



Myths of Bandwidth Optimization

Increased use of IT in business, disaster recovery efforts, cloud computing, data center consolidation projects, and international presence all contribute to corporations' back-end Internet traffic growing at an accelerated rate. But bandwidth growth generally doesn't match the rate of data transfer growth—and when it does, bandwidth is still not a guarantee of throughput improvement.

White Paper
by Don MacVittie



WHITE PAPER

Myths of Bandwidth Optimization

Introduction

Moore's Law states that data density doubles approximately every 24 months, and Metcalfe's Law says that the value of a network grows in proportion to the square of the number of users. Because these postulates have held true in practice, global enterprises have found it advantageous to embed information technology into every aspect of their operations. However, this has led to increased bandwidth consumption—according to Plunkett Research, LTD, the worldwide data communications services industry now generates revenue in excess of \$3.1 trillion annually.

Despite a growing worldwide thirst for bandwidth, supply has outpaced demand by a wide margin. During the rapid expansion of the Internet in the 1990s, the data communications industry created an infrastructure that could deliver cheap bandwidth in high volumes. In fact, bandwidth has become so plentiful that even the effects of Metcalfe's Law are insufficient to consume available capacity for many years to come. The result of this imbalance has been the commoditization of bandwidth, rapidly declining bandwidth prices, and a vendor environment that actively promotes the myth that high bandwidth can address almost any performance problem.

Faster Application and Network Access

55% of surveyed customers selected an F5 solution over the competition because of faster application and network access for users.

Source: TechValidate Survey of BIG-IP users TVID: A0A-B7B-546

But as enterprise application deployments have expanded to the wide area network (WAN) and increasingly to the cloud, an environment where bandwidth is sometimes as plentiful as on the LAN, IT managers have noted a dramatic decrease in application performance. Many of them wonder, why would two networks with identical bandwidth capacities, the LAN and the WAN, deliver such different performance results?



WHITE PAPER

Myths of Bandwidth Optimization

The answer is that application performance is affected by many factors, associated with both network and application logic, that must be addressed to achieve satisfactory application performance. At the network level, application performance is limited by high latency (the effect of physical distance and physical communications medium), jitter, packet loss, and congestion. At the application level, performance is further limited by the natural behavior of application protocols (especially when faced with latency, jitter, packet loss, and congestion at the network level); application protocols that engage in excessive handshaking across network links; and serialization of the applications themselves.

Common Application Performance Myths

Myth #1: Application Performance Depends Only on Bandwidth

Application performance and throughput are influenced by many factors. Latency and packet loss have a profound effect on application performance. Little's Law, the seminal description of queuing theory and an equation that models the effects of physical distance (latency) and packet loss, illustrates the impact of these two factors on application performance.

```
In an application to networking, this law states:  
. (throughput) = N (number of outstanding requests)  
  T (response time)  
In terms of IP-based protocols, this translates to:  
TCP throughput = congestion window size  
  round trip time
```

Therefore, as the round trip time (RTT) of each request increases, the congestion window must increase or TCP throughput will decrease. Unfortunately, TCP does not effectively manage large windows. As a result, even small amounts of latency and packet loss can quickly cause network performance to drop for a given application to a fraction of what would be expected. Even if bandwidth capacity were to be increased to 100 Mbps, the application would never consume more than 1 percent of the total capacity. Under these conditions, managers who add network capacity waste money on a resource that cannot be consumed. The ROI for adding network capacity that will sit idle is non-existent; more than just capacity is included when determining performance.

"The Macroscopic Behavior of the TCP Congestion Avoidance Algorithm"¹ provides a short and useful formula for the upper bound on the transfer rate:



WHITE PAPER

Myths of Bandwidth Optimization

$$\text{Rate} = (\text{MSS}/\text{RTT}) * (1 / \sqrt{p})$$

Where:

Rate is the TCP transfer or rate throughput

MSS is the maximum segment size (fixed for each Internet path, typically 1460 bytes)

RTT is the round trip time (as measured by TCP)

p is the packet loss rate

Figure 1 illustrates this point:

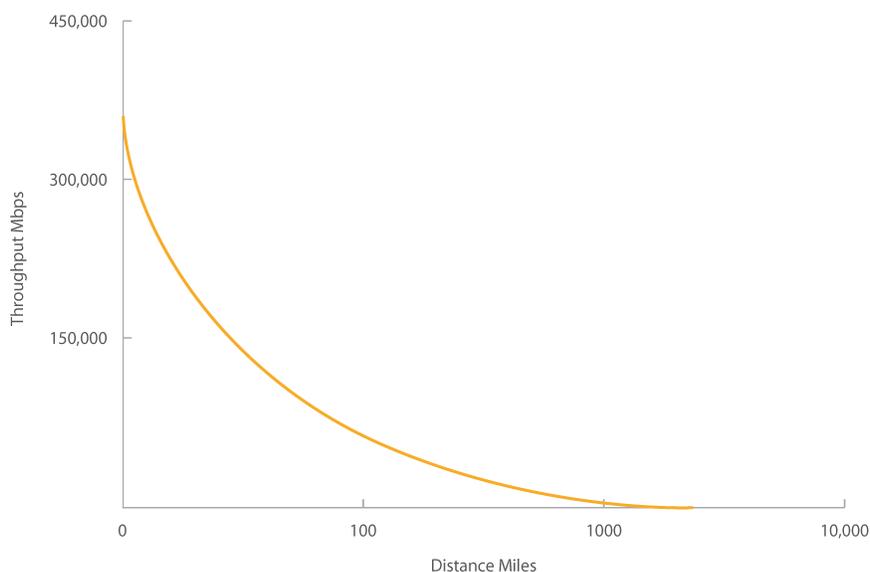


Figure 1: How TCP performance is affected by physical distance.

In WANs, sources of high round trip times (e.g., latency) include physical distance, inefficient network routing patterns, and network congestion—elements that are all present in abundance on the WAN.

Today, many TCP protocol stacks are highly inefficient when it comes to managing retransmissions. In fact, some stacks may have to retransmit the whole congestion window if a single packet is lost. They also tend to back off exponentially (i.e., reduce congestion windows and increase retransmission timers) in the face of network congestion—a behavior that is detected by TCP as packet loss. And while loss is often insignificant in frame relay networks (less than .01 percent on average), it is very significant in IP VPN networks that go into and out of certain markets like China, where loss rates commonly exceed 5 percent, and are often much higher. In the latter scenario, high loss rates can have a catastrophic effect on performance.



WHITE PAPER

Myths of Bandwidth Optimization

When packet loss and latency effects are combined, the performance drop-off is even more severe (see Figure 2).

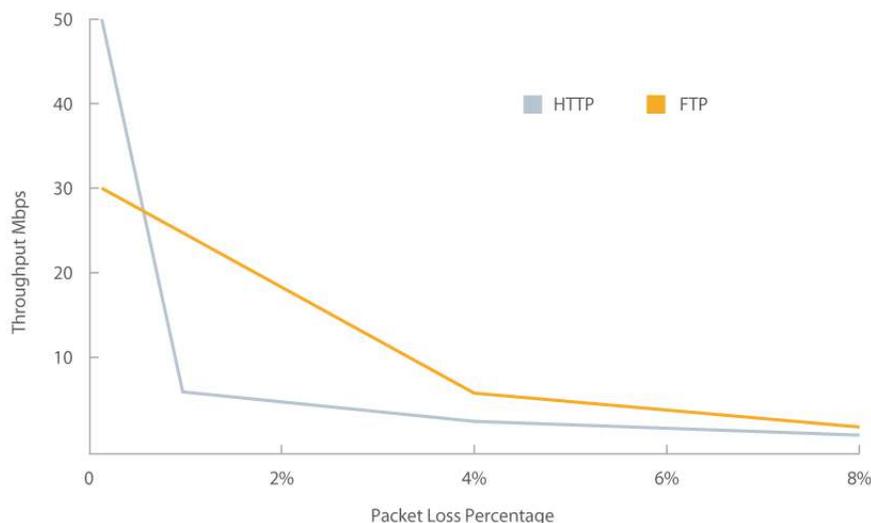


Figure 2: Sample TCP performance when packet loss is present.

Myth #2: TCP Requires Aggressive Back-Off to Ensure Fairness

Many network engineers believe that aggressive back-off in the face of congestion is necessary to keep network access fair. This is sometimes, but not always, true. Where congestion control is the responsibility of each host on a network—a network being an environment in which hosts have no knowledge of the other hosts' bandwidth needs—aggressive back-off is necessary to ensure fairness. However, if congestion is managed within the network infrastructure by a system that sees all traffic on a given WAN connection, then much greater and more efficient throughput is possible—and aggressive back-off is not required.

Standard protocol behavior specifies that when hosts consume bandwidth, they must do so independent of:

- The requirements of the application.
- The amount of available bandwidth.
- The amount of competition that exists for that bandwidth.

The result is that applications are often starved for bandwidth resources at the same time that the network is largely unused. This situation is obviously highly inefficient.



WHITE PAPER

Myths of Bandwidth Optimization

A much better solution to the TCP fairness problem is to allow individual hosts to consume as much bandwidth as they need, so long as all other hosts receive adequate service when they need it. This can be accomplished by implementing a single congestion window, shared by all hosts, that is managed within the network itself. The result is a system in which hosts get the bandwidth they need in periods of light competition, as well as when competition is more intense.

This single window method delivers consistently higher utilization and greater overall throughput than aggressive back-off. Hosts each see a clean, fast network that never loses packets (and therefore doesn't diminish TCP performance—see Myth #1), and cumulative traffic demands are matched to the overall buffering capability of the network. As a result, IT managers experience optimally utilized networks, under the broadest range of network latency and loss conditions.

Single window solutions can be constructed so that they are completely transparent to client systems. Components of such solutions may include TCP technologies such as selective acknowledgement, local congestion window management, improved retransmission algorithms, and packet dispersion. These capabilities are then combined with other technologies that match the throughput requirements of applications to the availability of network resources, and that track the bandwidth requirements of all hosts utilizing the network. By aggregating the throughput of multiple, parallel WAN links, single window technology can achieve even greater throughput and reliability.

Myth #3: Packet Compression Improves Application Performance

While common packet compression techniques can reduce the amount of traffic on the WAN, they tend to add latency to application transactions and therefore can impede application performance. These techniques require that packets be queued up, compressed, transmitted, decompressed on the receiver, and then retransmitted—all of which can consume substantial resources and add substantial latency, actually slowing down the very applications that need acceleration.

Next-generation application performance solutions combine protocol streamlining with transparent data reduction techniques. Compared to packet-based solutions, next-generation solutions dramatically reduce the amount of data that needs to be transmitted, eliminate latency that is introduced by protocol behavior in the face of physical distance, and can drive WAN performance at gigabit speeds. Transparent data reduction techniques often include multiple dictionaries where the level 1 dictionary is small and highly effective at reducing smaller patterns in data, and the level 2 dictionary is a multi-gigabyte space that can be used to reduce much larger patterns.



Myth #4: Quality of Service Technology Accelerates Applications

Quality of Service (QoS), if used properly, is a highly beneficial set of technologies that can be helpful for managing application performance. However, the only thing that QoS can do is divide existing bandwidth into multiple virtual channels. QoS does nothing to move more data or streamline protocol behavior. QoS simply decides, in an intelligent way, which packets to drop. And while it is better to drop packets in a controlled way than to leave it to chance, dropping packets does not accelerate applications.

Many QoS implementations rely on port numbers to track applications. Because applications often negotiate port assignments dynamically, these mechanisms must be configured to reserve a large port range to ensure coverage of the ports the application actually uses.

For QoS to be most effective, it should be dynamic. First-generation QoS implementations reduce large links into multiple smaller links, statically reserving bandwidth whether it is needed or not. "Channelizing" a network this way can ensure bandwidth availability for critical applications like voice, but actually wastes bandwidth as it is reserved for the specific application, even when the application is not in use.

Dynamic QoS solutions, on the other hand, ensure that bandwidth is reserved only when applications can use it. One common use of this technology is to extend enterprise backup windows by enabling continuous data backup when bandwidth becomes available.

F5 Brings It All Together

F5's WAN optimization solutions deliver dramatic replication performance and greatly reduced WAN costs. F5 provides these benefits by monitoring the limiting effects of network conditions, adjusting protocol behavior, and by managing all levels of the protocol stack, from the network layer through to the application layer.



WHITE PAPER

Myths of Bandwidth Optimization

Specifically, F5 BIG-IP WAN Optimization Manager (WOM) integrates advanced transport acceleration technologies such as adaptive TCP acceleration, Symmetric Adaptive Compression, Symmetric Data Deduplication, and session-aware rate shaping with best-of-class application acceleration technologies including SSL encryption offloading and termination. BIG-IP Local Traffic Manager (LTM)—and by extension, BIG-IP WOM—is supported by a statistics generation and monitoring engine that enables organizations to manage application network behavior in real time. Any communication to a remote BIG-IP WOM device is protected from prying eyes with IPsec tunnels, keeping corporate data safe even on the Internet.



Figure 3: Symmetric BIG-IP WOM devices optimize and secure traffic over the Internet.

F5 delivers LAN-like replication performance over the WAN. BIG-IP WOM optimizes replication solutions such as Oracle RMAN, Oracle Streams, Oracle Data Guard, and Oracle GoldenGate for database replication; Microsoft Exchange DAG for mailbox replication; and VMware vMotion over long distances for VM migration, file transfer, and other applications—all resulting in predictable, fast performance for all WAN users.

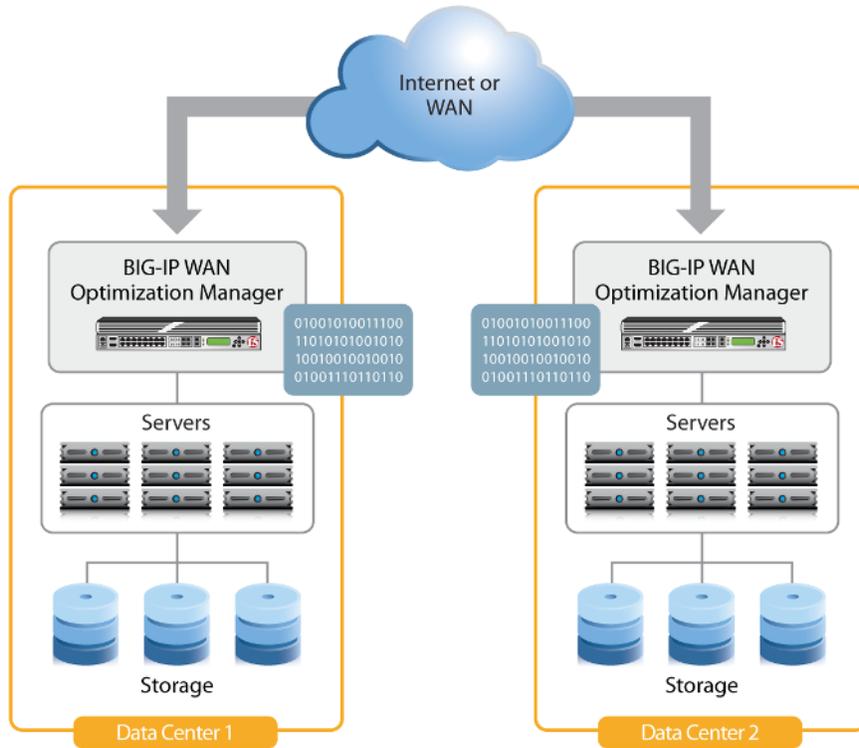


Figure 4: Accelerate all application traffic on the WAN.

F5 WAN optimization solutions are deployed on F5 hardware, which features fault tolerance, massive scalability, and unparalleled performance, or on virtual machines. For branch office deployments, BIG-IP Edge Gateway features functionality from BIG-IP WOM, BIG-IP WebAccelerator, and BIG-IP Access Policy Manager (APM).

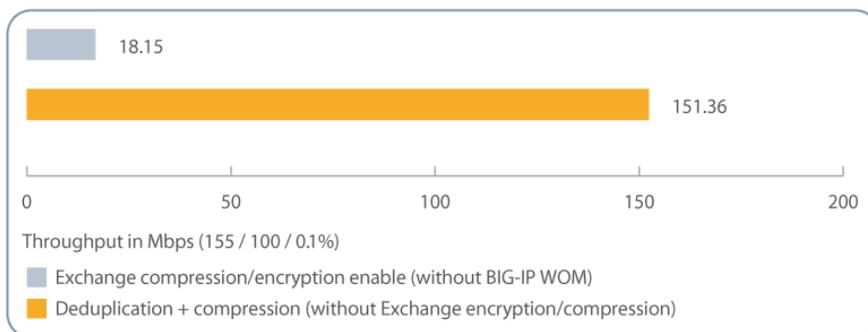


Figure 5: BIG-IP WOM improved throughput of Microsoft Exchange replication more than eight times more than Exchange compression and encryption.



WHITE PAPER

Myths of Bandwidth Optimization

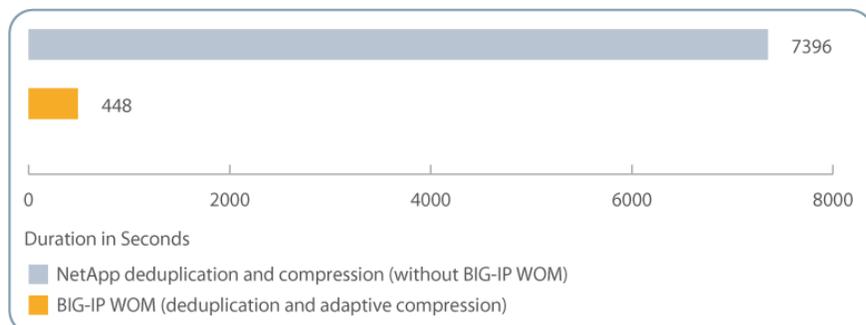


Figure 6: NetApp SnapMirror performance with BIG-IP WOM optimizations shows a 16x performance boost over NetApp deduplication and compression.

TCP performance with Symmetric Data Deduplication in BIG-IP WOM

When there is a lot of highly repetitive data flowing over the WAN, and that data is large enough to be caught by the BIG-IP WOM deduplication engine (a configuration setting controls the definition of "large enough"), then performance is massively improved. This is useful when someone is saving a lot of documents that contain the corporate logo or other commonly used set of bits.

The key part of deduplication is that while compression works on a single stream, deduplication can cross all of the streams being transferred across the WAN. This is important because duplication is often minimal within stream A, but when all of the streams being utilized across the WAN are considered, duplication rates are much higher. The amount of data reduction and the performance implications of that data reduction are heavily dependent on a given environment (the amount of duplication) and the configuration of the BIG-IP WOM device. The larger the number of bytes that must match, the fewer duplications administrators will receive; but the fewer entries must be checked for a duplication match, and more important, the longer a given set of bits will be saved for comparison. The smaller the number of bytes that must match, the faster the cache will fill up and older entries will be removed, but in theory the more matches administrators will find.



WHITE PAPER

Myths of Bandwidth Optimization

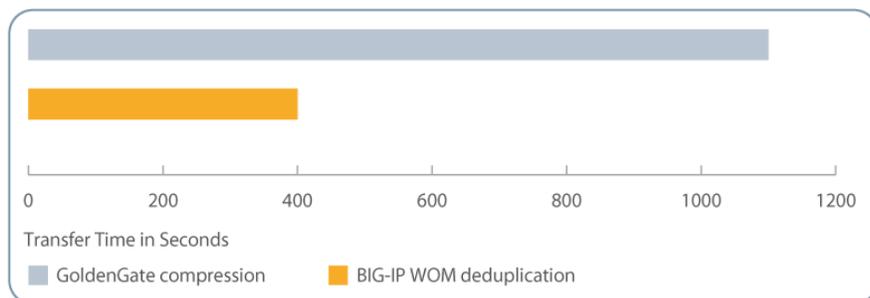


Figure 7: Oracle GoldenGate compression versus BIG-IP WOM deduplication functions show that with some datasets, deduplication can massively improve throughput.

Conclusion

Application performance on the WAN is affected by a large number of factors in addition to bandwidth. The notion that bandwidth solves all, or even most, application performance problems is, simply put, a myth. At the network level, application performance is limited by high latency, jitter, packet loss, and congestion. At the application level, performance is likewise limited by factors such as the natural behavior of application protocols that were not designed for WAN conditions; application protocols that engage in excessive handshaking; and the serialization of the applications themselves.

BIG-IP WOM recognizes the critical interdependence between application-level and transport-level behavior. It delivers predictable replication performance and increased throughput ranging from 3x to over 60x on networks as diverse as premium-quality, class-of-service managed networks to commodity, best efforts-based IP VPNs. The architectural advantages of F5 products result in replication and backup solutions that deliver best-of-class performance, massive scalability, and a return on investment that can be measured in months.

¹ Mahdavi, Jamshid; Mathis, Matthew; Ott, Teunis; and Semke, Jeffery. "The Macroscopic Behavior of the TCP Congestion Avoidance Algorithm." *Computer Communication Review*, a publication of ACM SIGCOMM, volume 27, number 3, July 1997. ISSN # 0146-4833. Accessed on November 14, 2011.

WHITE PAPER

Myths of Bandwidth Optimization



F5 Networks, Inc.
401 Elliott Avenue West, Seattle, WA 98119
888-882-4447 f5.com

Americas
info@f5.com

Asia-Pacific
apacinfo@f5.com

Europe/Middle-East/Africa
emeainfo@f5.com

Japan
f5j-info@f5.com

©2016 F5 Networks, Inc. All rights reserved. F5, F5 Networks, and the F5 logo are trademarks of F5 Networks, Inc. in the U.S. and in certain other countries. Other F5 trademarks are identified at f5.com. Any other products, services, or company names referenced herein may be trademarks of their respective owners with no endorsement or affiliation, express or implied, claimed by F5. 0113