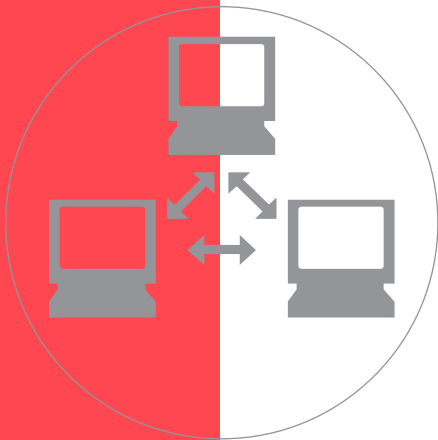# T R A N S I T I O N
### n e t w o r k s

**The Conversion Technology Experts**

# Quality of Service (QoS) in High-Priority Applications

**Abstract**

It is apparent that with the introduction of new technologies such as Voice over IP and digital video, network managers and administrators have a tough time keeping up with ever-increasing bandwidth requirements. Such technologies are brought with historically high expectations for reliability and quality. Today's networks must treat these services as high priority. These traditionally "best effort" Local Area Network protocols (Ethernet etc.) face a difficult time handling these High Priority requirements. Quality of Service (QoS) promises better handling of these new challenges; increasing reliability and quality.

Network administrators have two major types of QoS techniques available. They can attempt to negotiate, reserve and hard-set capacity for certain types of service (hard QoS), or just prioritize data without reserving any "capacity setting" (soft QoS). This paper will discuss both hard and soft QoS techniques including 802.1P, IP Precedence, Differentiated Services, Resource Reservation Protocol (RSVP) and ATM specific priority resources. The paper will also explain how to implement QoS features on Transition Networks' Management Aggregation Converter.

# QoS 101

QoS stands for Quality of Service. In QoS the bandwidth, error rates and latency can be monitored, sampled and possibly improved. QoS also delivers the set of tools to help deliver data efficiently by reducing the impact of delay during peak times when networks are approaching full capacity. QoS does not add capacity; nor does it multiplex the signals like WDM. It simply tries to manage data traffic better so that top priority traffic will not be compromised. QoS helps manage the use of bandwidth by applying a set of tools like priority scheme, so certain packets (mission critical – must go packets) will be forwarded first.

### QoS vs. Class of Service (CoS)

QoS is often used in conjunction with Class of Service. The shortest definition of CoS would be "a grouping". CoS defines groups of traffic with a specific type of service, QoS manages this type of service and assures that it is delivered. Similar types of data such as Voice, Live Video, or streaming video and large file transfer can be grouped together in a service class and treated with a same level of service priority.

### The need for QoS

Many users believe that more bandwidth will resolve the problem. Throwing more bandwidth though may not work anymore. Voice over IP Telephony and other new technologies such as a networked video security, remote monitoring, and recording over IP networks are becoming more popular. They have begun to penetrate traditionally data orientated networks, forcing network administrators and managers to employ measures such as QoS to accommodate these technologies efficiently and without any backslash to the performance of an existing network.

Multiservice traffic is difficult to handle efficiently because each type of traffic requires different transfer rate and and each has a different tolerance for delay or packet sequencing. The original best effort LAN protocols were designed for applications such as basic connectivity between stations, file transfer, e-mail, MRPs, and later on the Internet.
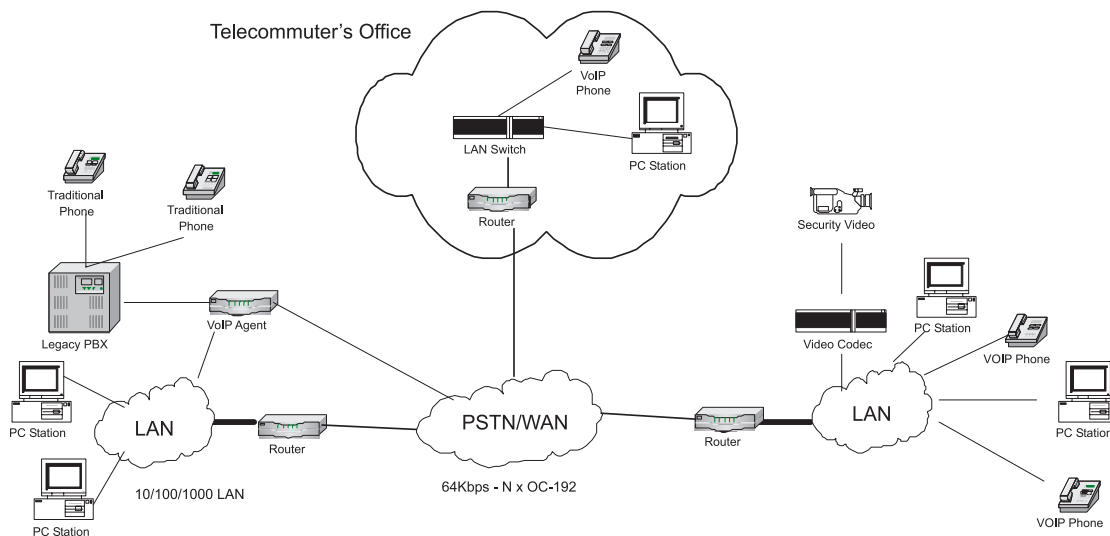


*Figure 1: VoIP Integration*

These applications are not compromised by packet delay so as long as connection was established and transfer of data happened in less then irritating manner, the network served its purpose. Also, multiservice traffic has to peacefully coexist with the infrastructure in place. In many instances VoIP has to be routed back to the PSTN (see Figure 1) in order to complete the necessary call or IP video has to be broadcast over the existing close circuit TV (CCTV).

The network bandwidth is still important, but it is no longer the only factor to consider for implementing future technologies. The new specific characteristics of this traffic (delay, jitter etc.) need to be read, understood, and implemented.

One of the keys in delivering voice or video over any media is the maintenance of a level of quality. The quality of voice or video may deteriorate as a result of three factors:

*Overcompression*
Compression ratios are inversely proportional to the quality of a voice signal that is transmitted over the network, and is inferior to what the user is accustomed with Plain Old Telephone (POTS). The lower the compression the higher the throughput necessary to transmit voice packets, increasing the possibility of network congestion and consequently the loss of quality. Compression can be easily controlled by users.

*Packet loss on the network*
Packets get lost on the network, which is not a problem for traditional applications. The quality of traditional applications such as file transfers is immune to packet loss because these losses are recognized by the network and retransmitted. VoIP products reconstruct the packets if the number is minimal. The rule of thumb is that no more than 10% of packets should be lost in VOIP networks otherwise the voice quality will be compromised.

*Latency*
Delay in data networks is not that critical. Waiting for a web page to load is not as irritating as a silence in your receiver when you are in the middle of an important conversation. A maximum delay of 150ms is the rule of thumb for one-way latency to achieve similar quality to POTS voice.

Network managers face a new challenge with voice and security applications. Traditional POTS is highly reliable in terms of transport and reliability. It is hard to imagine the situation when you have no dial tone in our phone even during the worst storm. While it was acceptable to wait 5 seconds loading a web page, it is impossible to tolerate such a delay during a conference call with the customer. It is impossible to accept a voice breaking off or any noticeable latency.

Such expectations are being brought to the "opportunistic – best effort" networks creating the need for QoS. First in First Out (FIFO) systems so commonly used in opportunistic networks have to be replaced by more sophisticated often dynamic resource allocation tools starting with 802.1P and all the way to RSVP.

One important condition to be met in order for the QoS to be successful is that it has to be employed and managed end to end, across several LANs and WANs (see Figure 2). This can guarantee all the bottlenecks are addressed and that voice/video will not be distorted. If QoS is employed only on the portion of the network, anything that has to go out of this network through the "bottleneck" will be treated and forwarded in the order it was received, at the available speed and with a possible delay.
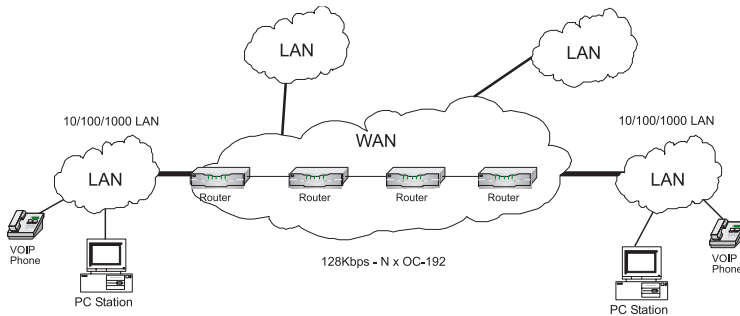


*Figure 2: End-to-End VoIP Application*

Protocols differ in their natural ability to properly handle high volumes of traffic and some offer traditionally higher "reliability". ATM is a very successful protocol in the multimedia applications because ATM can provide guaranteed rates and connections that are so valuable to voice or video transmission. ATM prioritizes the traffic by assigning it to one of four service classes. Each class can receive a priority level. There are the following four ATM service priorities/queues:

- Constant Bit Rate (CBR) absolute guarantee of a service level (VoIP or standard voice circuits).

- Variable Bit Rate (VBR) for variable burstable transmission rates with a very good throughput, but no guarantee as to the consistency over time (FTP, streaming).

- Available Bit Rate (ABR) offers a minimum guarantee.

- Unspecified Bit Rate (UBR) which makes no guarantees, whatever bandwidth is left can be used.

ATM can work along the priority settings done in the Ethernet LAN. ATM though, due to lower LAN penetration will not be able to solve all LAN QoS issues. This will have to be done by a "protocol of choice for LANs" - Ethernet.

Ethernet represents more of an opportunistic protocol. Ethernet is a connection less broadcast protocol and is "Best Effort - As soon as I can". Ethernet was designed to be less complex and hence less expensive. When data is transmitted Ethernet allocates the maximum possible bandwidth to this transfer until the network runs out of its bandwidth. Consequently the "critical traffic" is treated as any other transfer, so it pretty much drowns in the sea of less critical/significant data. This means that it will do just fine with voice and video in the time where there is no congestion.

**QoS in the Enterprise**

Networks are seldom designed for the worst case scenario (max overload) so QoS helps effectively manage what we have at our disposal without magically adding bandwidth.

Today's enterprise networks are experiencing increased complexity and packet equality becomes a song of the past. Below (Table 1) we define different types of traffic, their bandwidth requirements, and delay tolerance for each of them. The tolerance score will explain how tolerant users are towards each service when things do not go as smoothly as we would want them to go.

**Table 1: Performance Requirements by Media Type**

| Traffic | Bandwidth req. | Relative delay Tolerance | Time Factor |
|---|---|---|---|
| Internet browsing | 128 - 1mbps | High | Buffer tolerant |
| Data transfer (e-mail, fserv etc.) | 128 - 1mbps | High | Buffer tolerant |
| Fax | 28Kbps | High | Buffer tolerant |
| Customer Chat (text) | 28Kbps | Low | Low tolerance |
| Voice (phone) | 64Kbps | Extremely Low | REAL TIME |
| Voice (Teleconference) | 1Mbps | Low | REAL TIME |
| Voice & Video (MPEG-2) | 5Mbps | Medium / Low | REAL TIME |
| Video (security) | 256Kbps–1 Mbps | Low | Buffer tolerant |
| Video Streaming | 256Kbps–1 Mbps | Low | Buffer tolerant |

Our tolerance of delays dictates what kind of time factor to implement. Users' relatively high tolerance towards delay in such services as internet browsing, web hosting, data transfer or fax enables network designers to allow for transport buffering of such services. Real-time applications require establishing benchmarks for service.
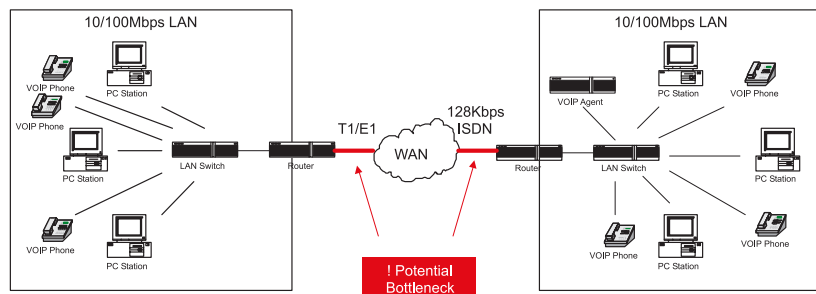


*Figure 3: Possible Bottlenecks in VoIP Implementation*

Clearly the bandwidth requirements mentioned above do not produce any hiccups within a small enterprise operation. 100Mbps Fast Ethernet can support these services, but what happens in bigger organizations? What about end-to-end, when data to go "outside" the enterprise LAN? Again, one of the key technical issues with QoS is that it must be supported end-to-end to be effective (see Figure 2). IP telephony and video conferencing, unlike basic Internet surfing, must have a minimum transfer rate guaranteed so that they can function properly.

A 100Mbps LAN connection cannot guarantee a voice connection with another LAN over 128Kbps WAN connection. Due to the nature and requirements of this communication, the connection has to be continuous and there is no room for voice buffering.

So in times of congestion - what can be done to ensure those critical pieces are flowing and the delay is minimal?

You need to first define what kind of traffic causes the bottleneck, and where the bottleneck is located. This may very well identify one of the following causes for congestion:

**Peak usage.** Too many packets are being sent over the network by users. A closer look at large traffic patterns can sometimes quickly identify the cause if it is deemed unnecessary. Necessary traffic can also cause congestion because the existing network cannot provide sufficient **switching and routing capability**. Segmenting the network can help.

If your network is still congested and you cannot throw more bandwidth at it, you can apply the set of the following tools:

- Prioritize Users

- Prioritize Segments

- Prioritize Applications

- Reserve/Limit the bandwidth for certain Users

- Reserve/Limit the bandwidth for certain Applications

- Select the applications that can be stopped

Having defined who and what gets priority, network administrators have a set of tools to implement these rules.. For instance, they can give priority to certain users based on their IP address (source address). Or they might prioritize by segment either through subnet mask or destination address. Prioritizing application means that all Voice over IP services get a higher priority than let's say e-mail.

Devices read the instructions as to the priority or bandwidth allocation and queue packets in the following four types of queues:

- Priority queuing from high priority queue to low priority. Packets are sent from the queues of higher priority first (as explained in the IEEE 802.1P).

- Weighted fair queuing. It allows for guaranteed bandwidth services, but over the same, shared link.

■ Class based Queuing divides user traffic into classes. These classes are assigned based on IP addresses, protocols and application types
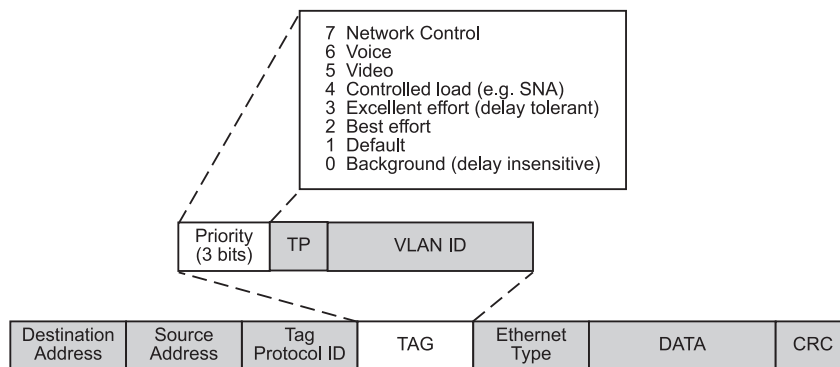


7 Network Control
6 Voice
5 Video
4 Controlled load (e.g. SNA)
3 Excellent effort (delay tolerant)
2 Best effort
1 Default
0 Background (delay insensitive)

| Priority (3 bits) | TP | VLAN ID |
|---|---|---|

| Destination Address | Source Address | Tag Protocol ID | TAG | Ethernet Type | DATA | CRC |
|---|---|---|---|---|---|---|

*Figure 4: IEEE® 802.1P*

### QoS IEEE® 802.1P Prioritization in LANs

The IEEE 802.1P is a signaling technique for prioritizing network traffic at the data-link/MAC sublayer (OSI Reference Model Layer 2). The 802.1P header includes a three-bit field for prioritization, which allows packets to be grouped into various traffic classes. The IEEE 802.1P compliant switches pick up on this tag (the packet contains a 32-bit tag header located after a destination and source address header), read it, and put the packet in the appropriate priority queue. No bandwidth is reserved nor requested by this technique.

There are eight levels (0-7) of priority and consequently eight queues that could be created (see Figure 4). Level Seven represents the highest priority. This will be assigned for mission-critical applications. Level 6 & 5 is designed for delay-sensitive applications such as interactive video and voice. Levels four and below, are suitable for regular enterprise data transfer, as well as streaming video. Level zero is assigned for a traffic that can tolerate all the drawbacks of a best-effort protocol.

The switch will analyze the packet based in the "P" tag and will place it in the appropriate priority Queue for sending. The user can have as many as eight priority queues. An adjustable algorithm is employed to choose how many packets are being sent from each queue before the packets in the lower priority queue are sent.

Transition Networks Management Aggregation Converter (MAC) is a converter that allows the remote end of the network to be managed. One of the MAC's many features includes supporting 802.1P packets. The MAC reads the 802.1P tag and places incoming packets in either a High Priority Queue or a Low Priority Queue.The network manager defines the priority level threshold (0-7) that determines if a packet is placed in the high priority queue or the low priority queue. . For example if the threshold is set to 4, a "P" tag of 5 will be forwarded to the High Priority Queue while a packet with a Tag of "3" will be placed in Low priority queue. The MAC converter also implements a user adjustable algorithm for packet queue selection. As shown in Figure 5, 15 packets from High

priority queue will be sent and then one packet from the Low priority queue will be sent before MAC converter comes back to high priority queue again.

The user can also set a high priority on a specific port (for IP phones, etc.) and this will automatically put all the packets on this port in a High Priority Queue. (see Figure 6)

In addition to queuing, MAC converters will also enable users to disable/enable Pause in higher priority applications so that real-time traffic (Voice) will not be paused in times of congestion. (see Figure 7)

All converter management can be performed by a fully SNMP compliant Graphical User interface (GUI) Software - Focal Point™ or it can also be managed via the web based management using any web browser.

QoS 802.1P is an efficient tool for prioritization within a LAN. QoS can also be accompanied by IP precedence or Differentiated Services - Layer 3 QoS mechanisms to achieve inter LAN prioritizing.

### IP Precedence
The IP protocol includes the Type of Service (ToS) an 8-bit field, intended for use in packet prioritization. It allocates three of the ToS bits to create up to 8 priority levels and three bits to describe delay sensitivity, as well as packet loss. Transition Networks' MAC converter is transparent to these packets.

### Differentiated Services
Another very popular method of QoS in the enterprise is Differentiated Services. It is an efficient method of managing traffic based on its class. Differentiated Services (Diffserv) prioritizes certain types of traffic like voice traffic over other types of communications. It works by categorizing IP packets into classes. The six bits in the type-of-service byte contained in the IP header of each packet, specifies a particular behavior type which determines the packet-forwarding scheme and priority.

Differentiated services can offer the following:

- Expedited Forwarding (EF), which defines minimum delay and jitter. Preferred mode for the VoIP.

- Assured Forwarding (AF), which introduces three selectable packet drop rates. During congestion, packets with a high drop precedence are discarded. Thus enabling the more important traffic marked with lower drop precedence to get through.

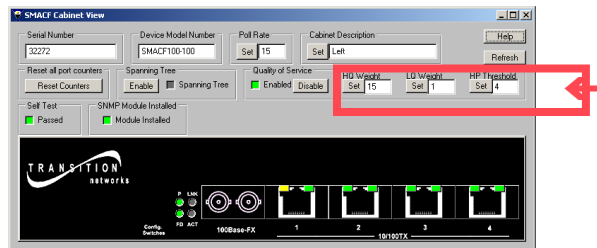- Best effort picks up the remains of the bandwidth not allocated to EF and AF.



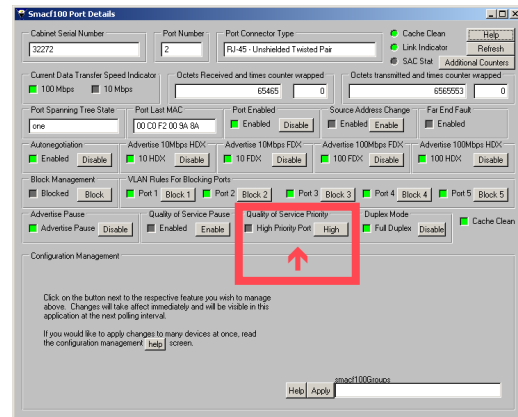*Figure 5: Adjustable algorythm for packet queuing*
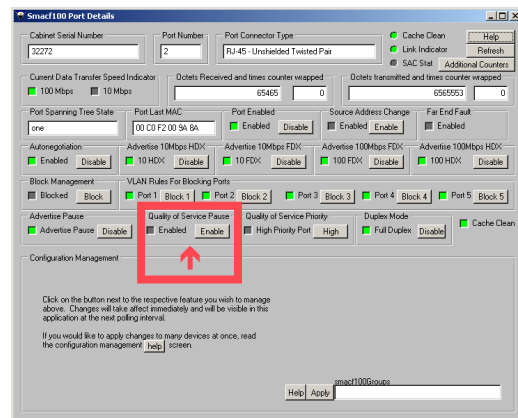


*Figure 6: Port Priority Setting*



*Figure 7: IEEE® 802.1P Pause*

DiffServ can be used as a QoS mechanism in enterprise networks. It is scalable. Almost all new router products as well as end-products such as VoIP phones support DiffServ and can tag the packets with the appropriate per-hop behavior type. Differentiated Services marking at the edge is read and understood at the core and the packets are forwarded based on the above mentioned priority schemes. Transition Networks' Mac Converter passes these packets transparently.

Such QoS services are not part of any negotiation or signaling between devices themselves. These rules are assigned by local network administrators who understand the above mentioned reasons for congestion and adjust priorities for users, applications or services accordingly. These assigned tags are passed in the packet and are NOT subject to change during the process of auto-negotiation or other forms of signaling. Such approach is called SOFT QoS. 802.1P, IP Precedence and DiffServ are the examples of soft QoS techniques.

### Hard QoS
Hard QoS describes the process during which the devices on the network through signaling can negotiate, request and adjust priority levels for different types of traffic based on the previously agreed values.

Hard QoS includes protocols such as Integrated Services/Resource Reservation Protocol

### Integrated Services/ Resource Reservation Protocol (RSVP)
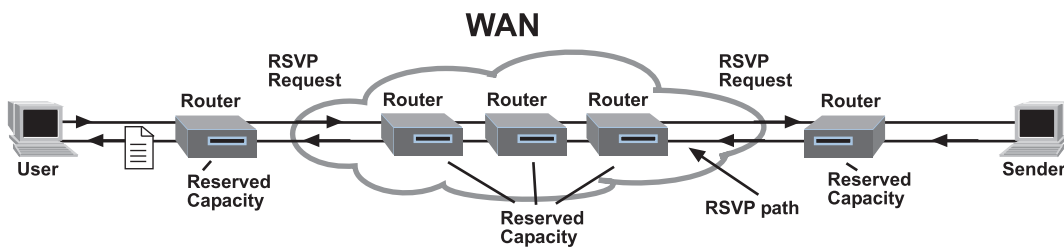


*Figure 8: RSVP*

RSVP enables network devices such as routers or switches to request the necessary /guaranteed bandwidth from other devices on the network for a particular traffic type (e.g. VoIP). Desired delay variances can also be defined in this approach. The RSVP sends a request to reserve specific bandwidth or switching/forwarding capability from other devices on the network. This requirement sent over the network is called flow specification. The requirements can result in three desired transfer types:

1. Traditional Best-effort

2. Rate-sensitive - VoIP requires a guaranteed bit-rate service established bandwidth for video streaming applications.

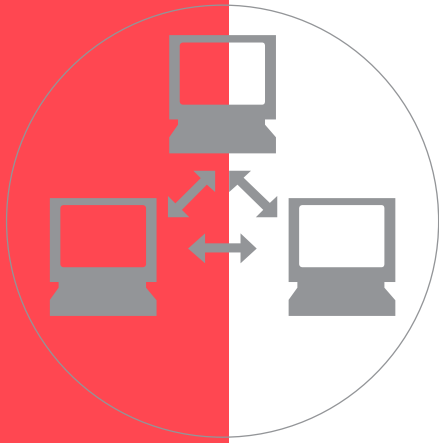3. Delay-sensitive - VoIP requires max delay to be defined and this maximum not allowed to be exceeded.

**Summary**

Voice, audio and video traffic put increasing pressure on both LAN and WAN networks. Users are accustomed to the high reliability and high quality of standard voice and video technologies. Although the transport medium is changing our expectations remain the same. VoIP will not be as fortunate as mobile technology where users sacrifice high quality for the convenience of having a phone on the road. LANs and their, in most cases, opportunistic/best effort protocols face a difficult time handling these high expectations and requirements.

On the other hand, the IP video and phone implementation projects will happen at the increased pace as companies are already starting to fully embrace the cost savings generated by these technologies. This accelerating process will not be matched by an adequate growth in network capacity. QoS is an attractive alternative to aimlessly adding bandwidth to the network.

When voice data becomes part of a network the priority has to be given to the voice packets to "meet" expected high quality of voice calls. ATM, has designed ability for successful QoS, yet since Ethernet runs on 85% of LANs QoS has to efficiently run on this platform as well. 802.1P, IP Precedence, and DiffServ – (soft QoS techniques) help administrators prioritize different types of traffic without any resource reservations. RSVP a Hard QoS technique will help reserve a required level of capacity to support QoS effort. None of these techniques are failure-proof. QoS has to be planned from end-to-end so the bottlenecks are identified and removed.

Finally, QoS will not do magic, and it will not relieve the responsibility of the network managers to plan, and allocate resources accordingly. But the various elements that comprise QoS can offer powerful tools to enable network managers to improve network performance.

**T R A N S I T I O N**
n e t w o r k s

02.03