# Traffic Management

*OR*

It's about the Application
not the bits!

**APM**
**ADVISORS**

August 2004

## Introduction

APM Advisors (APMA) practice is focused on the Application Performance Management (APM) marketplace. This paper was prepared as an update to the original Traffic Management paper published in March 2004. It is also an important extension to the Web Network Architecture paper published in February 2004.

There is a lot of action in what is commonly referred to as the traffic management space and therefore a need for the updates. We also continue to hear from more vendors, compare notes with enterprise architects and evolve our thinking accordingly.

## Preamble

> "A preamble defines a specific series of transmission pulses that is understood by communicating systems to mean 'someone is about to transmit data'. This ensures that systems receiving the information correctly interpret when the data transmission starts."

Overall we have tried to figure out why there is such enthusiasm around what we'll start calling Network Traffic Management (NTM), which is the compression, TCP optimizing, data caching, of network traffic. With some of the momentum some of these companies have experienced we were beginning to question our skepticism but then the fog began to clear. For the curious we have and continue to publish much of our papers independently rather than vendor sponsored. (Any exception to this will be clearly stated, but never will anything published by APMA be anything but unbiased.) This paper was not sponsored in any manner and represents our position as it has developed over time. In some cases our position and a set of solutions are quite closely aligned, which tends to be viewed as a bias. However, it would be our position that 'great minds think alike'!

As our practice has become further established and validated we have also engaged a number of architects within various IT organizations. This interaction has taken on a more formal structure, known as the APMA Infrastructure. Any architects, engineers or management within IT organizations are welcome to get involved by simply sending a note to apmainfrastructure@apmadvisors.com.

For those who have read our papers, we try hard to stay focused on the architectural issues of solutions and therefore what constitutes infrastructure. To become infrastructure, it must provide or contribute to the solution suite, which includes *information, resolution and control*. Additionally, the solution itself needs to be transparent or conform to the operational support standards common within most IT systems.

With that as a foundation NTM solutions are tactical and our recommendation is that they only be used sparingly rather than considering them as part of infrastructure. The following outlines the foundation for our position and as always we invite comment.

**Web Network Architecture (WNA)**

A few months ago we published the first draft of the WNA as a first cut at an architectural template for making APM a part of infrastructure within any system. The foundation of the WNA is that application architectures have a much more significant impact on the design of the underlying delivery infrastructure.

Therefore the infrastructure must have 'application continuity' rather than be 'application aware'. Application continuity is a concept that begins at the origin of the application (data center, etc), but is continued logically throughout the system. It's our position that between SSL and the dynamics of application services above Port 80, that it's impractical to rely on products or an architecture that need to perform 'deep packet inspection' to develop application awareness in the network.

Furthermore, any system architecture needs not only address what we consider APM *control* but also develop meaningful *information* and provide the means to *resolve* problems. As solutions like NTM become prevalent within the infrastructure the fundamentals that operations groups depend upon are disrupted.

So currently the bandwidth demand problem and the immediate relief from pain is driving the acceptance of NTM, but over the long term this short-term gratification will at best, be the easy way out. It's easy to scratch an itch!


**Network Traffic Management**

While 'traffic management' can and is used to describe a variety of solutions we add the 'network' descriptor to put more focus on a class of products. At a high level it's a product solution that is focused on managing traffic in alignment with the network. Therefore a box is put at a location to compress and manage traffic flows to improve application performance due to network constraints.

The excitement around these types of products is nothing new, since there have been proprietary compression products since there have been modems. Turning a bit pattern into a 'token' significantly reduces the number of bits that have to go on the network and therefore is instant gratification. One day there is a problem and the next there isn't. These types of solutions have a very clear value where multi-national companies have interconnected locations using very expensive bandwidth.

However, when looking at the infrastructure the enthusiasm recognized in a site-to-site solution isn't easily transferred to a large distributed system. Now before we go much further, we will acknowledge that the technology and the maturity of the vendor NTM solutions <u>can</u> work in a larger distributed environment. The arguments for not doing it though are based on architecture, scaling and operations issues that follow.

**WNA Infrastructure**

It's critical to approach infrastructure in a strategic versus tactical sense to assure the critical capital and manpower resources are aligned with the long-term architectural goals. Therefore in most distributed application environments the supporting infrastructure needs to align with the applications rather than the network.

While there is a significant base of non-browser based services the dollars being invested today are web focused. Quite often it's new demand on bandwidth generated by these applications that is taxing existing networking infrastructures and in many cases right along side of it is the move to consolidate servers. The long awaited adoption of VoIP is also an application service that can't be ignored while planning the infrastructure.

As architects there is also the stark reality of providing operational support to the new infrastructure, which is often looked at as a resource management issue, but the more important aspect is how well the new solution maintains transparency. The operational infrastructure is associated with monitoring and resolving problems efficiently, which in a distributed system is no small feat.

NTM Core Features

When looking at the features of an NTM, there are a lot of core technologies that deliver a significant amount of value in a network. With a focus on compression and the extensibility of it as a core technology into what is now known as 'data cache', these products can dramatically improve network efficiency.

*Compression*

As outlined earlier one of the key benefits of a NTM is compression of traffic, which is getting broader in definition everyday. The fundamentals of compression are pretty straightforward so we don't get too excited about who's product is x% better than the other, because it's still a question of architecture.

As the boxes identify data blocks, they assign 'tokens' which two or more boxes store in a reference 'dictionary'. The more repetitive the data between the two boxes, the more efficient the compression.

No matter how or what is said, that's the fundamentals of compression. To layer more value on the compression services, vendors are beginning to throw more memory (hard drives) into their products to take what was compression to a new level. With the additional memory these products can maintain more and larger dictionaries that can be used to limit the repetitive transmission of complete files.

While this obviously reduces the amount of data sent between sites for repetitive traffic it has limited value where the content between locations is more dynamic. The other consideration is that these devices are 'black holes' from an operations perspective. Other than having an IP address and the management services provided by the vendor they have no logical presence.

*TCP Optimization*

Furthermore depending upon the specific implementation the device may be terminating TCP on each end of the session. For networks with high latency or the mix of various traffic services the manipulation of TCP can achieve some very real value. However, the termination, emulation or restructuring of the TCP connection between a user and the target application system is something to carefully consider from an operational perspective. Having a clear understanding of the TCP services from end-to-end provides a significant amount of insight into performance, error conditions, latency, etc. Once that is compromised, either the product will need to integrate into established operations or vice versa.

*Network Optimization*

The orientation of these boxes toward managing traffic to align with the network as opposed to optimizing the application services to work efficiently over a distributed network is a concept important to understand. More often than not, the NTM is purchased for a specific bandwidth and therefore the hardware has a life span equal to the amount of bandwidth servicing the location. In fact some of the vendors have developed 're-deployment' schemes where a device used in a low bandwidth location is replaced and the original device is placed in the data center as part of a logical head-end.

*QoS or Prioritization*

Regardless of the bandwidth there is bandwidth contention that needs to be dealt with. Whether an email attachment or a VoIP session, there is a clear requirement to manage traffic flow. Another 'personality trait' of an NTM is the reliance in deep packet inspection to determine traffic type and then provide Layer 3/4 flow control. As we have pointed out repeatedly, the adoption SSL and granularity of application services over Port 80 make this a very challenging approach. Additionally this becomes a product specific maintenance issue, since the devices must be updated regularly to keep pace with application developments.

**WNA Alternative to NTM**

As outlined earlier, we do feel the NTM products have a role in certain network environments due to their clear ROI for various network services and capacity to improve performance over network facilities with high latency. However to use these products as part of infrastructure for a large distributed network we consider tactical rather than strategic.

From an architectural perspective as well as the foundations in standards, the methods to improve application performance and network efficiencies are fundamentally there. With that as a foundation there is a scalable and operationally transparent method to extend application services.


Application Traffic Management (ApplTM)

Fundamentally, ApplTM is based on establishing application continuity across the infrastructure from a logical perspective, that still delivers efficiencies and improves performance. To be clear we do, have and will support the value of compression, but it should be approached as an extension of an application rather than as the network. The ApplTM approach assures scalability and operational continuity that can't be achieved with NTM oriented solutions.

*Policy*

To have continuity across a system there has to be a common definition structure and the means to distribute rules across all infrastructure resources. Amongst some of the themes about the future of performance management that we've been vocal on, policy is very high on the list. The fundamental issue is that a system solution will require multiple vendor products and their perspective and configuration will need to be coordinated. (*APM Advisors is launching the Business Policy Initiative to put some focus on this issue with a number of innovative companies. To get involved* [info@apmadvisors.com](mailto:info@apmadvisors.com) *)*
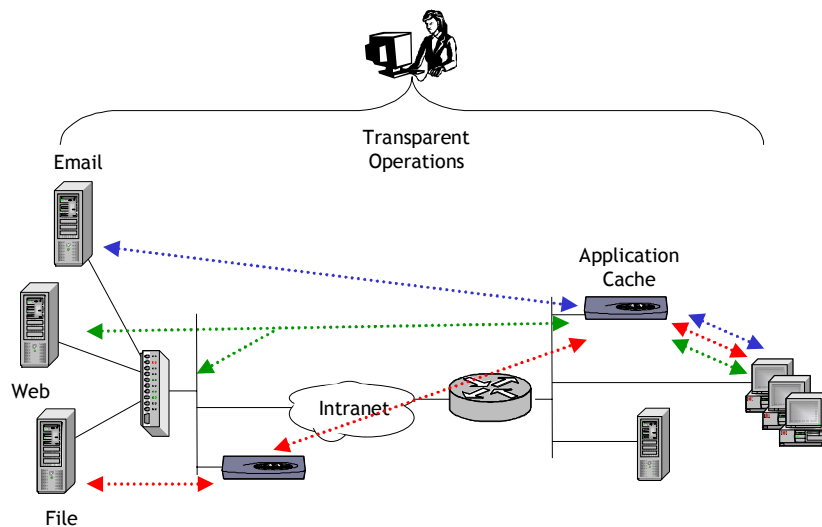
*HTTP Compression*

As outlined earlier, compression is using tokens to represent blocks of data to minimize redundant traffic on a network. Between a web server and browser, HTTP 1.1 included the ability to compress 'x'ML traffic, which is highly efficient since this traffic is predominately text. Embedded media such as images, video and audio is pre-compressed and therefore packaged for transmission.

For other embedded MIME traffic types the fundamentals of Gzip or Deflate can easily be applied between the ends with a little coordination. When a browser client links to a MIME file the specific function being performed on the client may bypass the decompression feature in the browser. Therefore with a little session coordination both ends could assure a successful compression delivery. (*Ahh yes, the session. Something that is basically a concept at this point, but there is always hope that it'll make a comeback.*)

*Application Cache*

Outside of browser driven bandwidth demand, the consolidation of servers is also placing a considerable demand on the networking infrastructure.  While the value of distributing servers to support data demand file and e-mail services was logical, but the capital and operational support issues are driving it toward extinction.

Taking an approach similar to the caching of web content, there are application oriented caching solutions coming to market that maintain the logical structure of the distributed application services, without the overhead (license and administration).  It's the logical or application orientation of these solutions that provide strong performance results and bring with it the operational transparency.



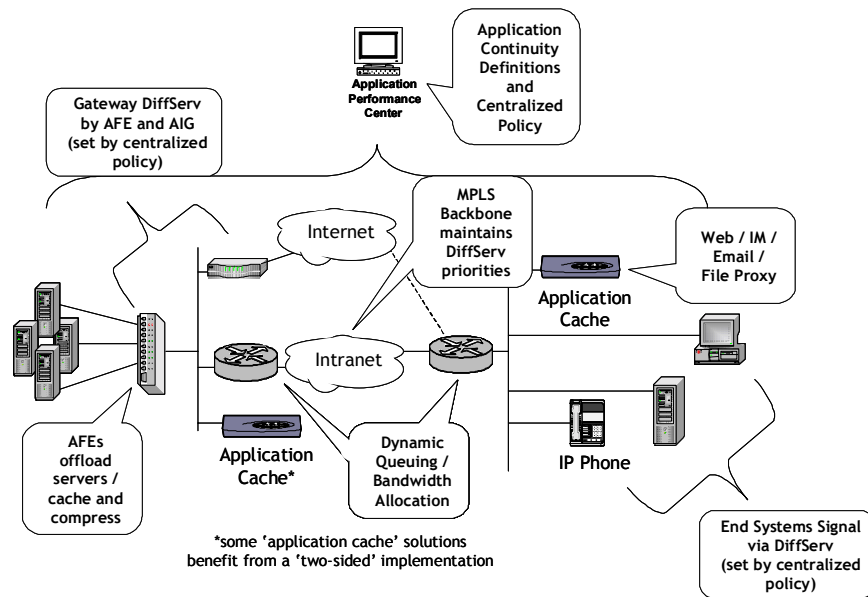*Application cache localizes remote application services*

While today these products are delivered as separate appliances from vendors, the consolidation of functionality will occur rapidly.  By looking like web servers, network drives and mail servers; these solutions will become a common element of infrastructure.

*QoS or Prioritization*

While all prioritization requires a network orientation and some packet analysis, the foundations built within today's end systems, edge routers and backbones are in place. Virtually all end systems can generate DiffServ settings, which are recognized by edge routers and passed into the core network.  The inhibitor as of today has been a definition of an application, the establishment and distribution of a common policy to end systems (host). *(Since the socket API allows the application to set the bits, there is a requirement to instrument the systems with a 'policy enforcer'.)*

## A Working WNA Model

As outlined earlier, all of the components needed to deliver on the WNA model are available, but the consolidation of functions and coordinated policy will develop over time.  Therefore in the review of the model the diagram is conceptual in respect to product, but representative of a viable solution.



*WNA is a scalable system architecture based on central control with distributed intelligence*

Based on the strength of vendors, trends in application architectures, economics and common sense the foundations of the WNA model outlined in the diagram will prove to be the model for large distributed networking services.  You heard it here first!

While on the path for a few years, some of the challenges encountered by architects, engineers and the vendors who cater to them, the development of the WNA model moving from 'foggy' to 'clear'.  We have collaborated with a number of enterprise architects who share a similar model for the future of their infrastructure.  Obviously there are exceptions to every rule and IT systems that don't fit this model, so again for clarity this simplified model aligns with large distributed networks.
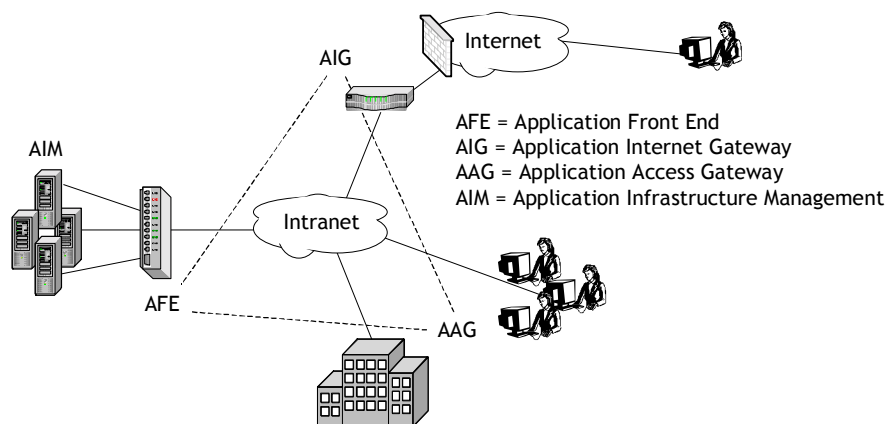
Fundamentals

A system has basically four 'service points' where intelligence can be reliably applied to meet the dynamics of demand and scale. As outlined in the following diagram, these product solution points are:

AIM – Application Infrastructure Management is the foundation for assuring that applications are available and are operating as efficiently as possible. At this point this is a very specialized area within the system due to the complexities of the n-tier application architecture. (see AIM - New Ways for New Times - http://www.apmadvisors.com )

AFE – Application Front Ends provide management of the user-to-application transaction by offloading many of the redundant protocol services from the application infrastructure and assuring that user connections are efficiently managed. From a data center 'outbound' perspective these devices encrypt traffic (SSL) and apply transaction appropriate compression. From an 'inbound' perspective these devices also provide a first / last line of defense against various intrusions that threaten the stability of the application infrastructure. (Watch for the AFE 2nd Edition scheduled for release September 8th 2004)

AIG – Application Internet Gateway is a critical service point to integrate into the architecture to assure application continuity within the intranet. From an APM perspective the primary service is to insulate the intranet from unwanted traffic, which after basic firewall services includes the enforcement of usage policies. The breadth of services required here is expanding rapidly and a significant amount of feature consolidation will continue.

AAG – Application Access Gateway is a logical extension of the core services provided by AFE's and AIG's to remote users, whether within or outside of the intranet. The goal in the delivery of those services is to optimize application performance while controlling cost and maintaining operational integrity.



AFE = Application Front End
AIG = Application Internet Gateway
AAG = Application Access Gateway
AIM = Application Infrastructure Management

*WNA 'service points' in a distributed system*

Architectural Foundation

Many of the challenges of obtaining the core services of Application Performance Management (information, resolution and control) are associated with the decoupled nature of the protocol stack. Therefore architects are now dealing with applying technology solutions that manage the relationships between the various layers of the stack.

Beyond that core problem, the good news is that a significant amount of the underlying application deliver infrastructure is established. This is enabling the shift from point products to system solutions for managing performance, which will be good for all involved.

*Deep Packet Inspection*

As we have been outlining it's going to be the application architecture that will drive the network or application delivery architecture. The dynamics of applications are just beginning to come to market with the breadth of application services that can be delivered via web technologies. This dynamic is developing a 'network on Port 80' and in many cases applications that defy recognition, either because of proprietary implementation or encryption (SSL).

Even if vendors are able to keep up with 'deep packet inspection' within the network there are some other issues that challenge the future of this approach as a strategy. Because application profiling is an ongoing process the vendors of these solutions need to continuously update their products, which becomes a maintenance issue for their customers. While it can be streamlined from a process perspective it is a cost of ownership issue.

*Application Cache*

The extension of application content to the user is one of the most powerful technologies to improve response times in a distributed system. The foundation of the WNA ApplTM approach is that the instrumentation of this 'data cache' needs to align with the application service it is extending. The 'application cache' needs to have a logical presence that emulates that service. This improves operational support by not forcing operations to rely on the management provided by the vendor.

Furthermore application cache services can provide a richer set of application-aligned services. For example an application cache that is extending file services can provide drive mapping services to map multiple remote drives to one local drive in a remote office. An email application cache pulls down only one attachment to the remote office even though all users are intended recipients.

*DiffServ and IP Queuing*

Other than the architectural considerations we've applied in the development of this model we have found that nearly all architects are or will use DiffServ signaling to segment traffic. This provides enough differentiation between traffic types, but there is also a clear understanding that more granularity will be required.

In distributed networks there is still lots of lower speed Frame Relay but for most IP/VPN – MPLS is around the corner. The bandwidth supports the payload / transaction ratio to maintain reasonable response times but contention is everyone's reality. Having the mechanism to prioritize traffic during contention is a requirement. As the same complexities of the application architecture become more commonplace in delivering application services to remote offices the end-systems need to get involved. While the core elements are built into the infrastructure (API and routers), there are only a couple of product solutions available to set common policy and set the bits, but that's short-term problem.

Obviously there will also be a requirement for the AFE, AIG and Application Cache to get involved in the support for DiffServ and adherence to a common policy. Here again some more validation for a common policy structure and definition of an application.
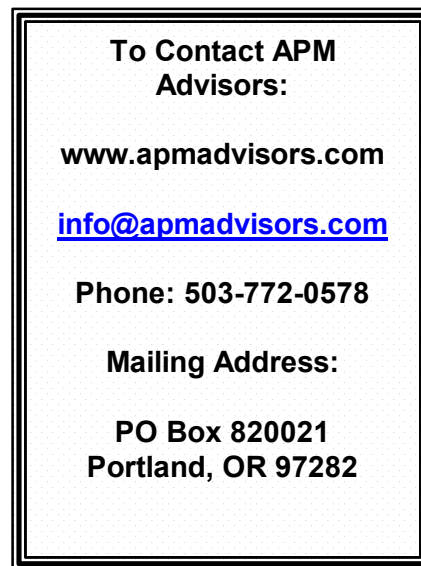
Carrying through the concept of continuity, there is no clearer example of how concept applies to reality. At the highest levels, an organization policy is established which is distributed across a system, and compliance is built into the infrastructure. As engineering and operations gather information and support the infrastructure, it doesn't take vendor specific management tools to verify compliance.

## Closing

The theme for well-architected solutions is maintaining application continuity, which requires a solution to have a logical application presence. The delivery of application services based on the WNA Application Traffic Management (ApplTM) provides highly scalable performance and aligns with established operational methodology.

While Network Traffic Management (NTM) solutions have their role in networks, they are not infrastructure for the majority of large distributed networks. The challenges with deep packet inspection in the network and lack of operational visibility will keep them relegated to limited tasks.

It has and always will be about the applications in IT, but for the first time the application architecture is truly driving change. As architects it's important to be on the front-end of that curve and enjoy the longevity of your efforts.

---

**To Contact APM Advisors:**

**www.apmadvisors.com**

**info@apmadvisors.com**

**Phone: 503-772-0578**

**Mailing Address:**

**PO Box 820021
Portland, OR 97282**

---

**Legal**

*We don't hire lawyers if we can avoid it, so let's keep it simple.*

*This document contains material, which APM Advisors has sole, and exclusive control over. Any use of this material without the written permission of APM Advisors is just plain wrong. It is only available via download from APM Advisors and any partner site solely selected by APM Advisors. No other distribution rights are implied. Any information contained in the document was developed with as much care as possible and sometimes even we make mistakes, so just let us know and we'll correct it.*