and availability of servers. To accomplish this goal, the AFE combines a number of functions, including:

- Server Load Balancing (SLB) to maximize the scalability and the availability of an application
- Layer 4 - Layer 7 switching to direct application queries to the most appropriate server
- Firewall functionality to ensure the integrity of the company's data and to provide application-specific security
- Off-loading computationally intensive tasks, such as the processing of SSL traffic

- Switch, accelerate and secure XML applications and web services

An AFE often contains functionality referred to as a *reverse cache*. One of the primary roles of a reverse cache is to store frequently accessed content. Since users retrieve this content from the reverse cache and not from the servers, the performance of the servers improves. However, this information still needs to transit the Network in order to be delivered to the end user. As a result, the value of a reverse cache reflects the overall value of an AFE. That value being that an AFE alleviates some of the concerns about the Host Data Center. An AFE, however, does not alleviate all of the concerns about the Host Data Center, nor does it alleviate the concerns about the Network in general, or the Middle Mile in particular.

## **The Middle Mile**

A previous section pointed out that in most ways, the Internet is superior to MPLS. The only significant weakness of the Internet is that it can't be counted on to provide low, predictable delay. This section of the white paper will discuss a set of techniques that can be used to improve the performance of the Middle Mile.

### Route Optimization

As was previously mentioned, one of the primary bottlenecks associated with the Middle Mile is the variability associated with the Internet. In particular, when BGP determines the routes from one subtending network to another, it does not take into account delay or packet loss. A few years ago, a new industry category appeared that is referred to as route optimization. The goal of route optimization is to improve the performance of IP networks. To accomplish this goal, a route optimization solution measures the performance of multiple paths through the network and chooses the optimum path from origin to destination.

### Transport Optimization

TCP is the most common transport protocol. Unfortunately TCP has a number of parameters that can cause poor application performance. One of these characteristics is TCP's retransmission timeout. This parameter controls how long the transmitting device waits for an acknowledgement from the receiving device before assuming that the packets were lost and need to be retransmitted. If this parameter is set too high, it introduces needless delay as the transmitting device sits idle waiting for the timeout to occur. Conversely, if the parameter is set too low, it can increase the congestion that was the likely cause of the timeout occurring.

Another important TCP parameter is the TCP slow start algorithm. The slow start algorithm is part of the TCP congestion control strategy and it calls for the initial data transfer between two communicating devices to be severely constrained. The algorithm calls for the data transfer rate to increase if there are no problems with the communications. In addition to the initial communications between two devices, the slow start algorithm is also applied in those situations in which a packet is dropped.

Both of these parameters are part of TCP because they provide value. As a result, these parameters can't be just blindly ignored. However, TCP performance can be significantly

improved if these parameters can be set dynamically based on the characteristics of the network such as the speed of the links and the distance between the transmitting and receiving devices.

There is a strong synergy between route optimization and transport optimization. For example, because route optimization chooses the optimum path through the Middle Mile, it is more likely than is BGP to choose a path that has minimum congestion and hence not experience the problems associated with the TCP slow start algorithm or with TCP's retransmission timeout. In addition, because the path is optimized, it is possible to get more aggressive with both the TCP slow start algorithm and TCP's retransmission timeout without incurring additional congestion.

### Application Specific Optimization

There are several optimization techniques that focus on the applications and the subtending protocols. This includes compression, caching, pre-fetching and differencing.

The goal of compression techniques is to reduce the amount of data that must be transmitted across the WAN. As such, implementing compression has the affect of a onetime increase in the capacity of the communication channel. Data compression algorithms find redundancy in a data stream and use encoding techniques to remove the redundancy and create a smaller file. A number of familiar lossless compression tools for binary data are based on Lempel-Ziv compression. This includes zip, PKZIP and gzip algorithms.

Similar to compression, the goal of caching is to reduce the amount of data that must be transmitted. One common example of caching is a Web cache appliance that serves as a proxy Web server, caching frequently accessed objects from Web pages that are hosted remotely over the Enterprise WAN or on the Internet. Client browsers are configured to connect to the proxy server rather than directly to the remote server. When the request for a remote object is made, the local cache is queried first. If the cache contains a current version of the object, the request can be satisfied locally at LAN speed and with minimal latency. Most of the latency involved in a cache hit results from the cache querying the remote server to ensure that the cached object is still current. If the cache does not contain a current version of the object, that object must be fetched by the cache, cached, and then forwarded to the requestor.

Understanding an application's data objects as well as the application's semantics enables an application optimization solution to recognize application data and anticipate its use. This capability allows the solution to make requests of a distant server prior to those requests being made by the client. This is typically referred to as pre-fetching. Pre-fetching allows an application optimization solution to respond locally when the client does make those requests.

The goal of differencing is to avoid sending an entire file from origin to destination. In particular, the goal of differencing is to only send the changes that have been made to the file since the last time it was sent. To achieve this goal, differencing algorithms describe a new version of a file as a set of changes to a previous version of the same file, where the previous version of the file is often referred to as the base version. All binary differencing algorithms partition a file into two classes of variable length byte strings: those strings that appear in the base version and those that are unique to the version being encoded. The latter strings are included in what is often referred to as a delta file. The delta file contains the minimum set of changes that the receiver at the destination needs in order to build the new version of the file.

## The Application Delivery Challenges

This section will identify some of the typical applications that an enterprise uses. Route optimization improves the performance of all of these applications as each of these applications runs better over a WAN with low delay and packet loss. In addition, every application that uses TCP will benefit from TCP optimization.

As part of the description of the typical applications that an enterprise uses, this section will describe some of the primary performance challenges that are associated with these applications. This is done both to exemplify why route and transport optimization add value as well as to identify some of the key ways in which other optimization techniques can improve the performance of these applications.

### HTTPS

HTTPS, which is based on HTTP, uses a different default TCP port (443) than does HTTP and adds either an SSL or a TLS (Transport Layer Security) encryption and authentication layer between HTTP and TCP. Because of the additional security, HTTPS is essential for eCommerce sites and to protect enterprise resources accessed over the Internet.

For HTTPS sites, the CPU intensity of the SSL handshake process as well as the encryption and decryption process adds server latency and increases CPU utilization. This reduces a server's capacity to serve pages, and to perform other functions. As was previously mentioned, some IT organizations deploy an AFE in their data centers in part to handle the key exchange as well as the encryption and decryption of transmitted data that is associated with HTTPS.

As was depicted in Figure 1, a typical web page consists of numerous objects that must be transmitted sequentially via HTTPS from the web server to the browser. If the Middle Mile has low latency and low packet loss this reduces the time it takes to transmit these objects. However, the performance of HTTPS is also enhanced by application optimization techniques such as caching, compression and pre-fetching.

### Citrix ICA

Citrix Presentation Server provides access to enterprise applications via a thin client on local or remote PCs. The Citrix ICA protocol runs over TCP/IP and similar to HTTP is a chatty protocol. For example, a typical Citrix ICA transaction may require tens or even hundreds of exchanges. Each exchange consists of a small window of packets (the ICA default is two packets per window) and requires a round trip across the network.

The performance of Citrix ICA, and other chatty protocols, is therefore extremely sensitive to the round trip latency of the Middle Mile. Hence, route optimization and transport optimization play a key role in improving the performance of ICA. However, unlike the web pages being transmitted over HTTPS, ICA is typically concerned with small amounts of information that is already compressed. As such, applying additional application specific techniques such as compression typically offer little value and could actually degrade performance.

## UDP

UDP is a connectionless transport layer protocol that does not have the connection setup, flow control, or retransmission features that are provided by TCP. RTP (Real-time Transport Protocol) running over a combination of UDP and IP is widely used by real-time applications, such as VoIP and IP video, which require steady streams of traffic to flow between the endpoints. UDP is also used by some data applications, such as NFS, that were designed to maximize the use of LAN bandwidth, but have migrated to the WAN.

Real-time UDP applications are very sensitive to latency and jitter, as well as packet loss beyond some small threshold. On the enterprise LAN and WAN, QoS features can be enabled to protect the service quality of real-time traffic. When running over the Middle Mile, UDP applications will work well as long as there is adequate bandwidth and no significant degree of congestion along the end-to-end path between sender and receiver. Congestion leads to increased latency and jitter and typically also leads to packet loss. Since it is not possible to recover missing packets in real-time streams, packet loss results in gaps in the voice or video that can result in complete disruption of the communication session.

## SSL Remote Access VPNs

Remote access VPNs based on SSL technology are rapidly increasing in popularity. As large enterprises adopt SSL VPNs, the performance and scalability of centralized solutions can become an issue. As noted earlier, the CPU intensity of SSL processing can lead to performance problems and this has driven some IT organizations to deploy an AFE.

The performance of virtually any application running over an SSL VPNs is enhanced by both route and TCP optimization. Based on what the remote user is doing, the use of application specific optimization techniques may also be helpful. For example, assume that the remote user is attempting to download the company's inventory file. This process will be greatly enhanced if instead of sending the entire file, only the differences since the last time the file was downloaded are sent. To further enhance this process, the differences can be compressed before being transmitted.

## IPsec Based VPNs

There are two general classes of IPsec-Based VPNs. One class is focused on remote users and the other class is focused on branch office connectivity. Unlike SSL-based VPNs, IPsec based VPNs require client software on each PC. In spite of the added complexity of managing this client software, IPsec-Based VPNs are quite popular, particularly for providing remote access services over the Internet. IPsec uses Public Key Infrastructure (PKI) to establish an IP tunnel between two endpoints. After initial authentication and tunnel setup, all the IP packets are encrypted between the two tunnel endpoints.

The performance concerns associated with IPsec-based VPNs are very similar to the performance concerns associated with SSL Remote Access VPNs. In particular, the processing overhead associated with IPsec-based VPNs has led some IT organizations to deploy appliances

to perform the overhead processing. As is the case with SSL VPNs, the performance of virtually any application running over an IPsec-Based VPN is enhanced by route optimization and would also be enhanced by TCP optimization if the VPN runs over TCP. The use of application specific optimization techniques may also be helpful. For example, assume that a remote access user was attempting to download a large PowerPoint presentation. Compressing the presentation will significantly reduce the amount of time that it takes to download the file.

## The Use of Managed Services

The phrase *managed services* refers to the use of a third party to provide some combination of PDIM (plan, design, implement, and manage) services for a wide range of IT functionality. As previously mentioned, the majority of IT organizations that have deployed MPLS have deployed a managed service whereby the carrier will provide full PDIM services for the router on the customer's premise.

Managed service providers (MSPs) have begun to deploy services that are applicable in two fundamentally different application delivery environments. One of these environments is the delivery of applications to branch offices using a WAN technology other than the Internet (i.e., Frame Relay, ATM, or MPLS). These services involve deploying a piece of equipment often referred to as a WAN Optimization Controller (WOC) in selected branch offices as well as in the data center. WOCs implement a variety of technologies such as compression, caching and protocol acceleration and do so in ways that are proprietary. As part of these services, the MSP typically provides full PDIM services for the WOCs that are deployed in both the branch offices as well as the data center.

The other environment being addressed by MSPs is the delivery of applications over the Internet. There are many advantages to using a managed service provider for delivering applications over the Internet. This includes using an MSP who has:

- More in-house expertise

The network, security and application issues associated with ensuring acceptable application performance over the Internet are quite complex. In most cases, MSPs will have more expertise than an enterprise IT organization will have.

- Better in-house technology

In many cases, MSPs will have developed technology designed specifically to solve the problems associated with ensuring acceptable application performance over the Internet. For example, an MSP may have developed technology that mitigates the weakness associated with BGP and which chooses the route through the Middle Mile that has the lowest delay and packet loss.

- Better processes

Because an MSP performs services for thousands of customers, they are likely to have spent the time and resources to have developed and implemented effective processes.

- Lower Costs

Even if an IT organization had all of the capabilities of an MSP, the MSP would likely have a lower cost structure. This lower cost structure is due to the economies of scale inherent in developing the capability to successfully delivering applications over the Internet and then using this capability to service thousands of enterprises.

When choosing a managed service provider for delivering applications over the Internet there are some specific technological capabilities that IT organizations should look for. For example, previous sections of this white paper discussed optimization techniques (route, TCP, and application specific) and demonstrated the importance of these techniques to the performance of the typical applications that run over the Internet. Given the importance of these techniques, it is important that the service provider be able to provide them in an optimal fashion.

In order to maximize the benefits of these optimization techniques, an intelligent distributed infrastructure is required. For example, the determination of the best route through the Middle Mile requires dynamic information on the performance characteristics of the path, not just between the origin and the destination, but also between numerous intermediary points that are between the origin and termination. This level of granular information is needed so that the best end-to-end path can be determined. An intelligent distributed infrastructure is not required in order to implement TCP optimization techniques. However, the beneficial impact of these techniques is magnified if the delay and packet loss of the Middle Mile is minimized. Minimizing this delay and packet loss requires an intelligent distributed infrastructure.

As was previously mentioned, many AFEs function as a reverse cache. However, the information that is stored in the reverse cache still has to transit the Network. A more effective way to implement application optimization techniques such as caching is to implement them as close to the user as possible. In this way, the information can be delivered to the user with minimum delay.

Storing the information as close as possible to the user requires an intelligent distributed infrastructure so that there are servers that are close to the end user. It also requires the ability to establish a connection between the user and the optimal server based on factors such as the real-time Internet conditions and the server load.

## **Summary and Conclusions**

The Internet has proven itself to be an extremely valuable enabler of business communications. Some businesses, however, have been reticent to use the Internet for applications that are business-critical and delay sensitive. Instead, these businesses often use a managed MPLS service for these applications.

MPLS is a well-established and widely respected technology. However, in many ways the Internet is notably superior to MPLS. This includes the fact that the Internet:

- Is less complex than MPLS
- Has a much shorter lead time for the deployment of new service
- Is much less expensive
- Is capable of supporting home and nomadic workers

The reason that some businesses are reticent to use the Internet for applications that are business-critical and delay sensitive is that the Internet has not traditionally been able to provide the low predictable delay that these applications require. However, over the last few years a number of products and technologies have been developed that can significantly improve application performance over the Internet. It is possible for enterprise IT organizations to deploy some of these products and technologies on their own. However, a number of these technologies require the deployment of an intelligent distributed infrastructure. This is beyond the ability of virtually all enterprise IT organizations.

MSPs have begun to deploy services that address the issues associated with running business critical, delay sensitive applications over the Internet. The use of an MSP can have several advantages, including that the MSP:

- Has more in-house expertise
- Has better in-house technology
- Has better processes
- Offers lower costs

Only a miniscule number of IT organizations have deployed their own MPLS network. Instead, virtually all IT organizations acquire MPLS as a service from a carrier. In addition, the majority of IT organizations choose to let the carrier manage the customer premise routers. Given how IT organizations acquire MPLS services, it should not be much of a leap for these organizations to gain the advantages of the Internet by acquiring Internet based services from an MSP.

## Appendix A

ATM	Asynchronous Transfer Mode
ADC	Application Delivery Controller
AFE	Application Front End
BGP	Border Gateway Protocol
CPU	Central Processing Unit
FEP	Front End Processor
HTTP	Hypertext Transfer Protocol
HTTPS	HTTP over SSL
ICA	Independent Computing Architecture
IP	Internet Protocol
IT	Information Technology
LAN	Local Area Network
MPLS	Multi-Protocol Label Switching
MSP	Managed Service Providers
NFS	Network File System
NNI	Network-to-Network Interface
PC	Personal Computer
PDIM	Plan, Design, Implement, Manage
PKI	Public Key Infrastructure
QoS	Quality of Service
RTP	Real-time Transport Protocol
SLA	Service Level Agreement
SLB	Server Load Balancing
SSL	Secure Sockets Layer
TCP	Transmission Control Protocol
TLS	Transport Layer Security
UDP	User Datagram Protocol
VoIP	Voice over IP
VPN	Virtual Private Network
WAN	Wide Area Network
WOC	WAN Optimization Controller
XML	eXtensible Markup Language

## Appendix B

The following table illustrates a representative set of MPLS service classes.

<b>Class of Service (CoS)</b>	<b>Description</b>
CoS 1	This class is indicated with DSCP Expedited Forwarding (EF) and is intended for real-time applications such as interactive voice or video.
CoS2 (In Contract)	This class is indicated with DSCP Assured Forwarding 31 (AF31) and is intended for time sensitive, mission critical, low bandwidth, bursty data applications.
CoS2 (Out of Contract)	This class is indicated with DSCP Assured Forwarding 32 (AF32) and is intended for time sensitive, low bandwidth, bursty data applications. CoS2/InContract and CoS2/OutOfContract are serviced via the same queue. As such, they will have the same delay characteristics across the network. The difference is that in the event of severe congestion within CoS2, 'Out of Contract' class packets will be dropped first, allowing 'In Contract' CoS2 applications to be maintained.
CoS3 (In Contract)	This class is indicated with DSCP Assured Forwarding 21 (AF21) and is intended for time sensitive, mission critical, bursty data applications.
CoS3 (Out of Contract)	This class is indicated with DSCP Assured Forwarding 22 (AF22) and is intended for time sensitive, bursty data applications. CoS3/InContract and CoS3/OutOfContract are serviced via the same queue. As such, they will have the same delay characteristics across the network. The difference is that in the event of severe congestion within CoS3, 'Out of Contract' class packets will be dropped first, allowing 'In Contract' CoS3 applications to be maintained.
CoS 4	This class is indicated with DSCP default (default). It is also referred to as the best-effort class and is intended for all bulk data applications and non-time critical applications.

### Representation MPLS Service Classes